

DA\_2020 L25

# PCA IN A NUTSHELL

(see Higgs & Attwood  
chap. 2)

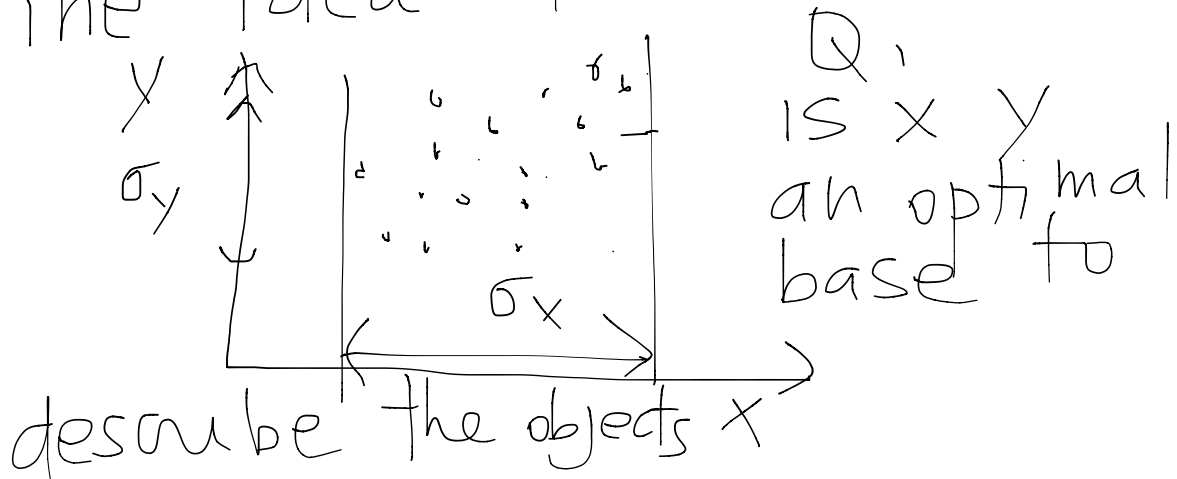
\* structured data ( $N \times P$ )  
matrix  $X$

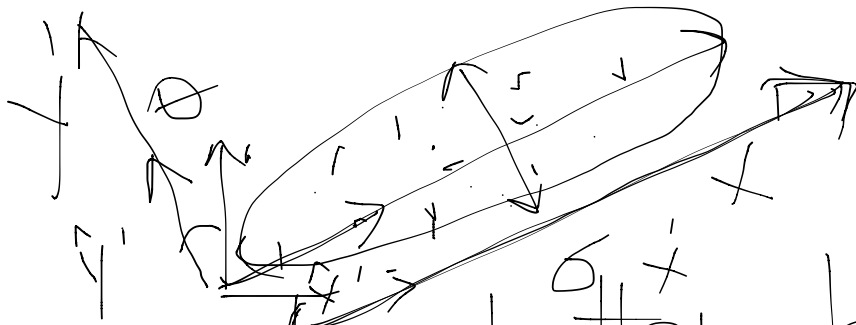
$X_{ij} \leftarrow$

$i = 1, \dots, N$  rows  $\approx$  objects  
 $j = 1, \dots, P$  measured quantities

We are collecting  $P$  dimensional  
data (rows)

The idea of PCA





This is better because  
 we have more variability  
 along  $x'$  than along  $x$

---

Let us find the optimal  
 base through linear  
 algebra (see H&A  
 box 2.2)

1. DATA  $X_{ij}$   $i = 1, \dots, N$   
 $j = 1, \dots, P$

2. AVERAGE  
Parameters over objects

$$\rightarrow \mu_j = \frac{1}{N} \sum_{i=1}^N X_{ij} \quad j = 1, \dots, P$$

3. MEASURE DISPERSIONS

$$\rightarrow \sigma_j = \left( \frac{1}{N} \sum_{i=1}^N (X_{ij} - \mu_j)^2 \right)^{1/2}$$

4. Z-transform  
(standardize data)

$$Z_{ij} \equiv \frac{X_{ij} - \mu_j}{\sigma_j}$$



5 change the base  
 Consider  $P$  vectors

$$\underline{V}_j = (V_{j1}, V_{j2}, \dots, V_{jp}) \leftarrow \text{row}$$

as a new base over which  
 to project the original  $Z_{ij}$   
 (coordinates) of the  $N$  data  
 vectors

\* to be a base a set of  
 vectors should be orthonormal

$$\left. \begin{array}{l} \sum_{k=1}^p V_{jk}^2 = 1 \\ (\forall j=1, \dots, P) \end{array} \right\} \rightarrow \sum_{k=1}^p V_{jk} V_{sk} = \delta_{js} \quad (V_{j,s} \equiv 1, \dots, P)$$

$$\overset{\text{new}}{y}_{ij} = \sum_{k=1}^p \underline{V}_{jk} Z_{ik} \quad \leftarrow \text{o/d}$$

$\textcircled{Z} \rightarrow y$  Transform  $j=1, \dots, P$   
 old new

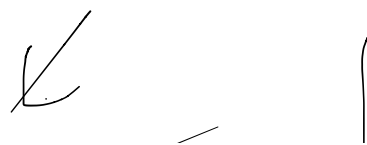
The new base is not arbitrary we want

That the data  $Z_{ij}$  projected over the first vector of the new base

HAVE MAX VARIANCE  
(max sum of residuals)

$$\text{MAX over all bases} \left( \frac{1}{N} \sum_{l=1}^N y_{1l}^2 \right)$$

along direction 1



We start with  $X_{ij} \in \mathbb{R}$   
 \* Define the Covariance matrix of the  $P$  properties

$$C_{jk} = \frac{1}{N} \sum_{i=1}^N (X_{ij} - \mu_j)(X_{ik} - \mu_k)$$

$j, k = 1, \dots, P$

$$= \frac{1}{N} \sum_{i=1}^N Z_{ij} Z_{ik} \leftarrow$$

The sample covariance matrix is the matrix of correlation coefficients

The matrix  $C_{jk}$  is real, symmetric. A theorem of linear algebra tells us

that I can always transform this type of matrices into diagonal form

$P \times P$

$$U \cdot C \cdot U^{-1} = \tilde{C} \quad V \cdot U^{-1} = 1$$

$$\tilde{C} = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_p \end{pmatrix}$$

In general  $\lambda_1, \lambda_2, \dots, \lambda_p$   
 (for generic  $U$  and  $U^{-1}$ ) -  
 ARE NOT IN ORDER

IN PCA we require  
 that  $\lambda_1 > \lambda_2 > \dots > \lambda_p$   
 (through permutation of  
 rows and columns of  
 $U \quad U^{-1}$ )

The  $(\lambda_s)$  are called  
 EIGENVALUES OF  $C_{jk}$  ✓  
 and the new base vectors  
 are called EIGENVECTORS  
 In PCA they are ORDERED  
 From max to min

$$\sum_{j=1}^p V_{nj} C_{jk} = \lambda_n V_{nk}$$

$\uparrow$     $\downarrow$     $\downarrow$     $\downarrow$   
 $n, j, k$     $1, \dots, p$

note if we had  
 defined vectors as  
 (VITTORIO)

$$\begin{pmatrix} V_{n1} \\ V_{n2} \\ \vdots \\ V_{np} \end{pmatrix}$$

we would have here  
 $\sum C_{jk} V_{nj}$  ✓

NOTESSENTIAL



The variance over the first (ordered) direction in the PCA space is

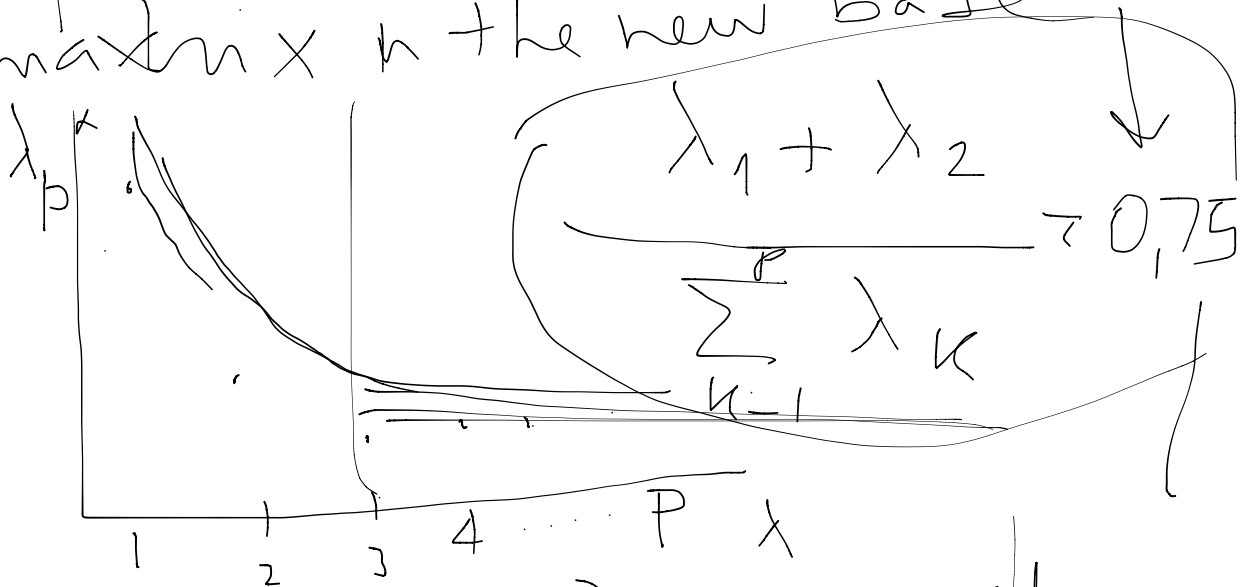
$$\begin{aligned}
 \frac{1}{N} \sum_{i=1}^N y_i^2 &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^P \sum_{k=1}^P V_{1j} z_{ij} V_{1k} z_{ik} \\
 &= \sum_{j=1}^P \sum_{k=1}^P V_{1j} \underbrace{C_{jk}}_{\lambda_1} V_{1k} = \sum_{k=1}^P \lambda_1 V_{1k}^2 \\
 &= \lambda_1 (1)
 \end{aligned}$$

---

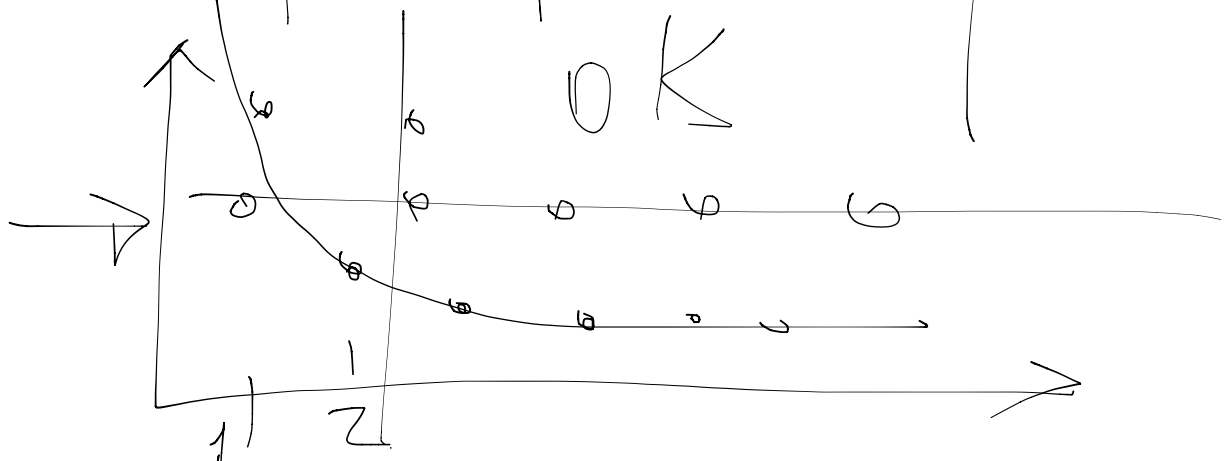
The sum of the eigenvalues is then the total variance seen from the new base

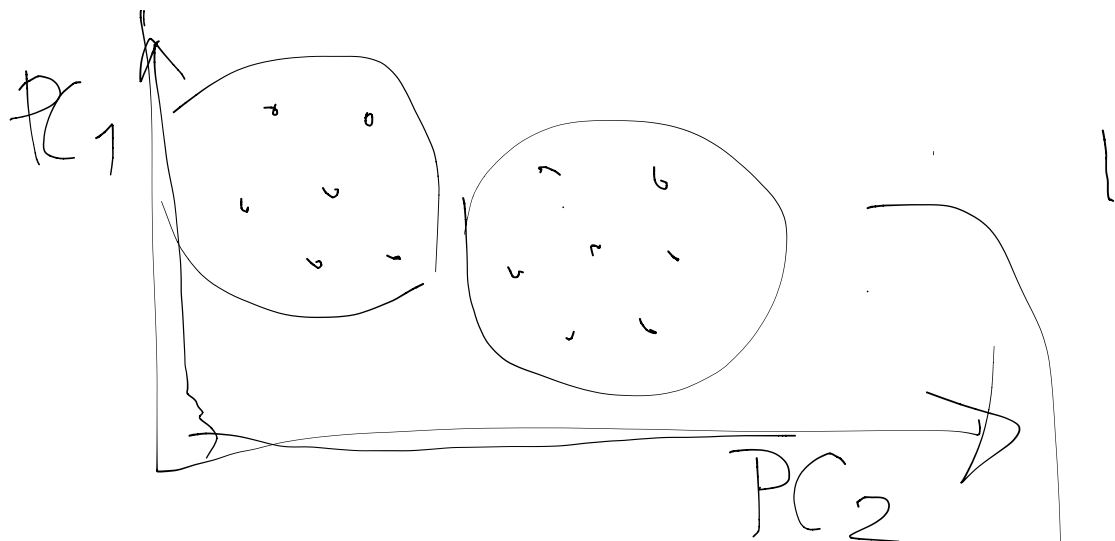
$\lambda_1 + \lambda_2 + \dots + \lambda_p$

Afterwards, if the 'spectrum' of the Covariance matrix  $\Sigma$  in the new base



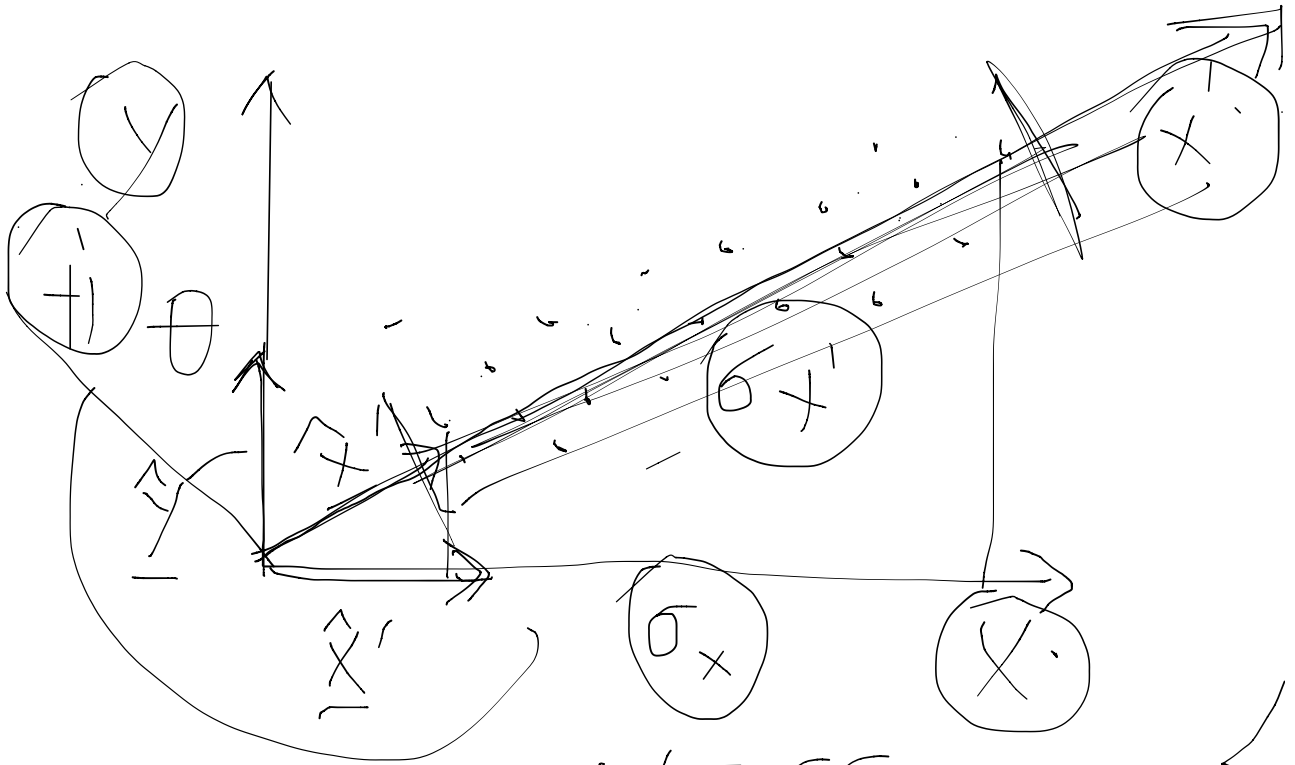
we get a Dimensional reduction we can take only the first 2 Principal directions





Possible natural clustering !

each point in this plane is the scalar product of the original  $\underline{Z}_i$  vectors  $\underline{V}_1$  and  $\underline{V}_2$



PCA IS LESS  
DIMENSIONS  
BUT MORE VARIABIL.

Concentrated varia  
bility