#### CS623: Introduction to Computing with Neural Nets *(lecture-16)*

Pushpak Bhattacharyya Computer Science and Engineering Department IIT Bombay

### **Principal Component Analysis**

# Eaample: IRIS Data (only 3 values out of 150)

ID	Petal Length (a <sub>1</sub> )	Petal Width (a <sub>2</sub> )	Sepal Length (a <sub>3</sub> )	Sepal Width (a <sub>4</sub> )	Classific ation
001	5.1	3.5	1.4	0.2	Iris- setosa
051	7.0	3.2	4.7	1.4,	Iris- versicol or
101	6.3	3.3	6.0	2.5	Iris- virginica

### **Training and Testing Data**

- Training: 80% of the data; 40 from each class: total 120
- Testing: Remaining 30
- Do we have to consider all the 4 attributes for classification?
- Do we have to have 4 neurons in the input layer?
- Less neurons in the input layer may reduce the overall size of the n/w and thereby reduce training time
- It will also likely increase the generalization performance (Occam Razor Hypothesis: A simpler hypothesis (i.e., the neural net) generalizes better

#### The multivariate data



#### Some preliminaries

- Sample mean vector:  $\langle \mu_1, \mu_2, \mu_3, ..., \mu_p \rangle$ For the *i*<sup>th</sup> variable:  $\mu_i = (\Sigma^n_{j=1} x_{ij})/n$
- Variance for the *i*<sup>th</sup> variable:

 $\sigma_i^2 = [\Sigma_{j=1}^n (x_{ij} - \mu_i)^2]/[n-1]$ 

Sample covariance:

 $c_{ab} = [\Sigma^{n}_{j=1} ((x_{aj} - \mu_{a})(x_{bj} - \mu_{b}))]/[n-1]$ This measures the correlation in the data In fact, the correlation coefficient

 $r_{ab} = c_{ab} / \sigma_a \sigma_b$ 

#### Standardize the variables

• For each variable  $x_{ij}$ Replace the values by  $y_{ij} = (x_{ij} - \mu_i)/\sigma_i^2$ 

**Correlation Matrix** 

$$R = \begin{bmatrix} 1 & r_{12} & r_{13} \dots & r_{1p} \\ r_{21} & 1 & r_{23} \dots & r_{2p} \\ & \vdots \\ r_{p1} & r_{p2} & r_{p3} \dots & 1 \end{bmatrix}$$

#### Short digression: Eigenvalues and Eigenvectors

#### $AX = \lambda X$

 $a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots a_{1p}x_p = \lambda x_1$  $a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots a_{2p}x_p = \lambda x_2$ 

 $a_{p1}x_1 + a_{p2}x_2 + a_{p3}x_3 + \dots a_{pp}x_p = \lambda x_p$ Here,  $\lambda$ s are eigenvalues and the solution  $\langle x_1, x_2, x_3, \dots x_p \rangle$ For each  $\lambda$  is the eigenvector

# Short digression: To find the Eigenvalues and Eigenvectors



#### Next step in finding the PCs

$$R = \begin{bmatrix} 1 & r_{12} & r_{13} \dots & r_{1p} \\ r_{21} & 1 & r_{23} \dots & r_{2p} \\ & \vdots \\ r_{p1} & r_{p2} & r_{p3} \dots & 1 \end{bmatrix}$$

Find the eigenvalues and eigenvectors of R

#### Example

#### (from "Multivariate Statistical Methods: A Primer, by Brian

Manly, 3<sup>rd</sup> edition, 1944)

49 birds: 21 survived in a storm and 28 died.
5 body characteristics given
X<sub>1</sub>: body length; X<sub>2</sub>: alar extent; X<sub>3</sub>: beak and head length
X<sub>4</sub>: humerus length; X<sub>5</sub>: keel length *Could we have predicted the fate from the body charateristic*

 $R = \begin{bmatrix} 1.000 \\ 0.735 & 1.000 \\ 0.662 & 0.674 & 1.000 \\ 0.645 & 0.769 & 0.763 & 1.000 \\ 0.605 & 0.529 & 0.526 & 0.607 & 1.000 \end{bmatrix}$ 

### Eigenvalues and Eigenvectors of R

Component	Eigen value	First Eigen- vector: V <sub>1</sub>	V <sub>2</sub>	V <sub>3</sub>	V <sub>4</sub>	V <sub>5</sub>
1	3.612	0.452	0.462	0.451	0.471	0.398
2	0.532	-0.051	0.300	0.325	0.185	-0.877
3	0.386	0.691	0.341	-0.455	-0.411	-0.179
4	0.302	-0.420	0.548	-0.606	0.388	0.069
5	0.165	0.374	-0.530	-0.343	0.652	-0.192

## Which principal components are important?

Total variance in the data=

 $λ_1 + λ_2 + λ_3 + λ_4 + λ_5$ = sum of diagonals of *R*= 5

- First eigenvalue= 3.616 ≈ 72% of total variance 5
- Second ≈ 10.6%, Third ≈ 7.7%, Fourth ≈ 6.0% and Fifth ≈ 3.3%
- First PC is the most important and sufficient for studying the classification

### Forming the PCs

- $Z_1 = 0.451X_1 + 0.462X_2 + 0.451X_3 + 0.471X_4 + 0.398X_5$
- $Z_2 = -0.051X_1 + 0.300X_2 + 0.325X_3 + 0.185X_4 0.877X_5$
- For all the 49 birds find the first two principal components
- This becomes the new data
- Classify using them

#### For the first bird

 $X_1$ =156,  $X_2$ =245,  $X_3$ =31.6,  $X_4$ =18.5,  $X_5$ =20.5 After standardizing  $Y_1$ =(156-157.98)/3.65=-0.54,  $Y_2$ =(245-241.33)/5.1=0.73,  $Y_3$ =(31.6-31.5)/0.8=0.17,  $Y_4$ =(18.5-18.46)/0.56=0.05,  $Y_5$ =(20.5-20.8)/0.99=-0.33

 $PC_1$  for the first bird=  $Z_1 = 0.45X(-0.54) + 0.46X(0.725) + 0.45X(0.17) + 0.47X(0.05) + 0.39X(-0.33)$  = 0.064Similarly,  $Z_2 = 0.602$ 

#### **Reduced Classification Data**

Instead of

Х <sub>1</sub>	Х <sub>2</sub>	Х <sub>3</sub>	X <sub>4</sub>	Х <sub>5</sub>
•		49 rows		

• Use



#### Other Multivariate Data Analysis Procedures

- Factor Analysis
- Discriminant Analysis
- Cluster Analysis

To be done gradually