

# Data Mining CSE 626

## Mid Term

**Dr.S.N.Srihari**  
**{srihari}@cedar.buffalo.edu**

### 1. LINEAR REGRESSION (25)

#### 1.1.

Let  $X$  be an independent variable and  $Y$  be a dependent variable. Given a set of pairs of values  $(X_i, Y_i)$  it is desired to calculate a best fit for the points  $Y_i$  using the method linear regression. Show the mathematical formulation of the problem and derive value for the regression coefficients in terms of the  $(X_i, Y_i)$  values. For sake of simplicity you can assume  $X$  and  $Y$  to be one dimensional variable.

(15)

#### 1.2.

Given the set of values  $X = (3, 9, 11, 5, 2)^T$  and  $Y = (1, 8, 11, 4, 3)^T$ . Evaluate the regression coefficients.

(5)

#### 1.3.

What is the value of variable  $Y$  when  $X = 7$ ?

(5)

## 2. MEASUREMENT AND DATA (35)

For Questions 2.1 2.2 2.3 and 2.4 use the data matrix given in table 2. There are 8 data points each having 2 dimensions X(1) and X(2).

X(1)	X(2)
1	5
3	7
2	6
4	8
1	8
4	5
3	6
2	7

**Table 1.** 8 data points in 2 dimensional space

**2.1.**

Compute the 2\*2 covariance matrix for the data matrix X.

(10)

**2.2.**

Pick all that apply: X(1) and X(2) are:

- (a) Independent                      (b) Dependent  
(c) Correlated                        (d) Uncorrelated

(5)

**2.3.**

The correlation coefficient between X(1) and X(2) is:

- (a) 1.42                      (b) 1.25  
(c) 130                        (d) 0

(5)

**2.4.**

The correlation coefficient between X(1) and X(1) is:

- (a) 1.00                      (b) 0  
(c) 1.42                        (c) Infinity

(5)

**2.5.**

Which of the following is/are metric(s):

- (a) Euclidian distance                      (b) Hamming distance

(c)  $\frac{S_{11}S_{00} - S_{10}S_{01}}{((S_{10} + S_{11}) + (S_{01} + S_{00}) + (S_{11} + S_{01}) + (S_{00} + S_{10}))^{1/2}}$                       (d)  $\frac{S_{11}}{S_{11} + S_{10} + S_{01}}$

(10)

### 3. PRINCIPAL COMPONENT ANALYSIS(40)

#### 3.1.

Given a square matrix  $V$ , write down the expression for the characteristic equation.

(2)

#### 3.2.

Explain briefly what the eigen vector corresponding the largest eigen value signify.

(3)

#### 3.3.

Mention briefly the steps involved in PCA starting with the data matrix  $X$  that is  $n \times p$ , where  $n$  is the number of data points and  $p$  is the number of dimensions. Include steps to show the reduction in dimensions from  $p$  to  $k$  where  $k$  is the reduced dimension for the data points.  $\{k < p\}$ .

(10)

#### 3.4.

Given an  $3 \times 3$  covariance matrix below  $\begin{pmatrix} 1 & 0.5 & 0 \\ 0.5 & 0.333 & 0 \\ 0 & 0 & 0 \end{pmatrix}$  calculate the eigen values and eigen vectors using the chatercteristic equation.

(15)

#### 3.5.

Using the eigen vectors obtained from 3.4, apply the last step of PCA on two data points  $X_1 = (2, 3, 3)^T$  and  $X_2 = (4, 1, 0)^T$  to reduce the dimensions from from 3 to 2. You can use the 2 eigen vectors corresponding to the top 2 eigen values.

(10)