

Data Mining CSE 626

Mid Term Solutions

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)

Solution

Formulation

$Y = \text{Dependent variable}$

$$Y = a + bX$$

Let Y' = data values

$$\begin{aligned} \text{Squared error difference } SSE &= \sum_{i=1}^N (Y' - Y)^2 \\ &= (Y' - (a + bX))^2 \end{aligned}$$

$$a, b = \arg \min_{a, b} (SSE)$$

Derivation

$$\begin{aligned} a &= \arg \min_a (SSE) \\ \frac{\partial SSE}{\partial a} &= \frac{\partial \sum_{i=1}^N (Y_i - a - bX_i)^2}{\partial a} \\ &= -2 \sum_{i=1}^N (Y_i - (a + bX_i)) \\ b &= \arg \min_b (SSE) \\ \frac{\partial SSE}{\partial b} &= \frac{\partial \sum_{i=1}^N (Y_i - a - bX_i)^2}{\partial b} \\ &= -2 \sum_{i=1}^N (Y_i - (a + bX_i))(X_i) \end{aligned}$$

Equate to 0 and rearrange

$$\begin{aligned} Na + b \sum_{i=1}^N X_i &= \sum_{i=1}^N Y_i \\ a \sum_{i=1}^N X_i + b \sum_{i=1}^N X_i^2 &= \sum_{i=1}^N X_i Y_i \end{aligned}$$

Solving for a and b

$$\begin{aligned} \mathbf{a} &= \frac{\sum_{i=1}^N Y_i}{N} - b \frac{\sum_{i=1}^N X_i}{N} \\ &= \mathbf{mean}_Y - \mathbf{b} \cdot \mathbf{mean}_X \\ \mathbf{b} &= \frac{\sum_{i=1}^N (X_i - \frac{\sum_{i=1}^N X_i}{N}) \cdot (Y_i - \frac{\sum_{i=1}^N Y_i}{N})}{\sum_{i=1}^N (X_i - \frac{\sum_{i=1}^N X_i}{N})^2} \\ &= \frac{\sum_{i=1}^N (\mathbf{X}_i - \mathbf{mean}_X) \cdot (\mathbf{Y}_i - \mathbf{mean}_Y)}{\sum_{i=1}^N (\mathbf{X}_i - \mathbf{mean}_X)^2} \end{aligned}$$

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)

Solution

$$\begin{aligned}
b &= \frac{\sum_{i=1}^N (\mathbf{X}_i - \mathbf{mean}_X) \cdot (\mathbf{Y}_i - \mathbf{mean}_Y)}{\sum_{i=1}^N (\mathbf{X}_i - \mathbf{mean}_X)^2} \\
&= \frac{\sum_{i=1}^5 (X_i - 6) \cdot (Y_i - 5.4)}{\sum_{i=1}^5 (X_i - 6)^2} \\
&= \mathbf{1.0} \\
a &= \mathbf{mean}_Y - b \cdot \mathbf{mean}_X \\
&= 5.4 - 1.0 * 6 \\
&= \mathbf{-0.6}
\end{aligned}$$

1.3.

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

(5)

Solution

$$\begin{aligned}
Y &= a + bX \\
&= -0.6 + 1 * 7 \\
&= \mathbf{6.4}
\end{aligned}$$

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)

Solution Mean center the data first.

$$\mu_{X(1)} = 2.5$$

$$\mu_{X(2)} = 6.5$$

Table 2.1 shows the mean centered data.

X(1)	X(2)
-1.5	-1.5
0.5	0.5
-0.5	-0.5
1.5	1.5
-1.5	1.5
1.5	-1.5
0.5	-0.5
-0.5	0.5

Table 2. Mean Centered 8 data points in 2 dimensional space

2 * 2 covariace matrix is now given by $\frac{X'X}{N}$.

$$\begin{aligned} \frac{X'X}{N} &= \frac{\begin{pmatrix} 10 & 0 \\ 0 & 10 \end{pmatrix}}{8} \\ &= \begin{pmatrix} 1.25 & 0 \\ 0 & 1.25 \end{pmatrix} \end{aligned}$$

2.2.

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

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

(5)

Solution Non diagonal entries of covariance matrix = 0

Hence (d) Uncorrelated

NOTE: Covariance matrix does not tell anything about independence. For instance assume a uniform distribution. To test independence $P(X(1),X(2))=P(X(1)) * P(X(2))$

Consider the first value, $X(1)_1 = 1$ and $X(2)_1 = 5$

$$\begin{aligned} P(X(1) = 1, X(2) = 5) &= P(1, 5) \\ \text{From table} &= \frac{1}{8} \\ P(X(1) = 1) &= \frac{2}{8} \\ P(X(2) = 5) &= \frac{2}{8} \\ P(X(1) = 1) * P(X(2) = 5) &= \frac{2}{8} \cdot \frac{2}{8} \\ &= \frac{1}{16} \end{aligned}$$

Hence they are not independent. Since no information is given about whether the data is uniformly distributed or not, **no comment can be made about independence.**

2.3.

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

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

(5)

Solution Since non diagonal entries are 0, correlation coefficient between X(1) and X(2) has to be 0.

2.4.

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

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

(5)

Solution

For any data, this is always equal to 1. The correlation between a variable and itself is always 1.

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)

Solution Take an example of two binary vectors $B_1 = 1, 0$ and $B_2 = 1, 1$.

Check for the three metric properties

1. $d(i,j)=0$ iff $i=j$
2. $d(i,j)=d(j,i)$
3. $d(i,j)+d(j,k) \geq d(i,k)$

Second property does not satisfy for (c), Here $d(B_1, B_2)=0$, But $B_1 \neq B_2$. Hence it is not a metric. Others are metrics.

3. PRINCIPAL COMPONENT ANALYSIS(40)

3.1.

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

(2)

Solution

For any square matrix V , the characteristic equation is given by $|V - \lambda I| = 0$

3.2.

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

(3)

Solution

The eigen vector corresponding to the largest eigen values signifies the direction of **Maximum variance** in the data.

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)

Solution

1. Mean center the data matrix $n \times p$.
2. Compute $p \times p$ covariance matrix.
3. Perform eigen vector decomposition and get eigen vectors $V_{p \times p}$ and eigen values. $\lambda_{p \times 1}$.
4. Choose the top k eigen vectors corresponding to the top k eigen values and form the matrix $V_{p \times k}$.
5. If data point is given by $y_{p \times 1}$, reduced dimension of $y = y'_{k \times 1} = V'_{k \times p} * y_{p \times 1}$.

3.4.

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

(15)

$$\text{Solution } \left| \begin{pmatrix} 1 & 0.5 & 0 \\ 0.5 & 0.333 & 0 \\ 0 & 0 & 0 \end{pmatrix} - \begin{pmatrix} \lambda & 0 & 0 \\ 0 & \lambda & 0 \\ 0 & 0 & \lambda \end{pmatrix} \right| = 0.$$

$$\left| \begin{matrix} 1 - \lambda & 0.5 & 0 \\ 0.5 & 0.333 - \lambda & 0 \\ 0 & 0 & 0 - \lambda \end{matrix} \right| = 0.$$

Solving $\lambda = \{0, 0.065, 1.26\}$

3.5.

The eigen vectors for 3.4 is given in matrix below $\begin{pmatrix} 0 & 0.4729 & -.8817 \\ 0 & -.8817 & -.4719 \\ 1 & 0 & 0 \end{pmatrix}$ where each column is an eigen vector,

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 3 to 2. You can use the 2 eigen vectors corresponding to the top 2 eigen values(The last two columns of the Eigen vector matrix).

(10)

Solution Last two columns of the eigen vector matrix is

$$V'_{3 \times 2} = \begin{pmatrix} 0.4729 & -.8817 \\ -.8817 & -.4719 \\ 0 & 0 \end{pmatrix}$$

For $X_1 = (2, 3, 3)^T$, reduced dimension is given by

$$V'^T \cdot X_1 = \begin{pmatrix} 0.4729 & -.8817 & 0 \\ -.8817 & -.4719 & 0 \end{pmatrix}_{2 \times 3} \cdot \begin{pmatrix} 2 \\ 3 \\ 3 \end{pmatrix}_{3 \times 1} = \begin{pmatrix} -1.7017 \\ -3.1787 \end{pmatrix}_{2 \times 1}$$

For $X_1 = (4, 1, 0)^T$, reduced dimension is given by

$$V'^T \cdot X_1 = \begin{pmatrix} 0.4729 & -.8817 & 0 \\ -.8817 & -.4719 & 0 \end{pmatrix}_{2 \times 3} \cdot \begin{pmatrix} 4 \\ 1 \\ 0 \end{pmatrix}_{3 \times 1} = \begin{pmatrix} -1.005 \\ -3.998 \end{pmatrix}_{2 \times 1}$$