

Q. Compute the linear discriminant projection for the following two dimensional dataset:

Samples for class  $w_1$ ,

$$x_1 = (x_1, x_2) = \{(4, 2), (2, 4), (2, 3), (3, 6), (4, 4)\}$$

Samples for class  $w_2$ ,

$$x_2 = (x_1, x_2) = \{(9, 10), (6, 8), (9, 5), (8, 7), (10, 8)\}$$

Our task is to find a new axis to project the data values.

→ Find class means,

$$\mu_1 = \frac{1}{N_1} \sum_{x \in w_1} x = \frac{1}{5} \begin{bmatrix} 15 \\ 19 \end{bmatrix} = \begin{bmatrix} 3 \\ 3.8 \end{bmatrix}$$

$$\mu_2 = \frac{1}{N_2} \sum_{x \in w_2} x = \frac{1}{5} \begin{bmatrix} 42 \\ 138 \end{bmatrix} = \begin{bmatrix} 8.4 \\ 7.6 \end{bmatrix}$$

Covariance matrix of the first class:

$$\Sigma_1 = \sum_{x \in w_1} (x - \mu_1)(x - \mu_1)^T$$

$$= \frac{1}{4} \left( \begin{bmatrix} 1 \\ -1.8 \end{bmatrix} [1 \ -1.8] + \begin{bmatrix} -1 \\ 0.2 \end{bmatrix} [-1 \ 0.2] + \begin{bmatrix} -1 \\ -0.8 \end{bmatrix} [-1 \ -0.8] + \right. \\ \left. \begin{bmatrix} 0 \\ 2.2 \end{bmatrix} [0 \ 2.2] + \begin{bmatrix} 1 \\ 0.2 \end{bmatrix} [1 \ 0.2] \right)$$

$$= \frac{1}{4} \left( \begin{bmatrix} 1 & -1.8 \\ -1.8 & 3.24 \end{bmatrix} + \begin{bmatrix} 1 & -0.2 \\ -0.2 & 0.04 \end{bmatrix} + \right. \\ \left. \begin{bmatrix} 1 & 0.8 \\ 0.8 & 0.64 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 4.84 \end{bmatrix} + \begin{bmatrix} 1 & 0.2 \\ 0.2 & 0.04 \end{bmatrix} \right)$$

$$\Rightarrow \frac{1}{4} \begin{bmatrix} 4 & -1 \\ -1 & 8.8 \end{bmatrix} = \begin{bmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{bmatrix}$$

~~∴~~ ∴  $S_1 = \begin{bmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{bmatrix}$

Covariance matrix of the second class,

$$S_2 = \frac{1}{N-1} \sum_{x \in U_2} (x_2 - \mu_2)(x_2 - \mu_2)^T$$

$$= \frac{1}{4} \left( \begin{bmatrix} 0.6 \\ 2.4 \end{bmatrix} \begin{bmatrix} 0.6 & 2.4 \end{bmatrix} + \begin{bmatrix} -2.4 \\ 0.4 \end{bmatrix} \begin{bmatrix} -2.4 & 0.4 \end{bmatrix} \right)$$

$$+ \begin{bmatrix} 0.6 \\ -2.6 \end{bmatrix} \begin{bmatrix} 0.6 & -2.6 \end{bmatrix} + \begin{bmatrix} -0.4 \\ -0.6 \end{bmatrix} \begin{bmatrix} -0.4 & -0.6 \end{bmatrix}$$

$$+ \begin{bmatrix} 1.6 \\ 0.4 \end{bmatrix} \begin{bmatrix} 1.6 & 0.4 \end{bmatrix} \right)$$

$$= \frac{1}{4} \left( \begin{bmatrix} 0.36 & 1.44 \\ 1.44 & 5.76 \end{bmatrix} + \begin{bmatrix} 5.76 & -0.96 \\ -0.96 & 0.16 \end{bmatrix} \right)$$

$$+ \begin{bmatrix} 0.36 & -1.56 \\ -1.56 & 6.76 \end{bmatrix} + \begin{bmatrix} 0.16 & 0.24 \\ 0.24 & 0.36 \end{bmatrix}$$

$$+ \begin{bmatrix} 2.56 & 0.64 \\ 0.64 & 0.16 \end{bmatrix} \Big)$$

$$S_2 = \begin{bmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{bmatrix}$$

$$S_W = S_1 + S_2 = \begin{bmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{bmatrix}$$

$$S_B^{-2} (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T = \begin{pmatrix} -5.4 \\ -3.8 \end{pmatrix} (-5.4 \quad -3.8)$$

$$= \begin{bmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{bmatrix}$$

$$S_W^{-1} S_B^{-2} w = \lambda$$

$$\begin{bmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1823 \end{bmatrix} \begin{bmatrix} -5.4 \\ 3.8 \end{bmatrix} = \begin{bmatrix} 1.5813 \\ -0.6 \end{bmatrix}$$

Q. Find LD for given dataset

$$C_1 = \{(4, 2), (2, 4), (2, 3)\}$$

$$C_2 = \{(3, 6), (4, 4), (5, 5)\}$$

$$\mu_1 = \frac{1}{3} \left[ \begin{pmatrix} 4 \\ 2 \end{pmatrix} + \begin{pmatrix} 2 \\ 4 \end{pmatrix} + \begin{pmatrix} 2 \\ 3 \end{pmatrix} \right] = \begin{pmatrix} 2.67 \\ 3 \end{pmatrix}$$

$$\mu_2 = \frac{1}{3} \left[ \begin{pmatrix} 3 \\ 6 \end{pmatrix} + \begin{pmatrix} 4 \\ 4 \end{pmatrix} + \begin{pmatrix} 5 \\ 5 \end{pmatrix} \right] = \begin{pmatrix} 4 \\ 5 \end{pmatrix}$$

$$\tilde{\sigma}_1^2 = \frac{1}{N-1} \sum_{i=1}^N \|x_i - \mu_1\|^2$$

$$\tilde{\sigma}_1^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_1)(x_i - \mu_1)^T$$

$$\Rightarrow \frac{1}{2} \left[ \begin{pmatrix} 1.33 \\ -1 \end{pmatrix} \begin{pmatrix} 1.33 & -1 \end{pmatrix}^T + \begin{pmatrix} -0.67 \\ 1 \end{pmatrix} \begin{pmatrix} -0.67 & 1 \end{pmatrix}^T \right]$$

$$+ \begin{pmatrix} -0.67 \\ 0 \end{pmatrix} \begin{pmatrix} -0.67 & 0 \end{pmatrix}^T \right]$$

$$\Rightarrow \frac{1}{2} \left[ \begin{pmatrix} 1.7689 & -1.33 \\ -1.33 & 1 \end{pmatrix} + \begin{pmatrix} 0.4489 & -0.67 \\ -0.67 & 1 \end{pmatrix} \right]$$

$$+ \begin{pmatrix} 0.4489 & 0 \\ 0 & 0 \end{pmatrix} \right]$$

$$= \frac{1}{2} \begin{bmatrix} 2.66 & -2 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 1.33 & -1 \\ -1 & 1 \end{bmatrix}$$

$$S_2 = \frac{(x_2 - \mu_2)(x_2 - \mu_2)^T}{N-1}$$

$$\Rightarrow \frac{1}{2} \left[ \begin{pmatrix} -1 \\ 1 \end{pmatrix} (-1 \ 1) + \begin{pmatrix} 0 \\ -1 \end{pmatrix} (0 \ -1) + \begin{pmatrix} 1 \\ 0 \end{pmatrix} (1 \ 0) \right]$$

$$\Rightarrow \frac{1}{2} \left[ \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \right]$$

$$\Rightarrow \frac{1}{2} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

$$S_W = S_1 + S_2 = \begin{bmatrix} 2.33 & -1.5 \\ -1.5 & 2 \end{bmatrix}$$

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T = \begin{bmatrix} -2.67 \\ -2 \end{bmatrix} \begin{bmatrix} -2.67 & -2 \end{bmatrix}$$

## Feature Selection

It can be supervised or unsupervised

unsupervised

### 1. Filter methods

- chi square test
- correlation coefficient
- mutual

### 2. Wrapper method

- Forward selection
- Backward elimination
- Recursive method

### 3. Embedded method

- Lasso

Regularization

## Unsupervised

1. Variance threshold
2. Dimensionality Reduction
3. Correlation Analysis

## Chi-Square Test

### Steps

1. Define the Null and Alternate Hypothesis.
2. Calculate the contingency table for our feature

e.g:

	Subscribed (0)	Not Subscribed	Total
Low Income	20	30	50
Medium Income	40	25	65
High Income	10	15	25
Total	70	70	140

3. calculate the expected frequencies.

$$\text{Expected Frequency} = \frac{\text{Row total} \times \text{Col total}}{\text{Total}}$$

$$E_{L,L} = \frac{50 \times 70}{140} = \underline{\underline{25}}$$

$$E_{L,N} = \frac{50 \times 70}{140} = \underline{\underline{25}}$$

$$E_{M,L} = \frac{65 \times 70}{140} = \underline{\underline{32.5}}$$

$$E_{M,N} = \frac{65 \times 70}{140} = \underline{\underline{32.5}}$$

$$E_{H,L} = \frac{25 \times 70}{140} = \underline{\underline{12.5}}$$

$$E_{H_1, NS} = \frac{25 \times 70}{140} = 12.5$$

4. Calculate the chi-square value, with critical value, *and compare*.

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$$\Rightarrow \frac{(20 - 25)^2}{25} + \frac{(30 - 25)^2}{25} + \frac{(40 - 32.5)^2}{32.5} +$$

$$\frac{(25 - 32.5)^2}{32.5} + \frac{(10 - 12.5)^2}{12.5} + \frac{(15 - 12.5)^2}{12.5}$$

$$\Rightarrow 1 + 1 + 1.7307 + 1.7307$$

5. Compare the chi square value or with critical value to accept or reject.

$$\chi^2 = \chi^2_{0.05} =$$

If the calculated value is higher than the one obtained from the table, then the feature is relevant and we reject the null hypothesis.

## Regularization

### Lasso Regularization Regression

- It uses L1 regularization technique called LASSO regression.
- It adds the absolute value of magnitude of the coefficient as a penalty term to the loss function.
- We use Lasso if many features in the dataset are irrelevant or not strongly related to the target.

$$\text{cost} \Rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m |w_i|$$

$\lambda \Rightarrow$  Regularization parameter  
 (how much penalty is added for the weights)

higher  $\lambda \Rightarrow$  it applies stronger regularization and shrinks more coefficients to zero.

### Ridge Regression

- It uses L2 regularization technique.
- It adds the "squared magnitude" of the coefficient as a penalty.
- Used when all features contribute to the

— / —

target but to varying degrees and when we want retain all features with smaller coefficients. i.e. if you have a dataset where all features have some relevance but we want to reduce the impact of noise or multicollinearity.

$$\text{cost} \Rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m w_i^2$$

### Elastic Net Regression

- Combination of L1 and L2 regularization.
- We add

d. Find PCA for the given dataset

$x_1$	2.5	0.5	2.2	1.9	3.1	2.3	2.0	1.0	1.5	1.1
$x_2$	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

$$\bar{x}_1 = \frac{1}{10} (2.5 + 0.5 + 2.2 + 1.9 + 3.1 + 2.3 + 2.0 + 1.0 + 1.5 + 1.1)$$

$$\Rightarrow \underline{\underline{1.81}}$$

$$\bar{x}_2 = \frac{1}{10} (2.4 + 0.7 + 2.9 + 2.2 + 3.0 + 2.7 + 1.6 + 1.1 + 1.6 + 0.9)$$

$$\Rightarrow \underline{\underline{1.91}}$$

Finding covariance matrix,

$$\Sigma = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{bmatrix}$$

$$\text{cov}(x_1, x_1) = \frac{1}{9} \left[ (0.69)^2 + (-1.31)^2 + (0.39)^2 + (0.09)^2 + (1.29)^2 + (0.49)^2 + (0.19)^2 + (-0.81)^2 + (-0.31)^2 + (-0.71)^2 \right]$$

$$\Rightarrow \frac{0}{9} = 0$$

$$\Rightarrow \frac{5.549}{9} = 0.6165$$

1, new	0.69	-1.31	0.39	0.09	1.29	0.49	0.19	-0.81	-0.31	-0.71
2, new	0.49	-1.21	0.99	0.29	1.09	0.79	-0.31	-0.81	-0.31	-1.01

$$\text{cov}(x_1, x_2) = \frac{1}{9} [5.539] = 0.6154$$

$$\text{cov}(x_2, x_1) = 0.6154$$

$$\text{cov}(x_2, x_2) = \frac{1}{9} [0.246.449] = 0.7165$$

$$\Sigma = \begin{bmatrix} 0.6165 & 0.6154 \\ 0.6154 & 0.7165 \end{bmatrix}$$

7 more based outliers on Anomaly Detection

Suppose we have a dataset of daily sales revenue for a retail store over the past 30 days.

[100, 100, 120, 125, 140, 130, 110, 135, 130, 150,  
140, 100, 95, 80, 120, 125, 130, 100, 140, 135,  
130, 145, 110, 120, 130, 135, 140, 125, 130, 120]

We want to identify any days where the sales revenue is significantly different from the other days, which may indicate an anomaly or outlier.

$$\mu = \frac{3740}{30} = \underline{\underline{124.67}}$$

$$z = \frac{x - \mu}{\sigma}$$

$$\sigma = \sqrt{\frac{8346.667}{9}} = \sqrt{922.407} = \underline{\underline{30.453}}$$

$$\underline{\underline{16.68}}$$

LOF scores

Q.

- Consider the 2D dataset representing points in a 2 dimensional space

$$A(1,2) \quad B(2,3) \quad C(3,4) \quad D(10,10)$$

Calculate LOF scores for each point using  $k=2$  nearest neighbours.

$$AB = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2} = 1.414$$

$$AC = \sqrt{(-2)^2 + (-2)^2} = \sqrt{8} = 2.828$$

$$AD = \sqrt{(-9)^2 + (-8)^2} = \sqrt{145} = 12.041$$

$$BC = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2} = 1.414$$

$$BA = 1.414$$

$$BC = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2} = 1.414$$

$$BD = \sqrt{(-8)^2 + (-7)^2} = \sqrt{113} = 10.630$$

$$CA = 2.828$$

$$CB = 1.414$$

$$CD = \sqrt{(-7)^2 + (-6)^2} = \sqrt{85} = 9.2195$$

Calculate reachability distance (before that find  $k=2$  nearest neighbours for each point)

for A,

$$B(1.414), \quad B; \quad AB = 1.414$$

$$C; \quad AC = 2.828$$

for B,

$$A, AB = 1.414$$

$$C, BC = 1.414$$

for C,

$$A, AC = 2.828$$

$$B, BC = 1.414$$

for D,

$$B, BD = 10.630$$

$$C, CD = 9.2195$$

Now calculate the Reachability distances,

for A, ~~BA~~

$$\Rightarrow \max(AB, BA) \Rightarrow \max(1.414, 1.414)$$

$$\text{mean, } \frac{1}{2} \cdot \Rightarrow 1.414$$

(nearest neighbour)  $\equiv$

$$\Rightarrow \max(AC, CB) \Rightarrow \max(2.828, 1.414)$$

$$\Rightarrow 2.828$$

$\equiv$

for B,

$$\Rightarrow \max(BA, AB) = \max(1.414, 1.414)$$

$$\Rightarrow 1.414$$

$$\Rightarrow \max(BC, CB) = \max(1.414, 1.414)$$

$$\Rightarrow 1.414$$

for C,

$$\Rightarrow \max(AC, AB) = \max(2.828, 1.414)$$

$$\Rightarrow 2.828$$

$\equiv$

$$\Rightarrow \max(BC, BA) = \max(1.414, 1.414) \\ \Rightarrow \underline{\underline{1.414}}$$

for D,

$$\begin{aligned} \Rightarrow \max(BD, CB) &= \max(9.2195, 1.414) \\ &\quad \text{by order} \qquad \qquad \qquad \Rightarrow \underline{\underline{9.2195}} \\ \Rightarrow \max(BD, BA) &= \max(10.630, 1.414) \\ &\quad \qquad \qquad \qquad \Rightarrow \underline{\underline{10.630}} \end{aligned}$$

### Modified z-score

usually we take  $\pm 3.5$

$$M_i = \frac{0.6745 \times (x_i - \text{Median})}{\text{MAD}}$$

where  $x_i = \text{datapoint}$

$$\text{MAD} = \text{Median}(|x_i - \text{Median}|)$$

(Median Absolute Deviation)

0.6745 is a scaling factor

- Q. Consider a small dataset with 7 points in 2D plane

$$X = \{A(1,1), B(2,2), C(2,3), D(3,3), E(3,4), F(8,8), G(100, 100)\}$$

11  
19602  
19208

1. Compute distance between the points.

$$AB = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2} = 1.414$$

9604  
9409  
19013

$$AC = \sqrt{(-1)^2 + (-2)^2} = \sqrt{5} = 2.236$$

2464

$$AD = \sqrt{(-2)^2 + (-2)^2} = \sqrt{8} = 2.828$$

$$AE = \sqrt{(-2)^2 + (-3)^2} = \sqrt{13} = 3.605$$

9409  
9216

$$AF = \sqrt{(-7)^2 + (-7)^2} = \sqrt{98} = 9.899$$

$$AG = \sqrt{(-99)^2 + (-99)^2} = \sqrt{19602} = 140.007$$

16928

$$BA = 1.414$$

$$BC = \sqrt{0^2 + (-1)^2} = 1$$

$$BD = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2} = 1.414$$

$$BE = \sqrt{(-1)^2 + (-2)^2} = \sqrt{5} = 2.236$$

$$BF = \sqrt{(-6)^2 + (-6)^2} = \sqrt{72} = 8.485$$

$$BG = \sqrt{(98)^2 + (-98)^2} = \sqrt{19208} = 138.592$$

$$CA = 2.236$$

$$CB = 1$$

$$CD = \sqrt{(-1)^2 + 0^2} = \sqrt{1} = 1$$

$$CE = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2} = 1.414$$

$$CF = \sqrt{(-6)^2 + (-5)^2} = \sqrt{61} = 7.810$$

$$CG = \sqrt{(98)^2 + (97)^2} = \sqrt{19013} = 137.88$$

$$DA = 2.828$$

$$DB = 1.414$$

$$DC = 1$$

$$DE = \sqrt{0^2 + (-1)^2} = \sqrt{1} = 1$$

$$DF = \sqrt{(-5)^2 + (-5)^2} = \sqrt{50} = 7.071$$

$$DG = \sqrt{(-97)^2 + (-97)^2} = \sqrt{18818} = 137.178$$

$$EA = 3.605$$

$$EB = 2.236$$

$$EC = 1.414$$

$$ED = 1$$

$$EF = \sqrt{(-5)^2 + (-4)^2} = 6.403$$

$$EG = \sqrt{(92)^2 + (-96)^2} = \sqrt{186125} = 136$$

$$GA = FA = 9.899$$

$$GB = FB = 8.485$$

$$GC = FC = 7.810$$

$$GD = FD = 7.071$$

$$GE = FE = 6.403$$

$$GF = FA = \sqrt{(-92)^2 + (-92)^2} = \sqrt{16928} = 130.107$$

$$GA = 140.007$$

$$GB = 138.592$$

$$GC = 137.88$$

$$GD = 137.178$$

$$GE = 136$$

$$GF = 130.107$$

Finding  $k=2$  nearest neighbours of ~~reachability~~ values  
for  $A \Rightarrow B, C$ .

$$AB, AC (1.414, 2.236)$$

$$R_N = \max(AB, BC)$$

$$\text{for } B \Rightarrow C, A (1, 1.414)$$

for  $C \Rightarrow B, D [1, 1]$

for  $D \Rightarrow C, E [1, 1]$

for  $E \Rightarrow C, D [1.414, 1]$

for  $F \Rightarrow E, D [6.403, 7.071]$

for  $G \Rightarrow E, F [136, 130.107]$

Finding Reachability values,

for  $A$

$$\Rightarrow \max(AB, BC) = \max(1.414, 1) \\ = \underline{\underline{1.414}}$$

$$\Rightarrow \max(AC, CB) = \max(2.236, 1) \\ = \underline{\underline{2.236}}$$

for  $B$

$$\Rightarrow \max(BE, CB) = \max(1, 1) \\ = \underline{\underline{1}}$$

$$\Rightarrow \max(BA, AB) = \max(1.414, 1.414) \\ = \underline{\underline{1.414}}$$

for C,

$$\Rightarrow \max(CB, BC) = \max(1, 1) = 1$$

$$\Rightarrow \max(CD, DC) = \max(1, 1) = 1$$

for D,

$$\Rightarrow \max(DC, CB) = \max(1, 1) = 1$$

$$\Rightarrow \max(DE, ED) = \max(1, 1) = 1$$

for E,

$$\Rightarrow \max(ED, DC) = \max(1, 1) = 1$$

$$\Rightarrow \max(EC, CB) = \max(1.414, 1) = 1.414$$

for F,

$$\Rightarrow \max(FE, ED) = \max(6.403, 1) = 6.403$$

$$\Rightarrow \max(FD, DC) = \max(7.071, 1) = 7.071$$

for G,

$$\Rightarrow \max(GE, ED) = \max(136, 1) = 136$$

$$\Rightarrow \max(GF, FE) = \max(130.107, 6.403) \\ = 130.107$$

Next we calculate LRD values.

$$\text{LRD} = \frac{1}{k} \sum_{q \in N_k(p)} \text{reach dist}(p, q)$$

$$\text{LRD of } A = \frac{1}{\frac{1}{2}(1.414 + 2.23)} = \frac{1}{1.822} = 0.54$$

$$\text{LRD of } B = \frac{1}{\frac{1}{2}(1 + 1.414)} = \frac{1}{1.207} = 0.8285 \text{ (-maybe)}$$

$$\text{LRD of } C = \frac{1}{\frac{1}{2}(1+1)} = \frac{1}{1}$$

$$\text{LRD of } D = \frac{1}{\frac{1}{2}(1+1)} = \frac{1}{1}$$

~~$$\text{LRD of } E = \frac{1}{\frac{1}{2}(1+1)} = \frac{1}{1}$$~~

$$\text{LRD of } E = \frac{1}{\frac{1}{2}(1 + 1.414)} = \frac{1}{1.207} = 0.8285$$

$$\text{LRD of } F = \frac{1}{\frac{1}{2}(6.403 + 7.071)} = \frac{1}{6.737} = 0.1484$$

$$\text{LRD of } G = \frac{1}{\frac{1}{2}(136 + 130.107)} = \frac{1}{133.0535} = 0.0075$$

(mostly outlier)

$$LOF = \frac{1}{k} \sum_{q \in N_k(p)} LRD(q)$$

$$\begin{aligned} LOF \text{ of } A &= \frac{\sum \text{LRRD of } A}{k \times \text{LRRD of } A} \\ &= \frac{\sum \text{LRRD's of } A's \text{ neighbours}}{k \times \text{LRRD of } A} \end{aligned}$$

point with highest LOF value is the outlier.

### Kernel induced Feature Expansion

Different types of kernel:

1) Linear Kernel

$$k(x, y) = x^T y$$

$x$  and  $y$  are vectors,

$$\text{Therefore, } k(x, y) = \phi(x) \cdot \phi(y) = x^T y$$

2) Polynomial Kernel

$$k(x, y) = (x^T y)^q$$

$q$  is the degree of the polynomial.

3) Gaussian Kernel (Radial Basis Function)

$$k(x, y) = e^{-\frac{(x-y)^2}{2\sigma^2}}$$

$$k(x, y) = e^{-\frac{(x-y)^2}{2\sigma^2}}$$

$\sigma$  is the constant

#### 4) Sigmoid Kernel

Consider 2 data points

$$x = (1, 2) \text{ and } y = (2, 3) \text{ with } \sigma = 1.$$

Apply RBF kernel and find the value of RBF kernel for these points.

$$k(x, y) = e^{-\frac{(x-y)^2}{2\sigma^2}}$$

$$(1-2)^2 + (2-3)^2 = (-1)^2 + (-1)^2 = 2,$$

$$\text{If } \sigma = 1,$$

$$k((1, 2), (2, 3)) = e^{-\frac{(1-2)^2}{2}} = e^{-\frac{(2)^2}{2(1)^2}}$$

Q Consider 2 data points  $x = (1, 2)$  and  $y = (2, 3)$  with  $c = 1$ . Apply linear, homogeneous and inhomogeneous kernels:-

#### 4) Inhomogeneous Kernel

$$k(x, y) = (c + x^T y)^q$$

$c$  is a constant and  $q$  is the degree of the polynomial.

linear kernel ( $q=1$ )

$$k(x, y) = \left( \begin{pmatrix} 1 \\ 2 \end{pmatrix}^T \begin{pmatrix} 2 & 3 \end{pmatrix} \right)^1 = 2+6 = 8,$$

when  $q=2$ ,

— / —

$$k(x, y) = \left( \begin{pmatrix} 1 \\ 2 \end{pmatrix}^T \begin{pmatrix} 2 & 3 \end{pmatrix} \right)^2 = 8^2 = \underline{\underline{64}}$$

If  $c=1$

$$k(x, y) = (c + x^T y)^2$$

$$= (1 + \begin{pmatrix} 1 \\ 2 \end{pmatrix}^T \begin{pmatrix} 2 & 3 \end{pmatrix})^2 = (1+8)^2 = \underline{\underline{81}}$$

Q Find PCA for the given dataset

	1	2	3	4	5	6	7	8	9	10
$x_1$	2.5	0.5	2.2	1.9	3.1	2.3	2.0	1.0	1.5	1.1
$x_2$	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

$$\text{No. of features } (n) = 2$$

$$\text{No. of samples } (N) = 10$$

$$\mu_{x_1} = 1.81$$

$$\mu_{x_2} = 1.91$$

Covariance matrix :

Ordered pair:  $(x, x)$   $(x, y)$   $(y, x)$   $(y, y)$

$$\text{cov}(x, x) = \frac{1}{N-1} \sum (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$$

$$= \frac{1}{9} [(2.5 - 1.81)^2 + (0.5 - 1.81)^2 + (2.28 - 1.81)^2 + (1.9 - 1.81)^2 + \\ (3.1 - 1.81)^2 + (2.3 - 1.81)^2 + (2 - 1.81)^2 + (1 - 1.81)^2 + \\ (1.5 - 1.81)^2 + (1.1 - 1.81)^2]$$

$$= \frac{1}{9} [0.4761 + 1.7161 + 0.1521 + 0.0081 + 1.6641 + \\ + 0.28014 + 0.0361 + 0.6561 + 0.0961 + 0.5041]$$

$$= \frac{1}{9} (5.549) = \underline{\underline{0.6165}} = \underline{\underline{0.62}}$$

$$\text{cov}(x, y) = \frac{1}{N-1} \sum (x_{ik} - \bar{x}_i)(y_{jk} - \bar{y}_j)$$

$$= \frac{1}{9} [0.34 + 1.6 + 0.386 + 0.03 + 1.41 + 0.387 \\ - 0.059 + 0.656 + 0.096 + 0.714]$$

$$\Rightarrow \frac{1}{9} (5.563) = 0.618 = \underline{\underline{0.62}}$$

$$\text{cov}(y, x) = 0.618 = \underline{\underline{0.62}}$$

$$\begin{aligned}\text{cov}(y, y) &= \frac{1}{9} [0.24 + 1.46 + 0.98 + 0.084 + 1.188 \\ &\quad + 0.624 + 0.096 + 0.66 + 0.016 + 1.02] \\ &= \underline{\underline{0.7164}}\end{aligned}$$

$$\text{cov} = \begin{bmatrix} 0.6185 & 0.62 \\ 0.62 & 0.7164 \end{bmatrix}$$

$$\text{Eigen values : } \det \begin{bmatrix} 0.62 - \lambda & 0.62 \\ 0.62 & 0.72 - \lambda \end{bmatrix}$$

$$\Rightarrow (0.62 - \lambda)(0.72 - \lambda) - 0.3844$$

$$\Rightarrow 0.4464 - 0.62\lambda - 0.72\lambda + \lambda^2 - 0.3844$$

$$\Rightarrow \lambda^2 - 1.34\lambda + 0.062$$

$$D = \lambda^2 - 1.34\lambda + 0.062$$

$$\sqrt{b^2 - 4ac} = 1.8 - 4(1)(0.062)$$

$$\Rightarrow \underline{\underline{1.552}}$$

$$\lambda = \frac{1.34 \pm 1.246}{2} = 0.67 \pm 1.246$$

$$\lambda_1 = \underline{\underline{1.9}}$$

$$\lambda_2 = \underline{\underline{-0.576}}$$

$$(\text{cov} - \lambda \mathbf{I}) \mathbf{U}_1 = 0 \Rightarrow \begin{bmatrix} -1.28 & 0.62 \\ 0.62 & -1.18 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = 0$$

$$1.28x - 0.62y = 0$$

$$-1.28x + 0.62y = 0$$

$$0.62x - 1.18y = 0$$

$$x = \frac{0.62y}{1.28}$$

$$x = 0.48y$$

$$1.28 - 0.62y = 0.62x - 1.18y \Rightarrow x(1.28 + 1.18)$$

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0.48 \\ 1 \end{bmatrix}$$

Normalize the eigen vector,

$$e_1 = \begin{bmatrix} 0.48 & \sqrt{0.48^2 + 1^2} \\ 1 & \sqrt{0.48^2 + 1^2} \end{bmatrix} = \begin{bmatrix} 0.43 \\ 0.9 \end{bmatrix}$$

Derive new data

$$PC_1 = e_1^T \begin{bmatrix} 2.5 - 1.81 \\ 2.4 - 1.91 \end{bmatrix} = [0.43 \ 0.9] \begin{bmatrix} 0.69 \\ 0.49 \end{bmatrix} = 0.7377$$

$$PC_2 = [0.43 \ 0.9] \begin{bmatrix} -1.31 \\ -1.21 \end{bmatrix} = -1.65$$

$$PCl_5 = [0.43 \quad 0.9] \begin{bmatrix} 0.369 \\ 0.99 \end{bmatrix} = 1.06$$

$$PCl_9 = [0.43 \quad 0.9] \begin{bmatrix} 0.09 \\ 0.29 \end{bmatrix} = 0.3$$

$$PCl_5 = [0.43 \quad 0.9] \begin{bmatrix} 1.29 \\ 1.09 \end{bmatrix} = 1.54$$

$$PCl_6 = [0.43 \quad 0.9] \begin{bmatrix} 0.49 \\ 0.79 \end{bmatrix} = 0.92$$

$$PCl_7 = [0.43 \quad 0.9] \begin{bmatrix} 0.19 \\ -0.31 \end{bmatrix} = -0.2$$

$$PCl_8 = [0.43 \quad 0.9] \begin{bmatrix} -0.81 \\ -0.81 \end{bmatrix} = -1.1$$

$$PCl_9 = [0.43 \quad 0.9] \begin{bmatrix} -0.31 \\ -0.31 \end{bmatrix} = -0.41$$

$$PCl_{10} = [0.43 \quad 0.9] \begin{bmatrix} -0.91 \\ -1.01 \end{bmatrix} = -1.2$$

$$\therefore PC = \begin{bmatrix} 0.7 & -1.65 & 1.06 & 0.3 & 1.54 & 0.92 & -0.2 & -1.1 \\ -0.41 & -1.2 \end{bmatrix}$$