

# DEE DEE

on Given data, apply different normalization / /

## \* Independent Component Analysis: 102.

Cocktail party problem.

Changing signals: Signals that change over time.

- Dimensionality reduction by identifying independent components.  
→ mixing matrix

$\mathbf{z} = \mathbf{A}\mathbf{x}$

observed signal → source → denoising matrix.

Recovering original:  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{z}$

$$(\mathbf{w} = \mathbf{A}^{-1})$$

## \* Data Normalization:

### \* Min-max normalization [0,1].

$$y = \frac{x - \min}{\max - \min}$$

Data:  $x = \frac{200}{\text{min}}; \frac{300}{\text{min}}; \frac{100}{\text{min}}; \frac{600}{\text{min}}; \frac{1000}{\text{min}}$ .

$$x' = \frac{200-200}{1000-200}; \frac{300-200}{1000-200}; \frac{100-200}{1000-200}; \frac{600-200}{1000-200}; \frac{1000-200}{1000-200}.$$

$$x' = 0; \frac{100}{800}; \frac{100}{800}; \frac{600}{800}; 1$$

$$x' = 0; 0.125; 0.25; 0.5; 1.$$

2. Z-Score normalization: Mean - standard deviation:

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

$$x = 200; 300; 400; 600; 1000 \quad \mu = 500$$

$$\sigma = \sqrt{\frac{(x - \mu)^2}{n}} = \sqrt{\frac{90000 + 10000 + 10000 + 10000 + 250000}{5}} = \sqrt{80000}.$$

$$\sigma = \underline{283}.$$

$$Z = \frac{200 - 500}{283}; \frac{300 - 500}{283}; \frac{400 - 500}{283}; \frac{600 - 500}{283}; \frac{1000 - 500}{283}.$$

$$Z = -1.06; -0.71; -0.35; 0.35; \underline{2.11}.$$

3. Z-score normalization: Mean - Absolute deviation:

$$Z = \frac{x - \mu}{A}.$$

$$A = |200 - 500| + |300 - 500| + |400 - 500| + |600 - 500| + |1000 - 500| / 5$$

$$A = 300 + 200 + 100 + 100 + 500 / 5 = 1200 / 5 = \underline{240}.$$

$$Z = \frac{200 - 500}{240}; \frac{300 - 500}{240}; \frac{400 - 500}{240}; \frac{600 - 500}{240}; \frac{1000 - 500}{240}.$$

$$Z = -1.25; -0.83; -0.42; 0.42; \underline{2.1}.$$

## 9. Normalization by Decimal Scaling:

- find value of  $g_j$ .
- smallest integer  $j$  such that  $\max\left(\frac{v_i}{10^j}\right) \leq 1$

$$\Rightarrow v = .200 ; .300 ; .100 ; .600 ; .1000$$

$$\frac{200}{10^3} \leq 1 = 0.2 \quad \frac{100}{10^3} \leq 1 = 0.1 \quad \frac{1000}{10^3} = 1$$

$$\frac{300}{10^3} \leq 1 = 0.3 \quad \frac{600}{10^3} \leq 1 = 0.6$$

$$\therefore x' = 0.2 ; 0.3 ; 0.1 ; 0.6 ; 1.$$

## \* Binning:

- A form of quantization.
- putting into buckets.
- reduces overfitting.
- smoothing data.

### 1- Equal frequency binning:

Input: [5, 10, 11, 13, 15, 35, 50, 55].

↪ Bin 1: [5, 10, 11, 13]

Bin 2: [15, 35, 50, 55].

### 2- Equal width binning: $w = \frac{\max - \min}{\text{no. of bins}}$

Bin 1:  $[\min + w]$ .

Bin 2:  $[\min + 2w]$ .

Bin n:  $[\min + nw]$ .

3 Customized binning:  
Outlier

\* Outliers:  
Deviated datapoints.

Graphically: Scatter plot; Boxplot.

\* IQR: Inter Quartile Range

1. sort the data from low to high.
2. Identify the first Quartile ( $Q_1$ ), the median, third Quartile ( $Q_3$ ).
3. Calculate  $IQR = Q_3 - Q_1$ .
4. Calculate upper fence =  $Q_3 + (1.5 \times IQR)$ .
5. Calculate lower fence =  $Q_1 - (1.5 \times IQR)$ .
6. Use fences to highlight any outliers, all values that fall outside the fences.

The outliers are those values that are greater than the upper fence and lesser than the lower fence.

on sorted data: 22 24 26

26 37 21 28 35 22 31 53 41 4 29

1. sorted: 22 24 26 28 29 31 35 37 41 53 4

2.  $Q_1 = 26$ ,  $Q_3 = \frac{41+37}{2} = 39$ ,

3.  $IQR = 39 - 26 = 13$ .

4.  $\text{Upper} = 39 + (1.5 \times 13) = 58.5$ ;  $\text{lower} = 26 - (1.5 \times 13) = 6.5$ .

5. Outliers = ~~64~~



29 Jan

11

## \* Z-score based Outliers or Anomaly Detection.

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

$Z_{\text{threshold}} = \pm 3, \pm 2, \pm 3.5, \pm 4.$

On:	100	150	120	125	140	130	140	135	130	150
	140	100	95	80	120	125	130	100	140	135
	130	145	110	120	130	135	140	125	130	120.

$$\mu = 124.64.$$

$$\sigma = 16.96. = \sqrt{\frac{\sum (x_i - \mu)}{N}}$$

$$z_{100} = \frac{100 - 124.64}{16.96} = -1.45. \quad z_{125} = 0.02.$$

$$z_{150} = 1.49. \quad z_{170} = -0.275. \quad z_{140} = 0.904.$$

$$z_{130} = 0.319. \quad z_{110} = -0.865. \quad z_{135} = 0.61.$$

$$z_{80} = -1.45. \quad z_{80} = -2.634. \quad z_{145} = 1.198.$$

Threshold  $\pm 3$ : No outliers.

$\pm 2$ : 80 is the outlier.

## \* LOF:

On:  $K = 2$  nearest neighbours.

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

① Euclidean distance formula:  $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

A(1,2)      B(2,3)      C(3,1)      D(10,10)

$$AB = \sqrt{(2-1)^2 + (3-2)^2} = \sqrt{1+1} = \sqrt{2} \approx 1.4.$$

$$BC = \sqrt{(3-2)^2 + (4-3)^2} = \sqrt{1+1} = \sqrt{2} \approx 1.4.$$

$$CD = \sqrt{(10-3)^2 + (10-4)^2} = \sqrt{49+36} = 9.22.$$

$$AC = \sqrt{(3-1)^2 + (4-2)^2} = \sqrt{1+1} = 2.83.$$

$$AD = \sqrt{(10-1)^2 + (10-2)^2} = \sqrt{81+64} = 12.04.$$

$$BD = \sqrt{(10-2)^2 + (10-3)^2} = \sqrt{64+49} = 10.63.$$

(a) finding nearest neighbors:

A: B & C.

C: B & A.

B: C & A.

D: C & B.

~~C~~

(b) Reachability:  $\text{reach-dist}(p,q) = \max(\text{dist}(p,q), k \text{ distance})$ .

$$\text{reach-dist}(A,B) = \max(1.4, 2.83) = 2.83.$$

$$\text{reach-dist}(A,C) = \max(2.83, 1.4) = 2.83.$$

$$\text{reach-dist}(A,D) = \max(12.04,$$

$$\text{reach-dist}(B,A) = \max(1.4, 2.83) = 2.83.$$

$$\text{reach-dist}(B,C) = \max(1.4, 2.83) = 2.83.$$

$$\text{reach-dist}(C,B) = \frac{1.4}{2.83}; \quad \text{reach-dist}(C,A) = 2.83.$$

$$\text{reach-dist}(D,C) = \max(9.22, 10.63) = 10.63;$$

$$\text{reach-dist}(D,B) = \max(10.63, 1.4) = 10.63;$$

?  
AD or CD?

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

$$LRD(A) = \frac{1}{\frac{1}{2}[1.4+2.83]} = \underline{\underline{0.44}}$$

$$LRD(B) = \frac{1}{\frac{1}{2}[1.4+2.83]} = \underline{\underline{0.44}}.$$

$$LRD(C) = \frac{1}{\frac{1}{2}[2.83+2.83]} = \underline{\underline{0.35}}.$$

$$LRD(D) = \frac{1}{\frac{1}{2}[9.22+10.63]} = \underline{\underline{0.1}}.$$

$$\textcircled{2} \quad LOF(p) = \frac{\frac{1}{k} \sum_{q \in N_k(p)} LRD(q)}{LRD(p)}$$

\* Modified Z-score:

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

where  $X_i$  - data point.

$MAD$  - Median ( $|X_i - \text{Median}|$ ) .

- Median Absolute Deviation.

0.6745 is a scaling factor.

Qn. Data points = 10 12 14 15 100.

Median = 14.

$$MHD = \text{Median} (|10-14|, |12-14|, |14-14|, |15-14|, |100-14|)$$

$$MAD = \text{Median} (1, 2, 0, 1, 86) \Rightarrow \text{Median} (0, 1, 2, 4, 86).$$

$$MAD = \underline{\underline{2}}.$$

$$\text{Misclassification score} = \frac{0.6745}{2} \times (10-14) = \underline{\underline{-1.349}}.$$

$$M_2 = 0.6745 \times (12-14) = \underline{\underline{-0.6745}}.$$

$$M_3 = 0. \quad M_4 = \frac{0.6745}{2} \times (15-14) = \underline{\underline{0.33725}}.$$

$$M_5 = 0.6745 \times (100-14) = \underline{\underline{29.0035}}.$$

Qn.  $n=7$ ; 2D plane.  $k=2$ .

$$X = A(1,1) \quad B(2,2) \quad C(2,3) \quad D(3,3)$$

$$E(3,4) \quad F(8,8) \quad G(100,100).$$

Q. Euclidean distance:

$$AB = \sqrt{(2-1)^2 + (2-1)^2} = 1.414. \quad AE = \sqrt{(8-1)^2 + (8-1)^2} = 3.6.$$

$$AC = \sqrt{(2-1)^2 + (3-1)^2} = 2.24. \quad AF = \sqrt{(8-1)^2 + (8-1)^2} = 9.9.$$

$$AD = \sqrt{(3-1)^2 + (3-1)^2} = 2.83. \quad AG = \sqrt{(100-1)^2 + (100-1)^2} = 140.007.$$

$$BC = \sqrt{0+1^2} = \sqrt{1} = 1. \quad BD = \sqrt{1^2 + 1^2} = 1.414. \quad BE = \sqrt{1^2 + 1^2} = 2.24.$$

$$BF = \sqrt{6^2 + 6^2} = 8.48. \quad BG = \cancel{138.6}$$

$$CD = \sqrt{1^2 + 0} = 1. \quad \cancel{1}$$

$$CE = \sqrt{1^2 + 1^2} = 1.414.$$

$$CF = \sqrt{6^2 + 5^2} = 7.81.$$

$$CG = \cancel{137.88}.$$

$$DE = \sqrt{0+12} = 1. \quad DF = \sqrt{5^2+5^2} = 7.07. \quad DB = 137.02.$$

$$EF = \sqrt{5^2+4^2} = 6.4. \quad PG = 136.47. \quad FG = 130.1.$$

② Nearest neighbours:

$$A = B, C = [1.414, 2.24].$$

$$B = A, C = [1.414, 1].$$

$$C = B, D = [1, 1].$$

$$D = C, E = [1, 1].$$

$$E = C, D = [1.414, 1].$$

$$F = D, E = [7.1, 6.4]$$

$$G = F, \cancel{P} = [136.47, 136.47]. [136.47, 130.1].$$

③ Reach distance:  $(p, q) = \max(\text{dist}(pq), q \text{'s neighbour})$

$$A \rightarrow B = \max(1.414, 1) = \underline{1.414}. \quad \text{[first occurrence]}$$

$$A \rightarrow C = \max(2.24, 2.24) = \underline{2.24}.$$

$$B \rightarrow A = 1.414; \quad B \rightarrow C = \max(1, \cancel{1.414}, 2.24) = \underline{2.24}.$$

$$C \rightarrow B = 1.414; \quad C \rightarrow D = \max(1, 1) = \underline{1}.$$

$$D \rightarrow C = \cancel{1.414}; \quad D \rightarrow E = \max(1, 1) = \underline{1}.$$

$$E \rightarrow C = \max(1.414, 2.24) = \cancel{2.24}; \quad E \rightarrow D = \underline{1}.$$

$$F \rightarrow D = \max(7.1, 1) = \underline{7.1}; \quad F \rightarrow E = \max(6.4, 1) = \underline{6.4}.$$

$$G \rightarrow E = \max(136.47, 1) = 136.47; \quad G \rightarrow F = \max(130.1, 130.1) = \underline{130.1}.$$

1 Feb

(4)  $LRD = \frac{1}{k} \sum_{i=1}^n \text{reach dist (p,q)}$

A:  $\frac{1}{2} / (1.414 + 2.29) = \frac{1}{1.821} = 0.544$

B:  $\frac{1}{2} / (1.414 + 2.29) = \frac{1}{1.821} = 0.544$

C:  $\frac{1}{2} / (1.414 + 1) = \frac{1}{1.204} = 0.8285$

D:  $\frac{1}{2} / (2.29 + 1) = \frac{1}{1.62} = 0.6143$

E:  $\frac{1}{2} / (2.29 + 1) = \frac{1}{1.62} = 0.6143$

F:  $\frac{1}{2} / (1.414 + 0.9) = \frac{1}{2.314} = 0.438$

G:  $\frac{1}{2} / (136.44 + 130.0) = \frac{1}{133.245} = 0.0075$

(5)  $LOF = \frac{1}{k} \sum_{i=1}^n \frac{LRD(p)}{LRD(q)}$

A:  $\frac{1}{2} \left[ \frac{LRD(B) + LRD(C)}{LRD(A)} \right] = \frac{1}{2} \left( \frac{0.544 + 0.83}{0.544} \right) = \underline{1.26}$

B:  $\frac{1}{2} \left( \frac{A+C}{B} \right) = \frac{1}{2} \left( \frac{0.544 + 0.83}{0.544} \right) = \underline{1.26}$

C:  $\frac{1}{2} \left( \frac{B+D}{C} \right) = \frac{1}{2} \left( \frac{0.544 + 0.6143}{0.8285} \right) = \underline{0.9}$

D:  $\frac{1}{2} \left( \frac{C+E}{D} \right) = \frac{1}{2} \left( \frac{0.8285 + 0.6143}{0.6143} \right) = \underline{1.171}$

E:  $\frac{1}{2} \left( \frac{D+F}{E} \right) = \frac{1}{2} \left( \frac{0.6143 + 0.438}{0.6143} \right) = \underline{1.171}$

F:  $\frac{1}{2} (D, E) / f = \frac{1}{2} (0.6143 + 0.438) / 0.198 = \underline{1.171}$

G:  $\frac{1}{2} (E, F) / G = \frac{1}{2} (0.6143 + 0.438) / 0.0075 = \underline{51.05}$

on LDA, PCA, performance evaluations, overfitting/underfitting-efficiency.

-/-

∴ The outliers are F & G. (mostly G).

\* ~~Type of Kernels:~~

Lower dimension  $\rightarrow$  higher dimension, and  
the linearly separate it. Used for classification.

- Types of Kernels:

① Linear Kernel:  $K(x, y) = x^T y$ .

$$K(x, y) = \phi(x) \cdot \phi(y) = x^T y$$

② Polynomial kernel:  $K(x, y) = (x^T y)^q$ .

$q$  - degree of polynomial.  
Homogeneous kernel.

For inhomogeneous kernel:

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

$c$  - constant.

$q$  - degree of polynomial.

③ Gaussian / RBF kernel: Radial basis Functions.

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

$\sigma$  - constant.

If  $y$  = small, RBF similar to SVM;

$y$  = large, kernel is influenced  
by more support vectors.

④ Sigmoid Kernel:  $K(x_i, y_i) = \tanh(\kappa x_i y_i - \sigma)$

Qn. Consider the 2 points  $x(1,2)$ ,  $y(2,3)$ ;  $\sigma = 1$ .

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

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

$$= e^{-\frac{(1,2)-(2,3))^2}{2}} = e^{-\frac{(1+1)^2}{2}} = e^{-\frac{(2,1)^2}{2}} \approx ?$$

$= 2/2 = 1.$

$$(1-2)^2 + (2-3)^2 = 2. \quad e^{-\frac{2}{2}} = \underline{\underline{0.36}}.$$

Qn.  $x(1,2)$      $y(2,3)$ .     $c=1$ .

Linear, homo and inhomogeneous

$$\text{Linear} = x^T y = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} 2 & 3 \end{bmatrix} = \underline{\underline{8}}.$$

$$\text{Homogeneous} = (x^T y)^2 = 8^2 = \underline{\underline{64}}.$$

$$\text{Inhomogeneous} = (c + x^T y)^2 = (1+8)^2 = 9^2 = \underline{\underline{81}}.$$

6 Feb.

Qn. i. Consider:  $y = \begin{matrix} 4 & 3 & 5 & 7 & 2 & 6 \\ \bar{y} = 28 & 52 & 6.9 & 21 & 54 \end{matrix}$

Compute the evaluation metric including MAE, MSE, RMSE & RelMSE.

ii. Apply LDA on:

$$x_1 = \{ (2,3) (3,4) (4,5) (5,6) \}$$

$$x_2 = \{ (7,8) (9,10) (10,11) \}$$

iii.  $y$ . true label: 1 0 1 0 1 0 1 0 1

$y$ . predict prob: 1 1 1 0 0 1 0 1 0 1

Evaluate confusion matrix, accuracy, precision, recall and F1 score.

$$\rightarrow i. \text{MAE} = \frac{1}{n} \sum_{i=0}^{n-1} |y_a - y_p| = \frac{1}{5} [0.2 + 0.2 + 0.1 + 0.1 + 0.5] \\ = \underline{\underline{0.18}}$$

$$\text{MSE} = \frac{1}{n} \sum_{i=0}^{n-1} (y_a - y_p)^2 = \frac{1}{5} (0.2^2 + 0.2^2 + 0.1^2 + 0.1^2 + 0.5^2) \\ = \frac{1}{5} (0.19) = \underline{\underline{0.038}}$$

$$\text{RMSE} = \sqrt{0.038} = \underline{\underline{0.195}}$$

Avg = 96.

$$\text{RelMSE} = \frac{\sum (y_a - y_p)^2}{\sum (y_a - \bar{y})^2} = \frac{0.19}{(2.56 + 0.16 + 5.76 + 6.76 + 1.96)} \\ = \frac{0.19}{19.2} = \underline{\underline{0.011}}$$

ii) LDA :  $\mathbf{x}_1 = \{ (2,3), (3,4), (4,5), (5,6) \}$   $N=4$   
 $\mathbf{x}_2 = \{ (4,8), (3,9), (9,10), (10,11) \}$ .

$$\textcircled{1} \quad \mu_1 = \frac{1}{N} \sum \mathbf{x}_i = \frac{1}{4} (11, 18) = (3.5, 4.5). \\ \mu_2 = \frac{1}{N} \sum \mathbf{x}_i = \frac{1}{4} (34, 38) = (8.5, 9.5).$$

$$\textcircled{2} \quad S_W = S_1 + S_2.$$

$$S_1 = \frac{1}{N-1} \sum (\mathbf{x} - \mu_1)(\mathbf{x} - \mu_1)^T.$$

$$= \frac{1}{3} \begin{bmatrix} [-1.5] & [-1.5 -1.5] \\ -1.5 & [0.5] & [0.5 -0.5] \\ & [0.5] & [0.5 0.5] \\ & & [1.5] & [1.5 1.5] \end{bmatrix} \\ = \frac{1}{3} \begin{bmatrix} 2.25 & 2.25 & [0.25 \cdot 0.25] \\ 2.25 & 2.25 & [0.25 0.25] \\ & & [0.25 0.25] \end{bmatrix} + \begin{bmatrix} 8.25 & 2.25 \\ 2.25 & 2.25 \end{bmatrix} \\ = \frac{1}{3} \begin{bmatrix} 5 & 5 \\ 5 & 5 \end{bmatrix} = \begin{bmatrix} 1.67 & 1.67 \\ 1.67 & 1.67 \end{bmatrix}.$$

$$S_2 = \frac{1}{N-1} \sum (\mathbf{x} - \mu_2)(\mathbf{x} - \mu_2)^T.$$

~~$$= \begin{bmatrix} 1.67 & 1.67 \\ 1.67 & 1.67 \end{bmatrix}$$~~

$$S_W = S_1 + S_2 = 2 \begin{bmatrix} 1.67 & 1.67 \\ 1.67 & 1.67 \end{bmatrix} = \begin{bmatrix} 3.34 & 3.34 \\ 3.34 & 3.34 \end{bmatrix}.$$

$$\textcircled{3} \quad S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T = \begin{bmatrix} -5 \\ -5 \end{bmatrix} \begin{bmatrix} -5 & -5 \end{bmatrix} \\ = \begin{bmatrix} 25 & 25 \\ 25 & 25 \end{bmatrix}.$$

$$\textcircled{4} \quad \text{Eigen values: } |S_W^{-1} S_B - \lambda I| = 0.$$

$|S_{ii}^{-1} S_{ii} - 1| \rightarrow$  Since  $\det(S_{ii}) = 0$ ,  $S_{ii}^{-1}$  does not exist.  $\therefore \tilde{S}_{ii}$

iii. y-true: 1 0 1 0 1 1 0 1 0 1.  
 y-pred: 1 1 1 0 0 1 0 1 0 1

Confusion matrix:

		True:			
		0	1		
Predicted	0	3 TN	1 FN	1	
	1	1 FP	5 TP	6	
				10	

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{3+5}{10} = \frac{8}{10} = 0.8 = 80\%.$$

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{5}{5+1} = \frac{5}{6} = 0.83 = 83\%.$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{5}{5+1} = \frac{5}{6} = 0.83 = 83\%.$$

$$\text{F1 Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \cdot 0.83 \cdot 0.83}{0.83 + 0.83} = \frac{0.83}{2} = 83\%.$$