

- I, Define the business Problem: Define the Problem that the company wants to solve using data mining.
- II, Build the data mining data base: collecting & integrating data from variety of sources to build a data mining database.
- III, Explore the data: using variety of data analysis techniques & visualize
- IV, Prepare the data for modeling (PreProcess): mostly feature engineering, cleaning data, removing NaN data types, etc.
- V, Build a model: using variety of machine-learning models & algorithms to identify patterns in the data that can be used to make prediction
- VI, Evaluate model: Evaluating the model's performance on a test set.
- VII, Act on the results: once the evaluation step's results are satisfying, it can be deployed to production.

Noise: unwanted signal that can interfere with the analysis of a dataset. 2

outlier: A data point that deviates significantly from the rest of the data.

a, Benefits of finding outliers: { Improved data quality, By removing outliers, we can improve the quality of ~~our~~ data
New Discoveries: By discovering them, we can gain better understanding of the data

b, Disadvantages of noise & outlier: { Reduced accuracy of data mining models.
Misleading results
Increased computational time

a, Ratio: { weight, Height. The difference between 30^{th} & 100^{th} is the same with 100^{th} & 130^{th}

Interval: { Temperature, Date & time: Same as Ratio, but, there's no true zero point on the Celsius scale, or date and time scale has no zero point.

Ordinal { Likert Scale: disagree, neither agree nor disagree, agree
Military rank: Private, Corporal, Sergeant, ...

3.61

Nominal { E.g. color, Zip Code. There is no natural order to the categories.

b, Student Number: { Tracking student progress over time
Matching students to records

Gender: { Identifying gender disparities
Studying gender differences

Median { Cannot be defined, Nominal
Can be defined, ordinal, Interval, Ratio

Mode { Cannot - - -
Can - - : all 4 categories

Average { Cannot - - : Nominal, ordinal
Can - - : Interval, Ratio (mean)

a, one-hot encoding: method of converting categorical var into binary vectors. For each category in the var, a new column is created in the dataset. The value of each column is 1 if the data point belongs to that category, and 0 otherwise.

Label encoding: method of converting categorical var into numeric values. Each category is assigned a unique integer value.

b, one-hot,

ID	Color: Blue	Green	red
1	1	0	0
2	0	0	1
3	1	0	0
4	0	1	0

Label:

ID	Color	Size	weight	shape
1	1	10	0.5	1
2	2	8	0.3	1
3	1	8	0.7	2
4	3	12	0.3	3

↑ integer encoding.

c, Both have their own advantages & disadvantages, But generally one-hot encoding is more suitable method for converting categorical into numerical values.

mean vector, $(V + (-V) + 2V + (-2V))/4 = [2.5, 0, -7.5]$, 5

$X \propto t \times$ Projection vector $= t \times [2.5, 0, -7.5] \rightarrow$ The line on which the data is is defined can be represented by this equation.

$$\beta_2 = \bar{y} - \beta_1 \bar{x}, \quad \beta_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hookrightarrow \beta_2 = \bar{y} - \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \bar{x}$$

$$\Rightarrow \text{Cov}(\beta_1, \beta_2) = E[(\beta_1 - E(\beta_1))(\beta_2 - E(\beta_2))]$$

$$= E \left\{ \left(\frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} - E \left(\frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \right) \right) \times \left(\bar{y} - \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \bar{x} - E \left(\bar{y} - \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \bar{x} \right) \right) \right\}$$

$$= \frac{E \left[\sum (y_i - \bar{y})(x_i - \bar{x}) - E \left[\sum (y_i - \bar{y})(x_i - \bar{x}) \right] \right]}{\sum (x_i - \bar{x})^2} \times \left(\bar{y} - \frac{E \left[\sum (y_i - \bar{y})(x_i - \bar{x}) \right]}{\sum (x_i - \bar{x})^2} \bar{x} - \bar{y} \right)$$

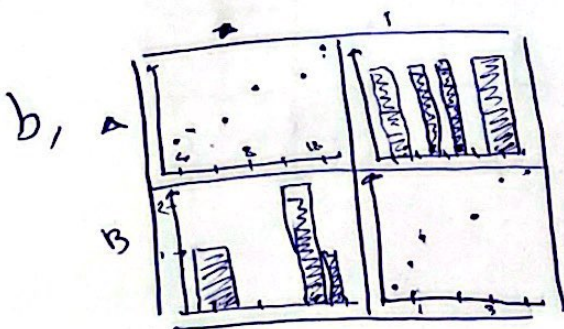
$$= \frac{\text{Cov}(x, y)}{\sum (x_i - \bar{x})^2} \cdot (-\bar{y}) = -\frac{\text{Cov}(x, y) \bar{y}}{\sum (x_i - \bar{x})^2} \rightarrow \text{if } x, y \text{ independent} = 0$$

a₁
Col 1: $4 - \frac{15.8}{Q_1, \text{median} = 8}$

Col 2: $0.5 - \frac{11.2}{Q_1, \text{median} = 2.1}$

$Q_1 = 5, Q_3 = 12$
IQR = 7, Lower bound = $5 - \frac{3}{2} \cdot 7 = -5.5$
Upper bound = $12 + \frac{3}{2} \cdot 7 = 22.5$

$Q_1 = 1, Q_3 = 3.2, \text{IQR} = 2.2$
Lower = $1 - \frac{3}{2} \cdot 2.2 = -2.3$
Upper = $3.2 + \frac{3}{2} \cdot 2.2 = 6.5$



$$SME, \frac{\text{No. of matching attributes}}{\text{Total No. of attributes}} = J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

For 2nd row $\rightarrow SME = \frac{1}{3}, J = \frac{f_{rr}}{f_{rr} + f_{rf} + f_{fr}} = 0$

b, Cosine similarity = $\frac{\sum x_i y_i}{\sqrt{\sum x_i^2} \sqrt{\sum y_i^2}} = \frac{0.1 + 1.1 + 0.1}{\sqrt{1+1+0} \sqrt{1+1+1}} = \frac{2}{\sqrt{2} \sqrt{3}} = 0.816$

Shannon's entropy $H(A) = -\log_e(\sum J(x_i, y_i)) = -\log_e(\sqrt{0.1} + \sqrt{1.1} + \sqrt{0.1})$
 $= -\log_e(2), \ln(2) = 0.693$

c, Correlation coefficient = $\frac{n \sum x_i y_i - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$
 $= \frac{4[0+0+1+1] - 3 \cdot 2}{\sqrt{(4 \cdot 3 - 9)(4 \cdot 2 - 4)}} = \frac{8-6}{\sqrt{3 \cdot 4}} = \frac{2}{2\sqrt{3}} = \frac{\sqrt{3}}{3} = 0.577$

but because it's binary data type, Pearson correlation coefficient is not appropriate. For this situation Phi correlation coefficient is used.

$\phi = \frac{ad - bc}{\sqrt{(a+b)(c+d)(a+c)(b+d)}}$, where $\begin{cases} a = \text{No. of times a1 & c2 are T} & = 2 \\ b = \text{No. of times a1 & c2 are F} & = 1 \\ c = \text{No. of times a2 & c2 are T} & = 0 \\ d = \text{No. of times a2 & c2 are F} & = 1 \end{cases}$

$\Rightarrow \phi = \frac{2 \cdot 1 - 1 \cdot 0}{\sqrt{3 \cdot 1 \cdot 2 \cdot 2}} = \frac{2}{\sqrt{12}} = 0.577 \text{ L.A.}$

d, $H(\text{class}) = -\left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4}\right) = 1$, $H(c1) = -\left(\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4}\right) = 0.56$

$H(c2) = -\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right) = 1$, $H(c3) = 1$

e, $P(m, n) = \begin{cases} e^{-(m+n)} & m > 0, n < 0 \\ 0 & \text{otherwise} \end{cases}$ $\Rightarrow \begin{cases} P(m), \int_{-\infty}^{\infty} P(m, n) dn = \int_{-\infty}^{\infty} e^{-(m+n)} dn = e^{-m} \\ P(n), \int_{-\infty}^{\infty} P(m, n) dm = \int_{-\infty}^{\infty} e^{-(m+n)} dm = e^{-n} \end{cases}$

$J(m, n) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P(m, n) \log\left(\frac{P(m, n)}{P(m)P(n)}\right) dm dn = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(m+n)} \log\left(\frac{e^{-(m+n)}}{e^{-m}e^{-n}}\right) dm dn = 0$