

Ch.2 Data Preprocessing.

① Chi-Squared Statistics.

	Observed (O)	Expected (E)	Total
Male	80	60	140
Female	40	60	100
	120	120	240

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

$$\chi^2 = \frac{(80-60)^2}{60} + \frac{(40-60)^2}{60} = 6.66 + 6.66 = 13.33$$

- ② Equi-depth partitioning - data [0, 2, 8, 10, 20, 21, 29, 29, 30]
- divide into k intervals of same size
 - divide into 3 intervals of same size = k = 3.

$$\text{width} = \frac{\text{max} - \text{min}}{k}$$

$$= \frac{30-0}{3} = \frac{30}{3} = 10$$

80 bins are (0, 10), (11, 20), (20, 30).

bin 1 → 0, 2, 8, 10

bin 2 → 20

bin 3 → 21, 29, 29, 30.

- ③ Equi-depth partitioning - data [0, 2, 8, 10, 20, 21, 29, 29, 30]
- number of values present in each bin are same
 - number of values in each bin k = 3.

bin 1 → 0, 2, 8

bin 2 → 10, 20, 21

bin 3 → 29, 29, 30.

④ Equi-depth (by means).

→ replace each value in bin with mean of that bin

bin 1 → 0, 2, 8, mean = 3.33 bin 1 → (3.33, 3.33, 3.33)

bin 2 → 10, 20, 21, mean = 17 bin 2 → (17, 17, 17)

bin 3 → 29, 29, 30, mean = 29.33 bin 3 → (29.33, 29.33, 29.33)

- ⑤ Equi-depth (by median) - Median is middle value of given range
- | | | |
|--------------------|-------------|----------------------|
| bin 1 → 0, 2, 8 | median = 2 | bin 1 → (2, 2, 2) |
| bin 2 → 10, 20, 21 | median = 20 | bin 2 → (20, 20, 20) |
| bin 3 → 29, 29, 30 | median = 29 | bin 3 → (29, 29, 29) |

- ⑥ Equi-depth (by boundaries)

Replace each value in bin with nearest boundary value.

bin 1 → 0, 2, 8	→	bin 1 → (0, 0, 8)
bin 2 → 10, 20, 21	→	bin 2 → (10, 20, 21)
bin 3 → 29, 29, 30	→	bin 3 → (29, 29, 30)

- ⑦ Co-selection analysis -

Check whether sale of ice-creams & sun-glasses are related?

ice cream sale sun-glass sale.

(A)

20
10
23
5

(B)

30
58
29
10

correlation const

$$\gamma_{A,B} = \frac{\sum (A-\bar{A})(B-\bar{B})}{(n-1) \sigma_A \sigma_B}$$

\bar{A} = mean of A $\bar{x} = \frac{\Sigma A}{n}$

σ_A = S.D. of A.

$$\sigma_A = \sqrt{\frac{\sum (A-\bar{A})^2}{n-1}}$$

$$A \quad n=4 \quad \bar{A} = \frac{58}{4} = 14.5 \quad \bar{B} = \frac{24}{4} = 18.5$$

A	B	$A-\bar{A}$	$(A-\bar{A})^2$	$B-\bar{B}$	$(B-\bar{B})^2$	$(A-\bar{A})(B-\bar{B})$
20	30	5.5	30.25	11.5	132.25	63.25
10	5	-4.5	20.25	-13.5	182.25	60.75
23	29	8.5	72.25	10.5	110.25	89.25
5	10	-9.5	90.25	-8.5	72.25	170.

$$\sigma_A = \sqrt{\frac{\sum (A-\bar{A})^2}{n-1}} = \sqrt{\frac{30.25 + 20.25 + 72.25 + 90.25}{3}} = \sqrt{\frac{219}{3}} = 8.426,$$

$$\sigma_B = \sqrt{\frac{\sum (B-\bar{B})^2}{n-1}} = \sqrt{\frac{132.25 + 182.25 + 110.25 + 72.25}{3}} = \sqrt{\frac{497}{3}} = 12.87$$

$$\gamma_{A,B} = \frac{\sum (A-\bar{A})(B-\bar{B})}{(n-1) \sigma_A \sigma_B} = \frac{63.25 + 60.75 + 89.25 + 170}{3 \times 8.426 \times 12.87} = 1.178$$

As $\gamma_{A,B} > 1$, sale of ice cream & sun glasses is positively correlated.

(8) Min-max normalization $v' = \frac{v - \min_A}{\max_A - \min_A} \{ \begin{matrix} \text{new_max}_A \\ \text{new_min}_A \end{matrix} \} + \text{new_min}_A$

income $\rightarrow \max = 98000, \min = 12000$

map income of 73600 to the range [0.0, 1.0]

$$\begin{aligned} v' &= \frac{v - \min_A}{\max_A - \min_A} (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A \\ &= \frac{73600 - 12000}{98000 - 12000} (1.0 - 0.0) + 0.0 = 0.716. \end{aligned}$$

(9) Z-score normalization (zero normalization)

$$v' = \frac{v - \bar{A}}{\sigma_A} \quad \bar{A} = \text{mean} \quad \sigma_A = \text{s.d.}$$

mean = 54000, s.d. = 15000, map 73600.

$$v' = \frac{73600 - 54000}{15000} = 1.225,$$

(10) Stratified sampling.

Group of 12 sales price records are sorted.

5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215.

Use sample of size 6 & strata, low, medium, high.

Use equi-width binning.

$k=3$, $\min=5$, $\max=215$, 6 samples needs to be drawn from 3 bins.

$$\text{width} = \frac{215 - 5}{3} = \frac{210}{3} = 70, \text{ intervals } [5, 70], [71, 140], [140, +]$$

bin 1 \rightarrow 5, 10, 11, 13, 15, 35, 50, 55

bin 2 \rightarrow 72, 92

bin 3 \rightarrow 204, 215

By monitoring value distribution in bin 1.

price 10/11 and 35/50 can be drawn as sample.

from bin 2 both records can be drawn as sample
from bin 3 both records can be drawn as sample

(11) Covariance. $\text{Cov}(A, B) = E(A \cdot B) - \bar{A} \cdot \bar{B}$.

Suppose two stocks A & B have following values in one week:

(2, 5), (3, 8), (5, 10), (4, 11), (6, 14)

Check if the stocks are affected by same industry trends, will their prices rise or fall together?

$$E(A) = \frac{2+3+5+4+6}{5} = \frac{20}{5} = 4$$

$$E(B) = \frac{5+8+10+11+14}{5} = \frac{48}{5} = 9.6$$

$$E(AB) = \frac{2 \times 5 + 3 \times 8 + 5 \times 10 + 4 \times 11 + 6 \times 14}{5} = \frac{212}{5} = 42.4$$

$$\text{Cov}(A, B) = 42.4 - 4 \times 9.6 = 4.$$

∴ Thus, A & B move together since $\text{Cov}(A, B) > 0$.

(12) Dispersion of data → .

following data points represent number of classes each teacher at BITS teaches: 5, 7, 5, 1, 1, 9, 4, 5, 3, 4, 6.

find out measures of dispersion of data.

→ first arrange in ascending order

1, 1, 3, 4, 4, 5, 5, 5, 6, 7, 9

Median = middle value if ~~the~~ number of data points are odd
= 5

Q₁ = middle value of data points ^{left side} of the median
25th percentile = 1, 1, 3, 4, 4

= 3.

Q₃ = middle value of data points ^{right side} of the median
75th percentile = 5, 5, 6, 7, 9

= 8

IQR = Q₃ - Q₁ = 8 - 3 = 5.

Inter quartile range

five number summary = min, Q₁, median, Q₃, max
= {1, 3, 5, 8, 9}.

Outliers = a value higher/lower than $1.5 \times \text{IQR}$

$$2 + 1.5 \times \text{IQR} = 1.5 \times 5 = 7.5$$

$$1.5 \times \text{IQR} = 1.5 \times 3 = 4.5$$

Outliers = value which is located $1.5 \times \text{IQR}$ or more below Q₁ or value which is located $1.5 \times \text{IQR}$ or more above Q₃.

All data values which are below $Q_1 - 1.5 \times \text{IQR}$ or are above $Q_3 + 1.5 \times \text{IQR}$ are outliers. ∴ No outlier is above dataset.

* The Type of Data

- Type of attribute - Qualitative, Quantitative
 - Nominal, Ordinal, Interval, Ratio.
 - Discrete, Continuous.
 - Asymmetric, Symmetric.
- Type of Datasets -
 - Recorded data - Market Basket, Data matrix, Sparse Data matrix
 - Graph based data - Data with "rel" among objects
Data with objects that are graph.
 - Ordered data - Sequential data, sequence data,
Time series data, Spatial data.

Characteristics of Data Sets -

- Dimensionality, Sparsity, Resolution.

* Data Quality

- Measurement & data collection issues.
- Noise & Artifact
- Precision, Bias, Accuracy.
- Outliers.
- Missing values
- Inconsistent values
- Duplicate data.

* Data Preprocessing.

Data cleaning.

Missing values - ignore tuple, use mean/mode, estimate using models

Noisy data - data smoothing - binning, clustering,
combined human & computer inspection,
regression

Inconsistent data - manual correction

knowledge engineering tools for constraint

Data Integration

entity identification - problems cust-ID & cust-no

redundancy - correlation.

conflict detection & resolution of data value conflicts

correlation between attribute A & B.

$$\rho_{A,B} = \frac{\sum (A-\bar{A})(B-\bar{B})}{(n-1) \sigma_A \sigma_B}$$

\bar{A} & \bar{B} are mean of A & B

σ_A & σ_B are s.d. of A & B.

$$\bar{A} = \frac{\sum A}{n}$$

$$\sigma_A = \sqrt{\frac{(A-\bar{A})^2}{n-1}}$$

- Data Transformation

smoothing

Aggregation

Generalization - concept hierarchy generation

Normalization - min-max, z-score, decimal scaling

Attribute (feature) construction

- Data Reduction

Data cube aggregation

Dimension reduction - stepwise fwd selection, backward elimination decision tree

Data compression - wavelet transform, PCA, factor analysis

Numerosity reduction - regression, histogram, clustering

Concept hierarchy generation Sampling

* Data Exploration - EDA

- Summary Stat

- Data similarity & dissimilarity measures

- Data visualization.

- Summary Stat

frequency & mode

Percentile

Measures of location - mean & median central points

Measures of spread - range & variation

Measures of spread - range & variation

Multivariate summary stat - mean, covariance, correlation

- Data similarity & dissimilarity measures.

Similarity, Dissimilarity, Proximity.

Similarity, Dissimilarity b/w simple attributes

Dissimilarity

$$\text{Nominal} \quad d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$$

$$s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$$

$$\text{Ordinal} \quad d = \frac{|x - y|}{n-1}$$

$$s = 1 - d$$

$$s = -d, s = \frac{1}{1+d}, s = e^{-\frac{d}{2}}$$

$$\text{Interval or Ratio} \quad d = |x - y|$$

- Dissimilarity between data objects

$$\text{Euclidean} \quad \text{Distances.} \quad d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

$$\text{Minkowski} \quad d(x, y) = \sqrt[n]{\sum_{k=1}^n (x_k - y_k)^\alpha}$$

$$\text{Manhattan} \quad d(x, y) = \sum_{k=1}^n |x_k - y_k|$$

$$\text{Supremum} \quad d(x, y) = \lim_{\alpha \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^\alpha \right)^{\frac{1}{\alpha}}$$

* Examples of Fuzzy Similarity Measures

Binary data. - similarity coefficients.

* Simple Matching Coeff SMC = $\frac{\text{number of matching attribute values}}{\text{no. of attributes}}$

$$SMC = \frac{f_{11} + f_{00}}{f_{01} + f_{10} + f_{11} + f_{00}}$$

- Jaccard Coeff. $J = \frac{\text{number of matching presences}}{\text{number of attributes not involved in few.}}$

Asymmetric attr.

$$= \frac{f_{11}}{f_{01} + f_{10} + f_{11}}$$

- Cosine Similarity - term-document metric.

$$\cos(\alpha, y) = \frac{\alpha \cdot y}{\|\alpha\| \|\gamma\|} \quad \alpha \cdot y = \sum_{k=1}^n \alpha_k y_k$$

$$\|\alpha\| = \sqrt{\sum_{k=1}^n \alpha_k^2}$$

- Extended Jacob Coeff Tanimoto Coeff. $EJ(\alpha, y) = \frac{\alpha \cdot y}{\|\alpha\|^2 + \|\gamma\|^2 - \alpha \cdot y}$

term documents memory with binary attributes.

- Correlation $-1 \rightarrow 1$

Pearson's correlation coeff.

$$\cos(\alpha, y) = \frac{\text{covariance } (\alpha, y)}{\text{s.d. } (\alpha) * \text{s.d. } (y)} = \frac{S_{\alpha y}}{S_{\alpha} \cdot S_y}$$

$$\text{covariance } (\alpha, y) = S_{\alpha y} = \frac{1}{n-1} \sum_{k=1}^n (\alpha_k - \bar{\alpha})(y_k - \bar{y})$$

$$\text{s.d. } (\alpha) = S_{\alpha} = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (\alpha_k - \bar{\alpha})^2}$$

$$\text{s.d. } (y) = S_y = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$\bar{\alpha} = \frac{1}{n} \sum_{k=1}^n \alpha_k \quad \bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$$

* Similarity for heterogeneous attributes -

① for attribute k, compute similarity $s_k(\alpha, y)$ in range [0, 1]

② Define indicator variable, d_k for kth attribute

③ Define indicator variable, d_k for kth attribute & both objects have a value of 0 or if one of objects has missing value

$$d_k = \begin{cases} 0 & \text{if } k^{\text{th}} \text{ attribute is asymmetric \& both objects have a value of 0 or if one of objects has missing value} \\ 1 & \text{otherwise.} \end{cases}$$

④ Compute overall similarity betw two data objects.

$$\text{similarity } (\alpha, y) = \frac{\sum_{k=1}^n d_k \cdot s_k(\alpha, y)}{\sum_{k=1}^n d_k}$$

If weights are used.

$$\text{similarity } (\alpha, y) = \frac{\sum_{k=1}^n d_k \cdot w_k \cdot s_k(\alpha, y)}{\sum_{k=1}^n d_k}$$

* Precision & Bias.

Precision - closeness of repeated measurements of same quantity to one another

Bias - systematic variation of measurements from the quantity being measured.

A device is being used to measure the speed of the vehicles.

Vehicles speed is 10 km/hr . On trial five measurements are taken for the same vehicle, which are $\{10.1, 10.3, 10.0, 10.2, 10.1\} \text{ km/hr}$.

Find precision & bias for the device being used for measurement.

$$\rightarrow \text{mean} = \frac{10.1 + 10.3 + 10.0 + 10.2 + 10.1}{5} = \frac{50.7}{5} = 10.14 \text{ km/hr.}$$

$$n = \text{no. of trials} = 5.$$

$$\text{s.d.} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2} = \sqrt{\frac{1}{4} \left[(10.1 - 10.14)^2 + (10.3 - 10.14)^2 + (10.0 - 10.14)^2 + (10.2 - 10.14)^2 + (10.1 - 10.14)^2 \right]} \\ = \sqrt{\frac{1}{4} \times 0.052} = 0.1140$$

$$\text{Precision} = \text{s.d.} = 0.1140.$$

Precision = s.d. \div mean value of quantity measured

Bias = difference betⁿ mean of values & known value of quantity measured

$$= 10.14 - 10 = 0.14$$

* Similarity & Dissimilarity between Simple Attributes

Attribute

Type,

Nominal.

Dissimilarity

$$d = \begin{cases} 0 & \text{if } x=y \\ 1 & \text{if } x \neq y \end{cases}$$

Similarity,

$$S = \begin{cases} 1 & x=y \\ 0 & x \neq y. \end{cases}$$

$$S = 1 - d.$$

$$S = -d.$$

$$S = \frac{1}{1+d}$$

$$S = e^{-d}.$$

Ordinal

$$d = \frac{|x-y|}{n-1}$$

Interval or Ratio

$$d = |x-y|.$$

Consider following dataset about students studying at BITE
find out dissimilarity & similarity for each attribute in it.

Id.	Name.	Grade-Midsem	Grade-Compre	Grade-Assessment
1	Pravin	10	8	Fair
2	Pawar	8	6	Excellent
3	Pravin	10	8	Good.

→ Name is nominal attribute, & categorical.

Grade-Midsem & Grade-Compre are intervals & numeric.

Grade-Assessment is categorical & ordinal.

for Name — $d(1,2) = 1$. as Pravin \neq Pawar $s(1,2) = 0$
 $d(1,3) = 0$ as Pravin = Pravin $s(1,3) = 1$
 $d(2,3) = 1$. as Pawar \neq Pravin $s(2,3) = 0$

for Grade-Midsem

$$\begin{array}{lll} d(1,2) = 2 & s = -d & s(1,2) = -2 \\ d(1,3) = 0 & \text{then} & s(1,3) = 0 \\ d(2,3) = 2 & & s(2,3) = -2 \end{array} \quad \begin{array}{l} s = \frac{1}{1+d} \\ s(1,2) = \frac{1}{1+2} = \frac{1}{3} = 0.33 \\ s(1,3) = \frac{1}{1+0} = 1 = 1 \\ s(2,3) = \frac{1}{1+2} = \frac{1}{3} = 0.33 \end{array}$$

for Grade-Compre.

$$\begin{array}{lll} d(1,2) = 2 & s = \frac{1}{1+d} & s(1,2) = \frac{1}{1+2} = \frac{1}{3} = 0.33 \\ d(1,3) = 0 & & s(1,3) = \frac{1}{1+0} = 1 = 1 \\ d(2,3) = 2 & & s(2,3) = \frac{1}{1+2} = \frac{1}{3} = 0.33 \end{array}$$

for Grade-Assessment

lets assume order as. Fair < Good < Excellent
 $\begin{array}{ccc} 1 & 2 & 3 \end{array}$

$$d(1,2) = \frac{|1-2|}{3-1} = \frac{1}{2} = 0.5 \quad s(1,2) = 1-d = 1-0.5 = 0.5$$

$$d(1,3) = \frac{|1-3|}{3-1} = \frac{2}{2} = 1 \quad s(1,3) = 1-d = 1-1 = 0$$

$$d(2,3) = \frac{|2-3|}{3-1} = \frac{1}{2} = 0.5 \quad s(2,3) = 1-d = 1-0.5 = 0.5$$

Similarity Measures

Distance between Vectors

$$x = \{x_1, x_2, x_3, \dots, x_n\}$$

$$y = \{y_1, y_2, y_3, \dots, y_n\}.$$

① Minkowski distⁿ.

$$d(x, y) = \left\{ \sum_{k=1}^n |x_k - y_k|^{\alpha} \right\}^{\frac{1}{\alpha}}$$

② Hamming distⁿ (Manhattan, L1) $\rightarrow \alpha = 1$.

$$\therefore d(x, y) = \sum_{k=1}^n |x_k - y_k|.$$

③ Euclidean distⁿ (L2) $\rightarrow \alpha = 2$

$$\therefore d(x, y) = \sqrt{\sum_{k=1}^n |x_k - y_k|^2}.$$

④ Supremum distⁿ

$$d(x, y) = \lim_{\alpha \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^{\alpha} \right)^{\frac{1}{\alpha}}.$$

$$\text{e.g. } x = \{1, 2, 3, 4\}$$

$$y = \{5, 3, 2, 1\}$$

Hamming distⁿ $\rightarrow \alpha = 1$.

$$\therefore d(x, y) = |1-5| + |2-3| + |3-2| + |4-1| \\ = 4 + 1 + 1 + 3 = 9$$

Euclidean distⁿ $\rightarrow \alpha = 2$

$$\therefore d(x, y) = \sqrt{|1-5|^2 + |2-3|^2 + |3-2|^2 + |4-1|^2} \\ = \sqrt{4^2 + 1^2 + 1^2 + 3^2} = \sqrt{27}.$$

Proximity Measures

Simple Matching Coefficient (SMC)

① Simple Matching Coefficient (SMC) consists of n binary attributes.

x & y are two objects that consists of n binary attributes.

f_{00} = number of cases where x is 0 & y is 0

$$f_{01} = \text{_____} \quad x \text{ is } 0 \& y \text{ is } 1$$

$$f_{10} = \text{_____} \quad x \text{ is } 1 \& y \text{ is } 0$$

$$f_{11} = \text{_____} \quad x \text{ is } 1 \& y \text{ is } 1$$

$$\text{SMC} = \frac{\text{number of matching attribute values}}{\text{number of attributes}}$$

$$= \frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}}$$

$$x = \{1, 0, 1, 0\}$$

$$f_{00} = 1$$

$$f_{01} = 1$$

$$y = \{0, 1, 1, 0\}$$

$$f_{10} = 1$$

$$f_{11} = 1$$

$$\text{SMC} = \frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}} = \frac{1+1}{4} = \frac{2}{4} = 0.5$$

② Jaccard coefficient

$J = \frac{\text{number of matching presences}}{\text{number of attributes not involved in no matches}}$

$$= \frac{f_{11}}{f_{01} + f_{10} + f_{11}}$$

$$\begin{matrix} x = \\ \{1, 0, 1, 0\} \end{matrix}$$

$$y = \{0, 1, 1, 0\}$$

$$J = \frac{f_{11}}{f_{01} + f_{10} + f_{11}} = \frac{1}{1+1+1} = \frac{1}{3} = 0.33.$$

$$\begin{matrix} f_{00} = 1 \\ f_{01} = 1 \\ f_{10} = 1 \\ f_{11} = 1 \end{matrix}$$

③ Cosine Similarity -

$$\cos(\alpha, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

where - indicates dot vector product

$$x \cdot y = \sum_{k=1}^n \alpha_k y_k.$$

$$\|x\| = \sqrt{\sum_{k=1}^n \alpha_k^2}.$$

$$x = \{3, 2, 0, 5\}$$

$$y = \{1, 0, 0, 0\}$$

$$x \cdot y = 3*1 + 2*0 + 0*0 + 5*0 = 3 + 0 + 0 + 0 = 3$$

$$\|x\| = \sqrt{3^2 + 2^2 + 0^2 + 5^2} = \sqrt{9+4+25} = \sqrt{38}$$

$$\|y\| = \sqrt{1^2 + 0^2 + 0^2 + 0^2} = \sqrt{1} = 1$$

$$\therefore \cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|} = \frac{3}{\sqrt{38} \cdot 1} = \frac{3}{\sqrt{38}} = 0.486$$

④ Extended Jaccard Coefficient (Tanimoto coefficient)

$$EJ(x, y) = \frac{x \cdot y}{\|x\|^2 + \|y\|^2 - x \cdot y}$$

$$x = \{3, 2, 0, 5\}$$

$$y = \{1, 0, 0, 0\}$$

~~$$EJ(x, y) = \frac{3}{(\sqrt{38})^2 - 3} = \frac{3}{38 - 3} = \frac{3}{35} = 0.0857$$~~

$$\frac{3}{36}$$

⑤ Correlation

$$\text{corr}(x, y) = \frac{\text{covariance } (\alpha, y)}{S.D.(\alpha) * S.D.(y)}$$

$$\text{covariance } (\alpha, y) = \frac{1}{n-1} \sum_{k=1}^n (\alpha_k - \bar{\alpha})(y_k - \bar{y})$$

$$\sigma(\alpha) = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (\alpha_k - \bar{\alpha})^2}$$

$$\sigma(y) = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$\bar{\alpha} = \frac{1}{n} \sum_{k=1}^n \alpha_k$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$$

- Q. Suppose two stocks A & B have following closing values in one week. $(2, 5), (3, 8), (5, 10), (4, 11), (6, 14)$.
- ④ Check if stocks are affected by same industry trends or not.
 - ⑤ Check if one stock can be used to predict value of other stock or not.

→ ④ Co-variance - stocks vary together or not.

If it's zero, they don't vary together.

If it's one, they vary together, will be affected by some trend.

$$\text{Covariance } (\bar{x}, \bar{y}) = \frac{1}{n-1} \sum_{k=1}^n (\bar{x}_k - \bar{x})(\bar{y}_k - \bar{y}).$$

$$n = 5.$$

$$\bar{x} = \{2, 3, 5, 4, 6\}, \quad \bar{x} = 4$$

$$\bar{y} = \{5, 8, 10, 11, 14\}, \quad \bar{y} = 9.6.$$

$$\begin{aligned} \text{Covariance } (\bar{x}, \bar{y}) &= \frac{1}{4} \left[(2-4)(5-9.6) + (3-4)(8-9.6) + (5-4)(10-9.6) \right. \\ &\quad \left. + (4-4)(11-9.6) + (6-4)(14-9.6) \right] \\ &= \frac{1}{4} \left[-4.8 + 1.6 + 0.4 + 0 + 8.8 \right] \\ &= 2.0 \end{aligned}$$

A positive value, hence both stocks vary together.

⑤ Correlation - stocks are related or not.

$$\text{Correlation } (\bar{x}, \bar{y}) = \frac{\text{Covariance } (\bar{x}, \bar{y})}{\text{s.d. } (\bar{x}) \cdot \text{s.d. } (\bar{y})}$$

$$\begin{aligned} \text{s.d. } (\bar{x}) &= \sqrt{\frac{1}{n-1} \sum_{k=1}^n (\bar{x}_k - \bar{x})^2} \\ &= \sqrt{\frac{1}{4} [(2-4)^2 + (3-4)^2 + (5-4)^2 + (4-4)^2 + (6-4)^2]} \\ &= \sqrt{\frac{1}{4} [4+1+1+0+4]} = \sqrt{2.5} = 1.58 \end{aligned}$$

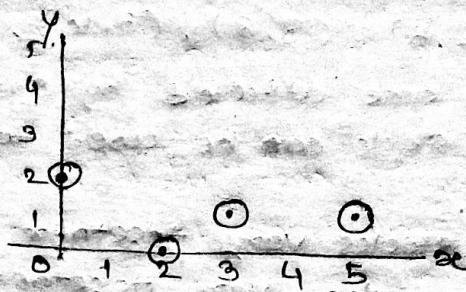
$$\begin{aligned} \text{s.d. } (\bar{y}) &= \sqrt{\frac{1}{n-1} \sum_{k=1}^n (\bar{y}_k - \bar{y})^2} \\ &= \sqrt{\frac{1}{4} [(5-9.6)^2 + (8-9.6)^2 + (10-9.6)^2 + (11-9.6)^2 + (14-9.6)^2]} \\ &= \sqrt{\frac{1}{4} [21.16 + 2.56 + 0.16 + 1.96 + 19.36]} = \sqrt{11.3} = 3.36 \end{aligned}$$

$$\text{Correlation } (\bar{x}, \bar{y}) = \frac{5}{1.58 \times 3.36} = 0.94$$

As correlation is 0.94 i.e. 94% the two stocks are strongly correlated with each other.

Distance calculation

Point	x	y
P1	0	2
P2	2	0
P3	3	1
P4	5	1



$$\text{Minkowski distance } d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^s \right)^{\frac{1}{s}}$$

$s=1$, Manhattan distⁿ / L, norm dist?

$$d(P_1, P_2) = \sum_{k=1}^2 |x_k - y_k| = |0-2| + |2-0| = 4$$

$$d(P_1, P_3) = |0-3| + |2-1| = 4$$

$$d(P_1, P_4) = |0-5| + |2-1| = 6$$

L ₁	P ₁	P ₂	P ₃	P ₄
P ₁	0	4	4	6
P ₂	4	0	2	4
P ₃	4	2	0	2
P ₄	6	4	2	0

$s=2$, Euclidean distⁿ

$$d(P_1, P_2) = \sqrt{\sum_{k=1}^2 (x_k - y_k)^2} = \sqrt{(0-2)^2 + (2-0)^2} = \sqrt{4+4} = \sqrt{8} = 2\cdot\sqrt{2}$$

$$d(P_1, P_3) = \sqrt{(0-3)^2 + (2-1)^2} = \sqrt{9+1} = \sqrt{10} = 3\cdot\sqrt{10}$$

$$d(P_2, P_3) = \sqrt{(2-3)^2 + (0-1)^2} = \sqrt{1+1} = \sqrt{2} = 1\cdot\sqrt{2}$$

$s=\infty$, supremum distⁿ

$$d(x, y) = \lim_{s \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^s \right)^{\frac{1}{s}} = \max |x_k - y_k|$$

$$d(P_1, P_2) = \max \{ |0-2|, |2-0| \} = \max \{ 2, 2 \} = 2$$

$$d(P_1, P_3) = \max \{ |0-3|, |2-1| \} = \max \{ 3, 1 \} = 3$$

$$d(P_2, P_3) = \max \{ |2-3|, |0-1| \} = \max \{ 1, 1 \} = 1$$

$$d(P_2, P_4) = \max \{ |2-5|, |0-1| \} = \max \{ 3, 1 \} = 3$$

Euclidean Distance.

L ₂	P ₁	P ₂	P ₃	P ₄
P ₁	0	2.8	3.2	5.1
P ₂	2.8	0	1.4	3.2
P ₃	3.2	1.4	0	2.0
P ₄	5.1	3.2	2.0	0

Supremum Dist

L _∞	P ₁	P ₂	P ₃	P ₄
P ₁	0	2	3	5
P ₂	2	0	1	3
P ₃	3	1	0	2
P ₄	5	3	2	0

★ Measures of spread : Range & Variance.

Assume that feature age contains following four values in it — $\{20, 25, 10, 20\}$. Compute all measures of spread.

$$\bar{x} = \text{age} = \{20, 25, 10, 20\}$$

$$n = \text{number of values} = 4$$

$$\bar{x} = \text{mean} = \frac{20+25+10+20}{4} = \frac{75}{4} = 18.75$$

$$\textcircled{1} \text{ Range: } \text{range}(\bar{x}) = \max(\bar{x}) - \min(\bar{x}) = \bar{x}_n - \bar{x}_1, \\ = 25 - 10 = 15.$$

$$\textcircled{2} \text{ Variance: } S_{\bar{x}}^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2$$

$$= \frac{1}{4-1} [(20-18.75)^2 + (25-18.75)^2 + (10-18.75)^2 + (20-18.75)^2]$$

$$= \frac{1}{3} (118.75) = 39.58.$$

$$\textcircled{3} \text{ Standard deviation: } \sigma(\bar{x}) = \sqrt{S_{\bar{x}}^2} \\ = \sqrt{39.58} = 6.241$$

$$\textcircled{4} \text{ Absolute Average Deviation (AAD):}$$

$$AAD(\bar{x}) = \frac{1}{n} \sum_{i=1}^n |\bar{x}_i - \bar{x}|$$

$$= \frac{1}{3} [|20-18.75| + |25-18.75| + |10-18.75| + |20-18.75|]$$

$$= \frac{1}{3} [1.25 + 6.25 + 8.75 + 1.25] = 1.071$$

$$\textcircled{5} \text{ Median Absolute Deviation (MAD):}$$

$$MAD(\bar{x}) = \text{median} (\{ |\bar{x}_1 - \bar{x}|, |\bar{x}_2 - \bar{x}|, \dots, |\bar{x}_n - \bar{x}| \})$$

$$= \text{median} (\{ |20-18.75|, |25-18.75|, |10-18.75|, |20-18.75| \})$$

$$= \text{median} (\{ 1.25, 6.25, 8.75, 1.25 \})$$

$$= \frac{6.25 + 8.75}{2} = 7.5$$

$$\textcircled{6} \text{ Interquartile Range IQR:}$$

$$IQR(\bar{x}) = \bar{x}_{75\%} - \bar{x}_{25\%}.$$

$$\text{Arrange } \bar{x} \text{ in increasing order} = \{10, 20, 20, 25\}.$$

$$\bar{x}_{25\%} = 10.5 = \frac{10+20}{2}$$

$$\bar{x}_{75\%} = 22.5 = \frac{20+25}{2}$$

$$IQR(\bar{x}) = 22.5 - 10.5 = 12$$