

Data transformation by normalization →

The Measurement unit used can affect the data Analysis. For ex:- Changing Measurement units from meters to inches for height, or from kilograms to pounds for weight, may lead to very different results. In general, expressing an attribute in smaller units will lead to larger range for that attribute, and thus tend to give such an attribute greater effect or "weight" to help avoid dependence on the choice of M/S units, the data should be normalized or standardized. This involves transforming the data to fall within a smaller or common range such as $[-1, 1]$ or $[0.0, 1.0]$.

Normalization of data attempts to give all attribute an equal weight. Normalization is particularly useful for classification Algo. involving neural networks or distance Measurements such as nearest neighbour classification and clustering. If using the neural network back Propagation algorithm for classification Mining, normalizing the I/P values for each attribute measured in the Train tuples will help speed up the learn phase.

For distance based Methods, normalizing helps prevent attributes with initially large ranges (e.g. income) from overweighing attributes with initially smaller ranges (e.g. binary attributes). It is also useful when given no prior knowledge of data.

There are many Methods for data normalization. We study min max normalization, z-score normalization, and normalization by decimal scaling. For our discussion, let A be a numeric attribute with n observed values, v_1, v_2, \dots, v_n .

Min-max normalization \rightarrow it performs a linear Transformation on the original data, suppose \min_A and \max_A are the minimum and maximum values of an attribute, A. min-max normalization maps a value, v_i of A to v_i' in the range $[\text{new-min}_A, \text{new-max}_A]$ by computing

$$v_i' = \frac{v_i - \min_A}{\max_A - \min_A} (\text{new-max}_A - \text{new-min}_A) + \text{new-min}_A$$

min-max normalization preserves the relationship among the original data values, it will encounter an "out of bounds" error if a future HP case for normalization falls outside of the original data range for A.

Ex: \rightarrow Suppose that the minimum and maximum values for the attribute Income are \$12,000 and \$98,000, respectively. we would like to map income to the range $[0.0, 1.0]$ by min-max normalization, a value of \$73,600 for income is

$$\text{Transformed to } \frac{73600 - 12000}{98000 - 12000} (1.0 - 0) + 0 = 0.716$$

in z-score normalization (or zero Mean normalization), the value for an attribute, A, are normalized based on the mean ~~the~~ and SD of A. A value, v_i of A is normalized to v_i' by computing

$$v_i' = \frac{v_i - \bar{A}}{\sigma_A}$$

where \bar{A} and σ_A are the Mean and Standard Deviation (SD), respectively of Attribute A. This method of normalization is useful when the actual minimum and maximum of attribute A are unknown, or when there are outliers that are distorted by the min-max normalization.

Ex:- Z-Score normalization:- Suppose that the Mean and SD of the values for the attribute Income are \$54,000 and \$16,000, respectively. With Z-score normalization, a value of 73,000 \$ for income is transformed to $\frac{73000 - 54000}{16000} = 1.225$

The Mean absolute deviation of A

$$S_A = \frac{1}{n} (|v_1 - \bar{A}| + |v_2 - \bar{A}| + |v_3 - \bar{A}| + \dots + |v_n - \bar{A}|)$$

Thus Z-score normalization using the Mean absolute deviation is $v_i' = \frac{v_i - \bar{A}}{S_A}$

The Mean Absolute deviation, S_A is more robust to outlier than the SD, σ_A , here effects of outlier are ~~less~~ somewhat reduced because the deviation from the mean (i.e. $|x_i - \bar{x}|$) are not squared.

Normalization by decimal scaling: \rightarrow it normalizes by

moving the decimal point of values of attribute A.

Thus number of decimal points depends on the maximum absolute value of A. A value v_i of A is normalized to v_i' by computing.

$$v_i' = \frac{v_i}{10^J}$$

where J is the smallest integer such that $\max(|v_i'|) < 1$

Suppose that a recorded values of A ranges from -986 to 917. The maximum absolute value of A is 986. To normalize by decimal scaling, we therefore divide each value by 1000 (i.e. $j=3$) so that -986 normalizes to -0.986 and 917 normalizes to 0.917

Question:- Use these Methods to normalize the following group of data:

200, 300, 400, 600, 1000

- min-max normalization by setting min=0, max=1
- z-score normalization
- z-score normalization using the mean absolute deviation instead of standard deviation.
- normalization by decimal scaling.

Sol:- ① Min max normalization

$$V_i' = \frac{V_i - \min A}{\max A - \min A} (\text{new max } A - \text{new min } A) + \text{new min } A$$

$$\min A = 200, \max A = 1000$$

$$\bar{A} = \frac{200 + 300 + 400 + 600 + 1000}{5} = \frac{2500}{5} = 500$$

$$V_1' = 0$$

$$V_2' = \frac{300 - 200}{1000 - 200} (1 - 0) + 0 = \frac{100}{800} = 0.125$$

$$V_3' = \frac{400 - 200}{(1000 - 200) / 200} = \frac{200}{200} = 0.25$$

$$V_4' = \frac{600 - 200}{800} (1) + 0 = \frac{400}{800} = 0.5$$

and so on.

$$V_5' = \frac{1000 - 200}{800} (1) + 0 = \frac{800}{800} = 1$$

② Z-Score Normalization

$$V_i' = \frac{V_i - \bar{A}}{\sigma_A}$$

$$\sigma_A = \sqrt{\frac{1}{n} \left(\sum_{i=1}^n x_i^2 \right) - \bar{x}^2} \Rightarrow \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$$

$$\bar{A} = 500$$

$$\sigma_A = \frac{1}{5} (200^2 + 300^2 + 400^2 + 600^2 + 1000^2)$$

$$- 500^2$$

$$\Rightarrow \frac{1}{5} (40000 + 90000 + 160000 + 360000 + 1000000)$$

$$\Rightarrow \frac{1}{5} (1650000) - 250000$$

$$\Rightarrow \frac{330000 - 250000}{5}$$

$$\Rightarrow \sqrt{80000} \Rightarrow 282.84$$

$$V_4' = \frac{600 - 500}{\sqrt{80000}} = \frac{100}{282.84} = 0.353$$

③ ~~Normal~~ Mean absolute deviation

mean absolute deviation

$$\begin{aligned} \sigma_A &= \frac{1}{N} (|200 - 500| + |300 - 500| \\ &\quad + |400 - 500| + |600 - 500| \\ &\quad + |1000 - 500|) \\ &= \frac{1}{N} (300 + 200 + 100 + 100 + 500) \\ &= \frac{1}{5} (1200) = 240 \end{aligned}$$

$$V_4' = \frac{600 - 500}{240} = \frac{100}{240} = 0.416$$

④ Normalization by decimal scaling \rightarrow

$$V_4' = \frac{600}{10^3} \Rightarrow \underline{0.6}$$

Graphic Displays of Basic Statistical Description of data:-

These include quantile plots, quantile-quantile plots, histograms and scatter plots. Such graphs are helpful for visual inspection of data, which is useful for data preprocessing. The first three of these show univariate distribution (i.e. data for one attribute), while scatter plots show bivariate distribution (i.e. involving two attributes).

Histograms:-

Histograms (or frequency histograms) are at least a century old and widely used method. "Histo" means pole or mast, and "gram" means chart, so a histogram is a chart of poles. Plotting histograms is a graphical method for summarizing the distribution of a given attribute, X . If X is nominal, such as automobile-model or item-type, then a pole or vertical bar is drawn for each known value of X . The height of the bar indicates the frequency (i.e. count) of that X -value. The resulting graph is more commonly known as a bar chart.

If X is numeric, the term histogram is preferred. The range of value for X is partitioned into disjoint consecutive subranges. The subranges, referred to as bucket or bins, are disjoint subsets of the data distribution for X . The length of a bucket is known as width. Typically, the buckets are of equal width.

For example, a price attribute with a value range of \$1 to \$200 can be partitioned into subranges 1 to 20, 21 to 40, 41 to 60

For each subrange, a bar is drawn with a height that represents the total count of items observed within the subrange.

Table 1: — A set of unit price data for items sold at a branch of All electronics

Unit Price \$	Count of items sold
40	275
43	300
47	250
—	—
74	360
75	515
78	540
—	—
115	320
117	270
120	350

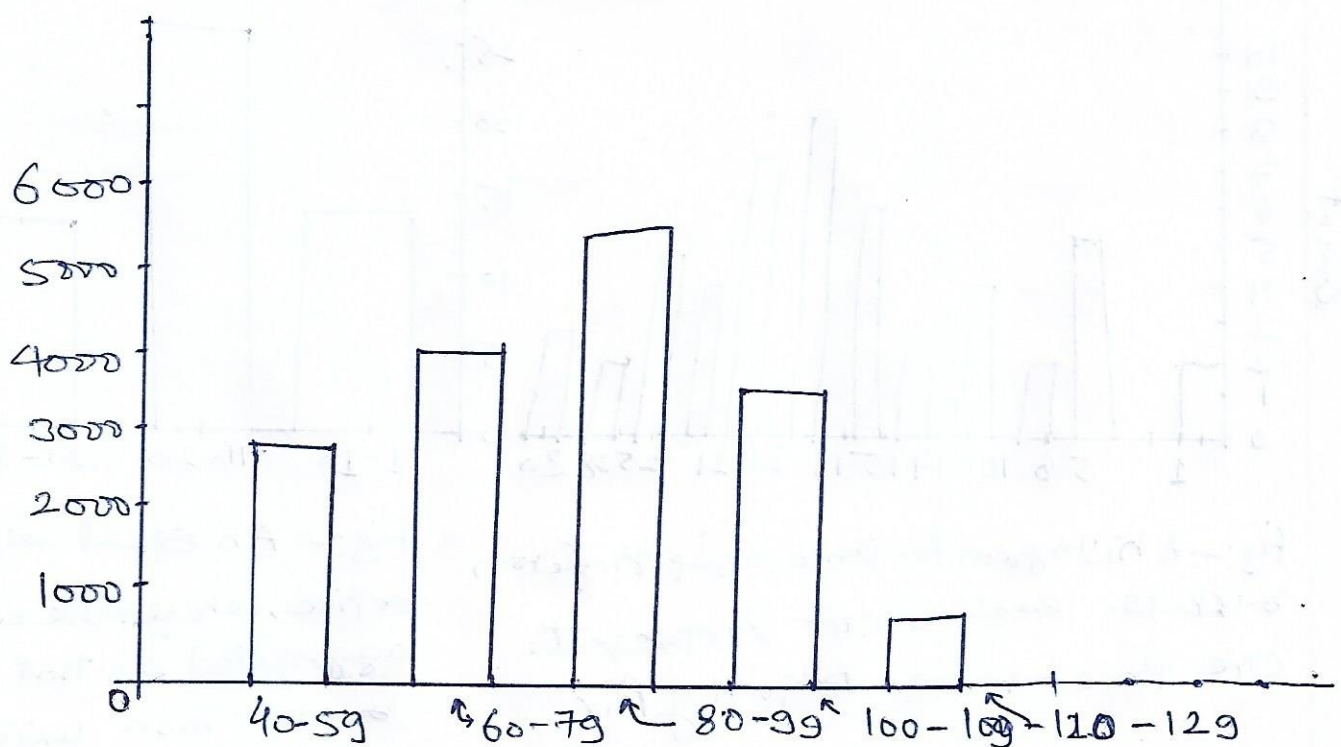


Fig:- A histogram for above data table.

histograms use binning to approximate data distributions and are a popular form of data reduction. A histogram for any attribute, A, partitions the data distribution of A into disjoint subsets referred to as bucket or bins. If each bucket represents only a single attribute-value/frequency pair, the bucket are called singleton buckets, bucket instead represents continuous ranges for the given attribute.

for Ex:- 1, 1, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30.

Fig (a) shows a histogram for the data using singleton buckets. To further reduce the data, it is common to have each bucket denote a continuous value range for the given attribute. in Fig (b) each bucket represents a different 10\$ range of Price.

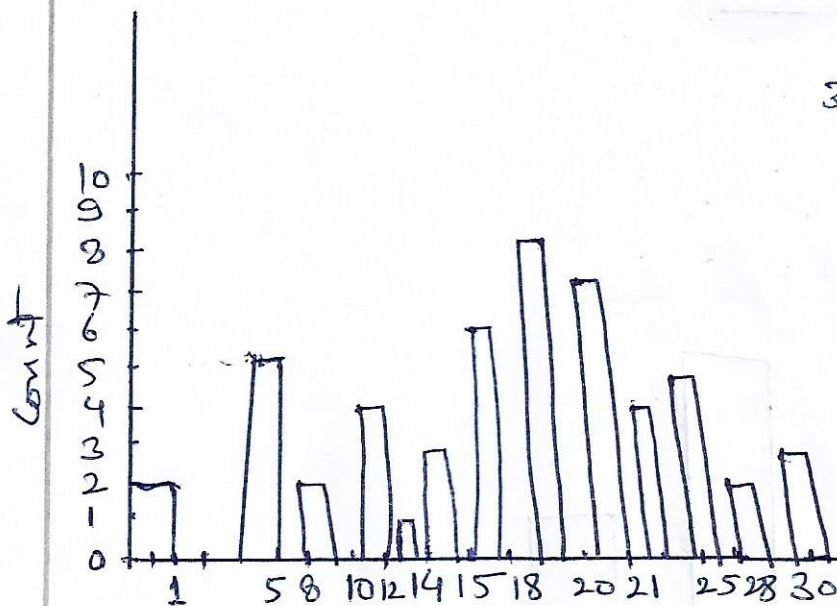


Fig:- A histogram for Price using singleton buckets- each bucket represents one Price-value frequency pair

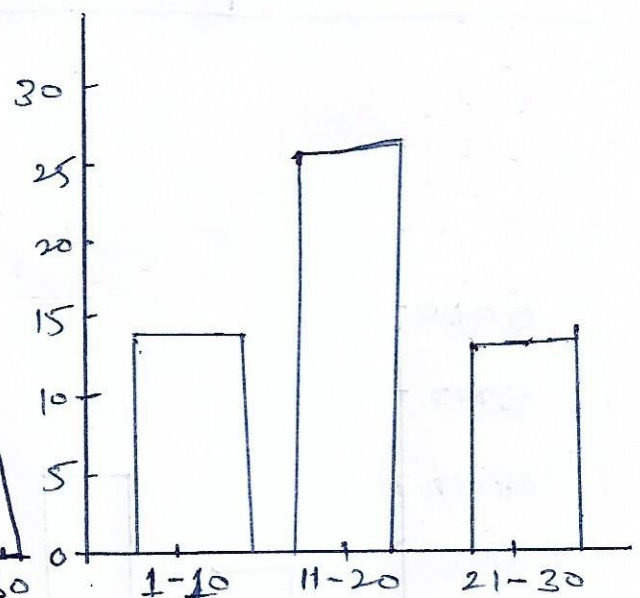


Fig:- An equal width histogram for Price, where value are aggregated so that each bucket has a uniform width of \$10.

Scatter plots and Data Correlation! → A Scatter plot is

one of the most effective graphical methods for determining if there appears to be a relationship, pattern, or trend between two numeric attributes. To construct a scatter plot, each pair of values is treated as a pair of co-ordinates in an algebraic sense and plotted as points in the plane. Fig (a) shows a scatter plot for the set of data in table.

The scatter plot is a useful method for providing a first look at bivariate data to see clusters of points and outliers, or to explore the possibility of correlation relationship.

Two attributes, X , and Y are correlated if one attribute implies the other. Correlations can be +ve, -ve or null (uncorrelated).

Fig (b) shows examples of +ve and -ve correlations b/w two attributes. If the plotted points pattern slopes from lower left to upper right, this means that the value of X increases as the values of Y increase, suggesting a positive correlation b(i). If the pattern of plotted points slopes from upper left to lower right, the values of X increase as the value of Y decrease, suggesting a -ve correlation (Fig b(ii)).

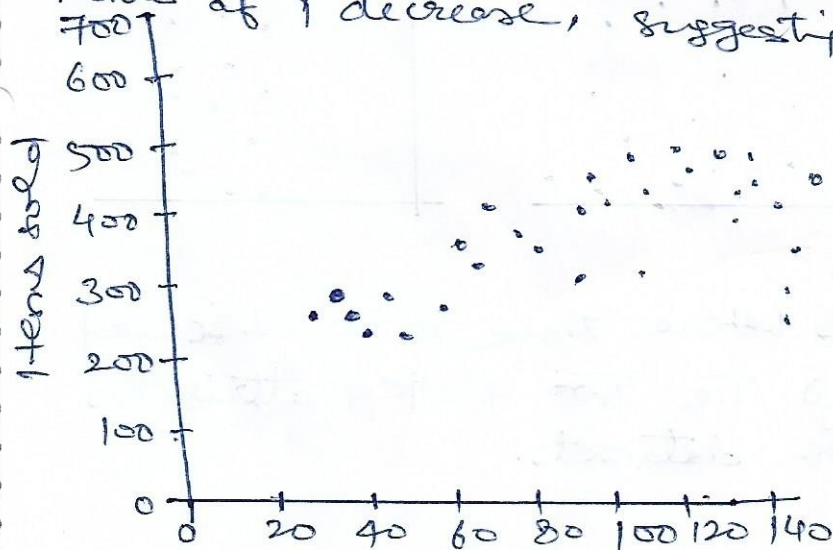
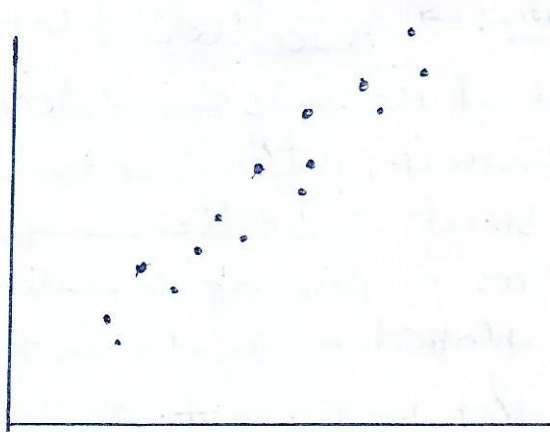
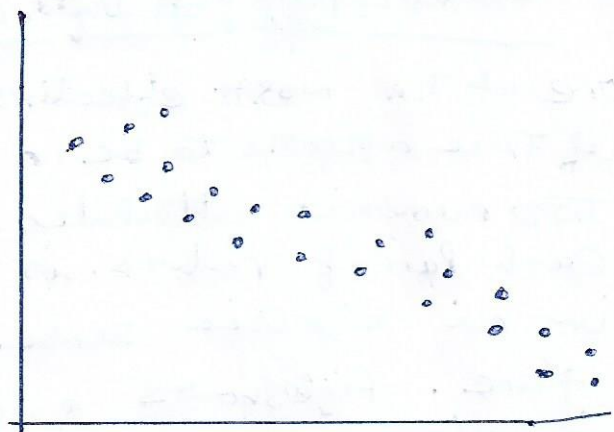


Fig (a) ! A histogram for the table.



(b) (i)



(b)(ii)

Fig:- scatter plots can be used to find (a) positive or
(b) Negative Correlations b/w attributes

A line of best fit can be drawn to study the correlation between the variables.

Fig (c) shows three cases for which there is no correlation relationship between the two attributes in each of the given data sets.

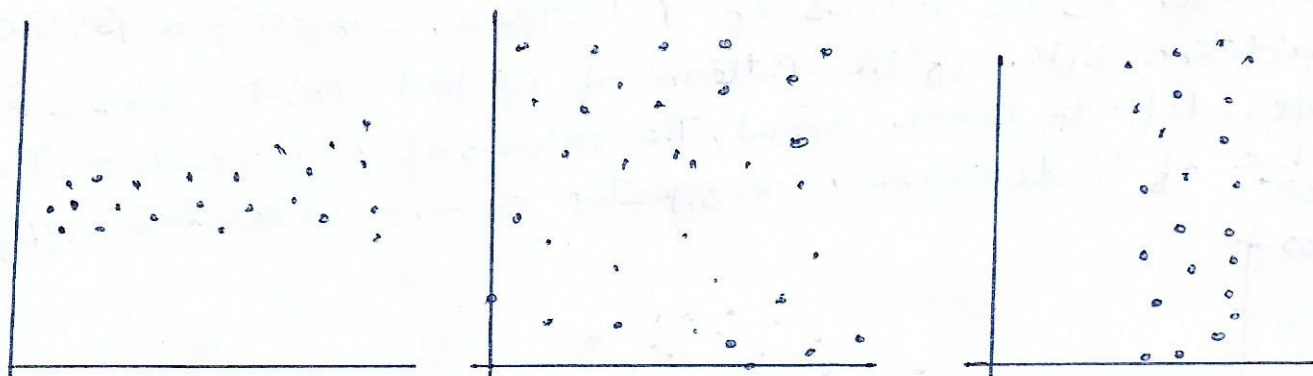


Fig:- Three cases where there is no observed correlation b/w the two plotted attributes in each of the datasets.