# Getting to know your data: →

## Date objects and Attribute types: →

Data sets are made up of data objects. A data object represents an entity — in a sales database, the objects may be customers, store items, and sales; in a Medical database, the objects may be patients;

in a university database, The objects may be Students, Professors, and Courses.

Data objects are typically described by attributes. Data objects can also be referred to as Samples, examples, instances, data points, or objects. If The data objects are stored in a databases, They are data tuples. That is, The row of a database correspond to the data objects, and The columns correspond to The attributes.

What is an attribute: — An attribute is a data field, representing a characteristic or feature of a data object. The noun attribute, dimension, feature, and Variable are often used interchangeably in The literature. the term dimension is commonly used in datawarehousing. machine learning literature tend to use the term feature, while statisticians prefer the term variable. Data mining and database Professionals commonly use the term attribute. Attributes describing a customer object / customer database can include, customer Id, name and address. The distribution of data involving one attribute (or variable) is called univariate. A bivariate distribution involves Two attributes, and so on.

The type of an attribute is determined by the set of possible values — nominal, binary, or numeric — the attribute can have.

## Nominal Attributes (=, ≠) :—

Nominal means "relating to names". The values a nominal attribute are symbols or names of things. Each value represents some kind of category, code, or state, and so nominal attributes are also referred to as categorical. The value do not have any meaningful order.

Ex :— Suppose that hair-color and marital-status are two attributes describing Person objects. In our application, possible values for hair-color are black, brown, blond, red, auburn, gray and white.

The attribute marital-status can take on the value single, married, divorced, and widowed. Both hair color and marital-status are nominal attributes.

Another example of a nominal attribute is occupation, with the values teacher, dentist, Programmer, farmer, and so on.

Although, we said that the values of a nominal attribute are symbols or "names of things". It is possible to represent such symbols or "names" with numbers. With hair color, for instances, we can assign a code of 0 for black, 1 for brown, and so on. But on these numbers we can not perform operation like numerical number (quantitative number).

Since, nominal attribute values do not have any meaningful order about them and are not quantitative, it makes no sense to find mean value or median value

## Binary Attributes:—

attribute with only two categories or states: 0 or 1 where 0 typically means that the attribute is absent, and 1 means that it is present. Binary attributes are referred to as Boolean if they the two states correspond to True and False.

### For EX:—

Given the attribute smoker describing a patient object, 1 indicates that the patient smokes, while 0 indicates that the patient does not. Similarly, suppose that Patient undergoes a medical test that has two possible outcomes. The attribute medical-test is binary, where a value of 1 means the results of the test for the Patient is positive, while 0 means the result of Patient is negative.

A binary attribute is <u>symmetric</u> if both of its states are equally valuable and carry the same weight; that is, there is no preference on which outcome should be coded as 0 or 1. One such example could be the attribute gender having the states. male and female. (Balanced data set)

A binary attribute is asymmetric if both of the states are not equally important, such as the positive and negative outcomes of a medical test for HIV. By convention, we code the most important outcome, which is usually the rarest one, by 1 (eg. HIV positive) and the other by 0 (eg. HIV negative)

## ordinal Attributes: $(<, >)$

An ordinal attribute is an attribute with possible values that have a meaningful order or ranking among them, but the magnitude b/w successive values is not known.

for Ex:- Suppose that drink-size corresponds to the size of drinks available at a fast-food restaurant. This nominal attribute has three possible values: small, medium, and large. The value have a meaningful sequence (which corresponds to increasing drink size); however, we cannot tell from the values how much bigger, say, a medium is than a large. Other examples of ordinal attributes include grade (eg. A+, A, A-, B+, and so on) and professional rank (assistant, associate, and Professors).

ordinal attributes are useful for registering subjective assessments of qualities that cannot be measured objectively; Thus ordinal attributes are often used in surveys for rating. In one survey, participants were asked to rate how satisfied they were as customers. customer satisfaction had the following ordinal categories:

0: very satisfied

1: somewhat dissatisfied

2: neutral

3: satisfied

4: very satisfied.

ordinal attributes may be obtained from the discretization of numeric quantities by splitting the value range into a finite number of ordered categories.

## Numeric Attributes :→

A numeric attribute is quantitative; that is, it is measurable quantity, represented in integer or real values. Numeric attributes can be interval-scaled or Ratio-scaled.

### Interval-scaled attributes :→

interval-scaled attributes are measured on a scale of equal-size units. The values of interval-scaled attributes have order and can be positive, 0, or negative. Thus, in addition to providing a ranking of values, such attributes allow us to compare and quantify the difference b/w values.

A temperature attribute is interval-scaled. Suppose that we have the outdoor temperature value for a number of different days, where each day is an object. By ordering the values, we obtain a ranking of the objects with respect to temperature. In addition we can quantify the difference b/w values. For Example, a temperature of 20°C is five degrees higher Than a temperature of 15°C. Calendar dates are another example.

For Instance, The years 2002 and 2010 are eight years apart.

* For interval attributes, The differences between values are Meaningful, i.e, a unit of measurement exists. (+, -)

for Ex:- Calendar dates, temperature in celsius or fahrenheit.

## Ratio-scaled Attributes :→

A ratio-scaled attributes is a numeric attribute with an inherent zero point. That is, if a measurement is ratio-scaled, we can speak of a value as being a multiple (or ratio) of another value. In addition, the values are ordered, and we can also compute the difference b/w values, as well as the mean, median and mode.

* for ratio variables, both differences and ratio are meaningful (*, /)

EX: temperature in kelvin, monetary quantities, counts, age, mass, length, electrical current;

## Discrete versus continuous Attributes :→

A discrete attribute has a finite or countably infinite set of values, which may or may not be represented as integers. The attributes hair-color, smoker, medical-test, and drink-size each have numeric values, such as 0 and 1 for binary attributes, or the values 0 to 110 for the attribute age.

An attribute is countably infinite if the set of possible values is infinite but the values can be put in a one-to-one correspondance with natural numbers. for Example, the attribute customer-id is countably infinite. The number of customers can grow infinity, but in reality, The actual set of values is countable (where the values can be put in one-to-one correspondence with the set of integers). zipcode are another example.

If an attribute is not discrete, it is continuous. The numeric attribute can be either integers or real number, but continuous value are only real numbers, where continuous attributes are typically represented as floating-point variables.