

What is Data?

Collection of data objects and their attributes.

Attribute is a Variable, field, characteristic, **feature**
property of an object

Collection of attributes describes an object
Record, point, case, sample, entity, instance, **observation**

Type of attributes

Discrete	Finite (countably)	Integer	<i>Zip, Counts</i>
Continuous	Real numbers	Floating points	<i>Temp., height, weight</i>

Hierarchy of attributes types

Qualitative	Nominal	Category (=, !=)	<i>ID, zip, eye, color</i>
	Ordinal	Ranked (>, <)	<i>Grades, {low, med., high}</i>
Quantitative	Interval	Distance (+, -)	<i>Dates, temp (C/F)</i>
	Ratio	Zero means absence (*, /)	<i>Length, time, temp(K)</i>

Type of data sets

Record	Collection of dataobjects and their attributes	Table
Relational	Collection of data objects and their relation	Graph
Ordered	Ordered collection of data objects	Sequence

Data quality

High quality	Are fit for their intended use Correctly represent the phenomena they correspond to
Problems	Noise Outliers Missing values

Noise

Definition	Unwanted perturbation to a signal Unwanted data
Reasons	Limits in measurement accuracy Interference from other signals Measurement of attributes not related to the data modeling task
Handling	Exclude noisy attributes Remove noise by filtering Include a model of noise

Outliers

Definition	Data objects which are significantly different from most others
Reasons	Measurement errors Natural property of data
Handling	Identify & exclude outliers Model the outliers

Missing values

Definition	No value is stored for an attribute in a data object
Reasons	Information is not collected <i>People decline to give their age</i> Attribute is not applicable <i>Annual income is not applicable to children</i>
Handling	Eliminate data objects Estimate missing values <i>e.g. average</i> Ignore the missing value in analysis

