# Unit II :
# Preprocessing and Extracting meaning from Data

- Identifying Missing values and approaches
- Noisy Data Extraction
- Data Cleaning as a process
- Data reduction
- Data Transformation and Discretization :
- Data Transformation by Normalization,

  Discretization by Binning

  Discretization by Histogram Analysis,

  Discretization by Cluster,
- Decision Tree, and Correlation and Regression analysis reasons to choose and cautions

# Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview

  - Data Quality

  - Major Tasks in Data Preprocessing

- Data Cleaning

- Data Reduction

- Data Transformation and Data Discretization

- Summary

# Data Quality: Why Preprocess the Data?

- Measures for data quality: A multidimensional view

    - Accuracy: correct or wrong, accurate or not

    - Completeness: not recorded, unavailable, …

    - Consistency: some modified but some not, dangling, …

    - Timeliness: timely update?

    - Believability: how trustable the data are correct?

    - Interpretability: how easily the data can be understood?

# Major Tasks in Data Preprocessing

- **Data cleaning**
  - Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies
- **Data reduction**
  - Dimensionality reduction
  - Numerosity reduction
  - Data compression
- **Data transformation and data discretization**
  - Normalization
  - Concept hierarchy generation

# Data Cleaning

- Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, transmission error
  - incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
    - e.g., *Occupation*=" " (missing data)
  - noisy: containing noise, errors, or outliers
    - e.g., *Salary*="−10" (an error)
  - inconsistent: containing discrepancies in codes or names, e.g.,
    - *Age*="42", *Birthday*="03/07/2010"
    - Was rating "1, 2, 3", now rating "A, B, C"
    - discrepancy between duplicate records
  - Intentional (e.g., *disguised missing* data)
    - Jan. 1 as everyone's birthday?

# Incomplete (Missing) Data

- Data is not always available
    - E.g., many tuples have no recorded value for several attributes, such as customer income in sales data
- Missing data may be due to
    - equipment malfunction
    - inconsistent with other recorded data and thus deleted
    - data not entered due to misunderstanding
    - certain data may not be considered important at the time of entry
    - not register history or changes of the data
- Missing data may need to be inferred

# How to Handle Missing Data?

- **Ignore the tuple**: usually done when class label is missing (when doing classification)—not effective when the % of missing values per attribute varies considerably

- **Fill in the missing value manually**: tedious + infeasible?

- **Fill in it automatically** with

  - a global constant : e.g., "unknown", a new class?!

  - the attribute mean

  - the attribute mean for all samples belonging to the same class: smarter

  - the most probable value: inference-based such as Bayesian formula or decision tree

# Noisy Data

- Noise: random error or variance in a measured variable
- Incorrect attribute values may be due to
    - faulty data collection instruments
    - data entry problems
    - data transmission problems
    - technology limitation
    - inconsistency in naming convention
- Other data problems which require data cleaning
    - duplicate records
    - incomplete data
    - inconsistent data

# How to Handle Noisy Data?

- Binning
  - first sort data and partition into (equal-frequency) bins
  - then one can smooth by bin means, smooth by bin median, smooth by bin boundaries, etc.
- Regression
  - smooth by fitting the data into regression functions
- Clustering
  - detect and remove outliers
- Combined computer and human inspection
  - detect suspicious values and check by human (e.g., deal with possible outliers)

# Data Cleaning as a Process

- First step in Data cleaning is Data discrepancy detection
    - Use metadata (e.g., domain, range, dependency, distribution): Data about data
    - Check field overloading
        - Check uniqueness rule – each value should be different than other value
        - consecutive rule – there can be no missing values between lowest and highest value of attr.
        - null rule – use of blanks, question marks etc.
    - Use commercial tools
        - Data scrubbing: use simple domain knowledge (e.g., postal code, spell-check) to detect errors and make corrections
        - Data auditing: by analyzing data to discover rules and relationship to detect violators (e.g., correlation and clustering to find outliers)

# Data Reduction Strategies

- **Data reduction**:
    - Obtain a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results
- Why data reduction? — A database/data warehouse may store terabytes of data. Complex data analysis may take a very long time to run on the complete data set.
- Data reduction strategies
    - Dimensionality reduction, e.g., remove unimportant attributes
        - Wavelet transforms
        - Principal Components Analysis (PCA)
        - Feature subset selection, feature creation
    - Numerosity reduction (some simply call it: Data Reduction)
        - Regression and Log-Linear Models
        - Histograms, clustering, sampling
        - Data cube aggregation
    - Data compression : Lossless

# Data Reduction 1: Dimensionality Reduction

- **Curse of dimensionality**
  - When dimensionality increases, data becomes increasingly sparse
  - Density and distance between points, which is critical to clustering, outlier analysis, becomes less meaningful
  - The possible combinations of subspaces will grow exponentially
- **Dimensionality reduction**
  - Avoid the curse of dimensionality
  - Help eliminate irrelevant features and reduce noise
  - Reduce time and space required in data mining
  - Allow easier visualization
- **Dimensionality reduction techniques**
  - Wavelet transforms
  - **Principal Component Analysis**
  - Supervised and nonlinear techniques (e.g., feature selection)

# Attribute Subset Selection

- **Another way to reduce dimensionality of data**
- Redundant attributes
  - Duplicate much or all of the information contained in one or more other attributes
  - E.g., purchase price of a product and the amount of sales tax paid
- Irrelevant attributes
  - Contain no information that is useful for the data mining task at hand
  - E.g., students' ID is often irrelevant to the task of predicting students' GPA

# Heuristic Search in Attribute Selection

- There are $2^d$ possible attribute combinations of $d$ attributes
- Typical heuristic attribute selection methods:
  - Best single attribute under the attribute independence assumption: choose by significance tests
  - Best step-wise feature selection: Forward approach
    - The best single-attribute is picked first
    - Then next best attribute condition to the first, …
  - Step-wise attribute elimination: Backward approach
    - Repeatedly eliminate the worst attribute
  - Best combined attribute selection and elimination
  - Optimal branch and bound:
    - Use attribute elimination and backtracking

# Attribute Creation (Feature Generation)

- Create new attributes (features) that can capture the important information in a data set more effectively than the original ones
- Three general methodologies
  - Attribute extraction
    - Domain-specific
  - Mapping data to new space (see: data reduction)
    - E.g., Fourier transformation, wavelet transformation, manifold approaches (not covered)
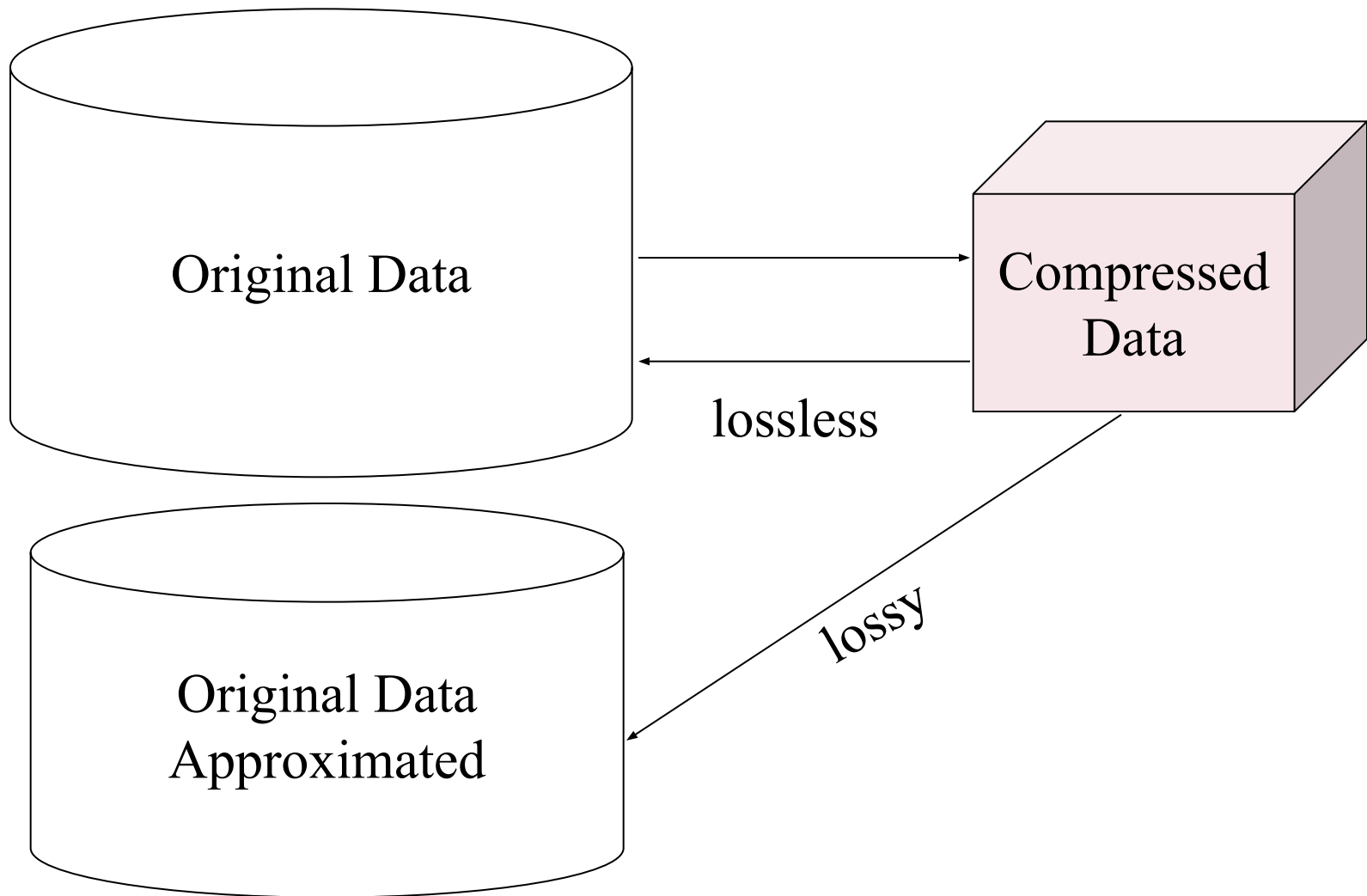  - Attribute construction
    - Data discretization

# Data Reduction 3: Data Compression

- String compression
  - There are extensive theories and well-tuned algorithms
  - Typically lossless, but only limited manipulation is possible without expansion
- Audio/video compression
  - Typically lossy compression, with progressive refinement
  - Sometimes small fragments of signal can be reconstructed without reconstructing the whole
- Time sequence is not audio
  - Typically short and vary slowly with time
- Dimensionality and numerosity reduction may also be considered as forms of data compression

# Data Compression



Original Data → Compressed Data

Compressed Data → Original Data (lossless)

Compressed Data → Original Data Approximated (lossy)

# Data Transformation

- A function that maps the entire set of values of a given attribute to a new set of replacement values s.t. each old value can be identified with one of the new values

- Methods

  - Smoothing: Remove noise from data [Binning,regression,clustering]

  - Attribute/feature construction

    - New attributes constructed from the given ones

  - Aggregation: Summarization, data cube construction[OLAP]

  - Normalization: Scaled to fall within a smaller, specified range

    - min-max normalization

    - z-score normalization

    - normalization by decimal scaling

  - Discretization: Concept hierarchy climbing

# Normalization

- **Min-max normalization**: to [new_min$_A$, new_max$_A$]

$$v' = \frac{v - min_A}{max_A - min_A}(new\_max_A - new\_min_A) + new\_min_A$$

- Ex. Let income range \$12,000 to \$98,000 normalized to [0.0, 1.0]. Then \$73,600 is mapped to $\frac{73,600 - 12,000}{98,000 - 12,000}(1.0 - 0) + 0 = 0.716$

- **Z-score normalization** (μ: mean, σ: standard deviation):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

- **Normalization by decimal scaling**

$$v' = \frac{v}{10^j}$$ Where $j$ is the smallest integer such that Max($|v'|$) < 1

# Discretization

- Three types of attributes
    - Nominal—values from an unordered set, e.g., color, profession
    - Ordinal—values from an ordered set, e.g., military or academic rank
    - Numeric—real numbers, e.g., integer or real numbers
- Discretization: Divide the range of a continuous attribute into intervals
    - Interval labels can then be used to replace actual data values
    - Reduce data size by discretization
    - Supervised vs. unsupervised
    - Split (top-down) vs. merge (bottom-up)
    - Discretization can be performed recursively on an attribute
    - Prepare for further analysis, e.g., classification

# Data Discretization Methods

- Typical methods: All the methods can be applied recursively
  - Binning
    - Top-down split, unsupervised
  - Histogram analysis
    - Top-down split, unsupervised
  - Clustering analysis (unsupervised, top-down split or bottom-up merge)
  - Decision-tree analysis (supervised, top-down split)
  - Correlation (e.g., $\chi^2$) analysis (unsupervised, bottom-up merge)
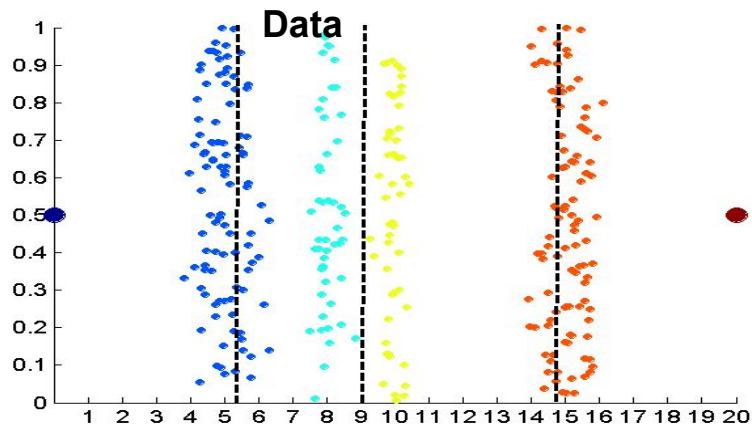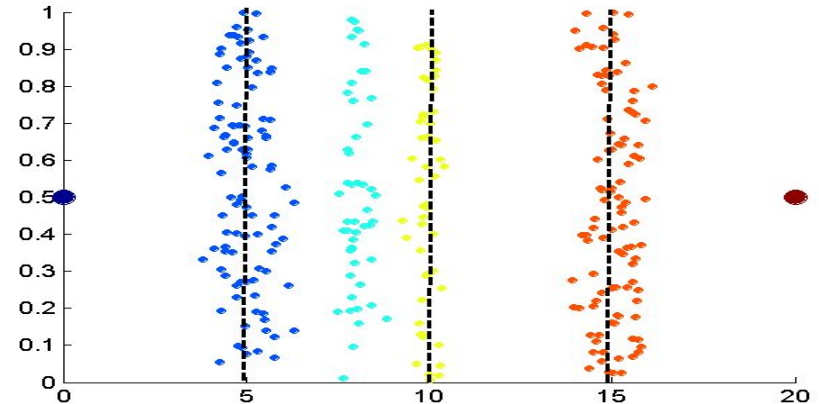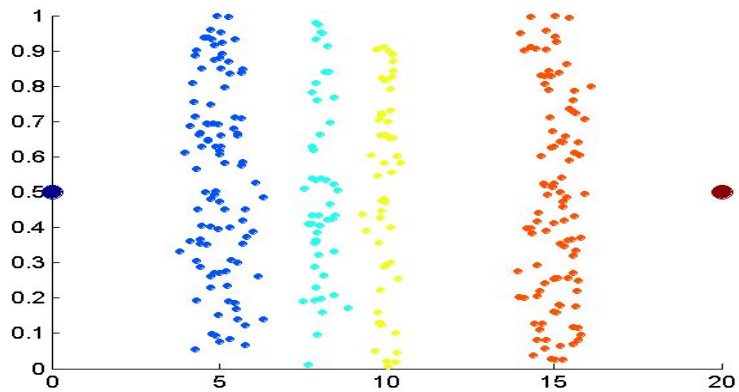
# Simple Discretization: Binning

- Equal-width (distance) partitioning
    - Divides the range into *N* intervals of equal size: uniform grid
    - if *A* and *B* are the lowest and highest values of the attribute, the width of intervals will be: $W = (B - A)/N.$
    - The most straightforward, but outliers may dominate presentation
    - Skewed data is not handled well

- Equal-depth (frequency) partitioning
    - Divides the range into *N* intervals, each containing approximately same number of samples
    - Good data scaling
    - Managing categorical attributes can be tricky
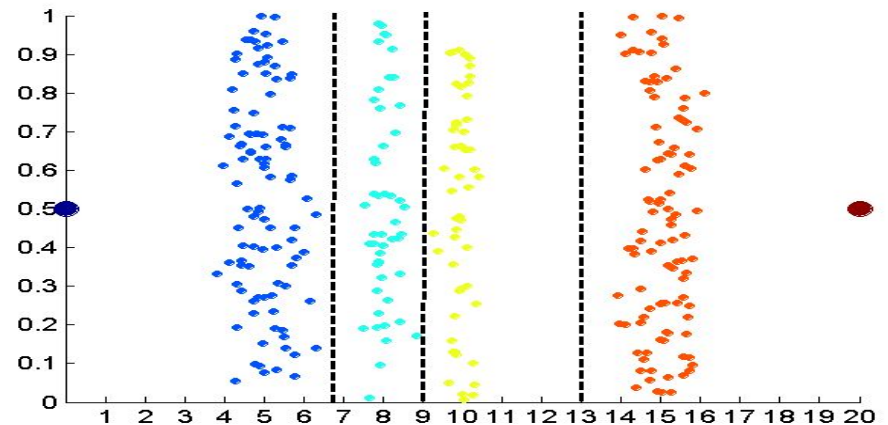
# Binning Methods for Data Smoothing

❑ Sorted data for price (in dollars): 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34

\* Partition into equal-frequency (**equi-depth**) bins:

    - Bin 1: 4, 8, 9, 15

    - Bin 2: 21, 21, 24, 25

    - Bin 3: 26, 28, 29, 34

\* Smoothing by **bin means**:

    - Bin 1: 9, 9, 9, 9

    - Bin 2: 23, 23, 23, 23

    - Bin 3: 29, 29, 29, 29

\* Smoothing by **bin boundaries**:

    - Bin 1: 4, 4, 4, 15

    - Bin 2: 21, 21, 25, 25

    - Bin 3: 26, 26, 26, 34

# Discretization Without Using Class Labels
## (Binning vs. Clustering)



**Data**

**Equal frequency (binning)**

**K-means clustering leads to better results**

# Discretization by Classification & Correlation Analysis

- Classification (e.g., decision tree analysis)

  - Supervised: Given class labels, e.g., Patients Vs Symptoms

  Class distribution information is used in the calculation and determination of split-points (data values for partitioning an attribute range). Intuitively, the main idea is to select split-points so that a given

  resulting partition contains as many tuples of the same class as possible.

  - Using *entropy* to determine split point (discretization point)

  - Top-down, recursive split

- Correlation analysis Measures of correlation can be used for discretization (e.g., Chi-merge: $\chi^2$-based discretization)

  - Supervised: use class information

  - Bottom-up merge: find the best neighboring intervals (those having similar distributions of classes, i.e., low $\chi^2$ values) to merge

  - Merge performed recursively, until a predefined stopping condition

- ChiMerge proceeds as follows. Initially, each distinct value of a numeric attribute $A$ is considered to be one interval. $X^2$ tests are performed for every pair of adjacent intervals. Adjacent intervals with the least $X^2$ values are merged together, because low values for a pair indicate similar class distributions. This merging process proceeds recursively until a predefined stopping criterion is met.

# Concept Hierarchy Generation

- **Concept hierarchy** organizes concepts (i.e., attribute values) hierarchically and is usually associated with each dimension in a data warehouse

- Concept hierarchies facilitate <u>drilling and rolling</u> in data warehouses to view data in multiple granularity

- Concept hierarchy formation: Recursively reduce the data by collecting and replacing low level concepts (such as numeric values for *age*) by higher level concepts (such as *youth, adult*, or *senior*)

- Concept hierarchies can be explicitly specified by domain experts and/or data warehouse designers

- Concept hierarchy can be automatically formed for both numeric and nominal data.  For numeric data, use discretization methods shown.
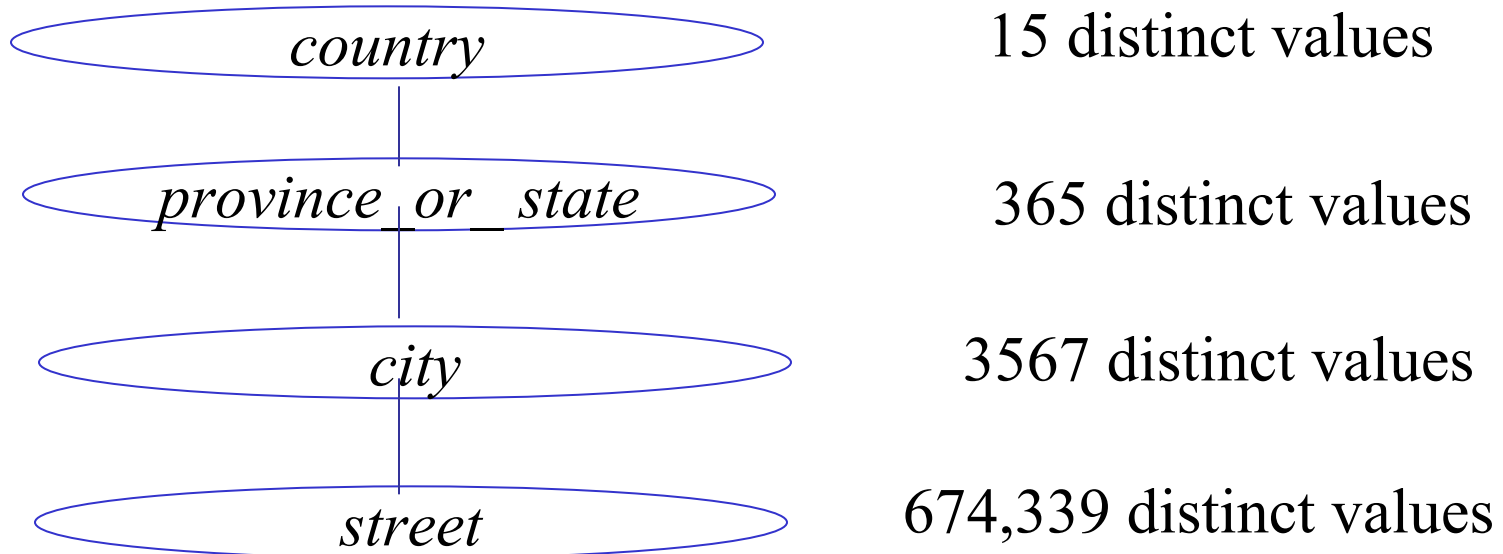
# Concept Hierarchy Generation for Nominal Data

- Nominal attributes have a finite (but possibly large) number of distinct values, with no ordering among the values. Examples include *geographic location*, *job category*, and *item type*.

- Specification of a partial/total ordering Concept hierarchies for nominal attributes or dimensions typically involve a group of attributes. A user or expert can easily define a concept hierarchy by specifying a partial or total ordering of the attributes at the schema level.

  - *street < city < state < country*

- Specification of a hierarchy for a set of values we can easily specify explicit groupings for a small portion of intermediate-level data {Mumbai,Pune ,Nagpur} < Maharashtra

- <span style="color:red">Specification of only a partial set of attributes</span>
  - E.g., only *street < city*, not others
- Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values
  - E.g., for a set of attributes: {*street, city, state, country*}

# Automatic Concept Hierarchy Generation

- Some hierarchies can be automatically generated based on the analysis of the number of distinct values per attribute in the data set
  - The attribute with the most distinct values is placed at the lowest level of the hierarchy
  - Exceptions, e.g., weekday, month, quarter, year

| | |
|---|---|
| *country* | 15 distinct values |
| *province_or_state* | 365 distinct values |
| *city* | 3567 distinct values |
| *street* | 674,339 distinct values |

# Summary

- **Data quality**: accuracy, completeness, consistency, timeliness, believability, interpretability
- **Data cleaning**: e.g. missing/noisy values, outliers
- **Data reduction**
    - Dimensionality reduction
    - Numerosity reduction
    - Data compression
- **Data transformation and data discretization**
    - Normalization
    - Concept hierarchy generation

- https://www.youtube.com/watch?v=P_iMSYQnqac