# Major Tasks in Data Preprocessing

- **Data cleaning**
  - Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies
- **Data integration**
  - Integration of multiple databases, data cubes, or files
- **Data Transformation**
  - Normalization
- **Data reduction (Coming soon)**
  - Dimensionality reduction
  - Numerosity reduction
  - Data compression
  - Feature Selection
- **Data discretization (Coming soon)**
  - Concept hierarchy generation

# Data Cleaning

# Data Cleaning

- Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, transmission error
  - <u>incomplete</u>: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
    - e.g., *Occupation*=" " (missing data)
  - <u>noisy</u>: containing noise, errors, or outliers
    - e.g., *Salary*="−10" (an error)
  - <u>inconsistent</u>: 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
  - <u>Intentional</u> (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
- Missing data may need to be inferred

# Missing Value Imputation

- 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

https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration

# 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
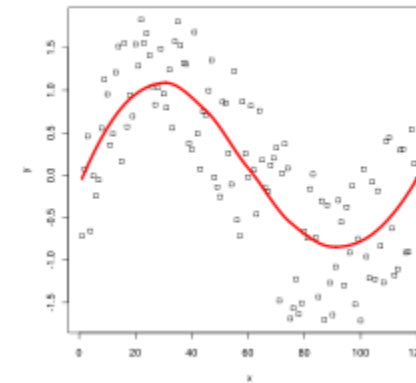
# How to Handle Noisy Data?

- **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)

Sorted data for Age: 3, 7, 8, 13,    22, 22, 22, 26,    26, 28, 30, 37

| equal frequency bins | bin means | bin boundaries' |
|---|---|---|
| Bin 1: 3, 7, 8, 13 | Bin 1: 8, 8, 8, 8 | Bin 1: 3, 3, 3, 13 |
| Bin 2: 22, 22, 22, 26 | Bin 2: 23, 23, 23, 23 | Bin 2: 22, 22, 22, 26 |
| Bin 3: 26, 28, 30, 37 | Bin 3: 30, 30, 30, 30 | Bin 3: 26, 26, 26, 37 |

https://T4Tutorials.com

Outliers

9

# Data Integration

# Data Integration

- **Data integration**:
    - Combines data from multiple sources into a coherent store
- Schema integration: e.g., A.cust-id $\equiv$ B.cust-#
    - Integrate metadata from different sources
- <span style="color:red">Entity identification problem</span>:
    - Identify real world entities from multiple data sources, e.g., Bill Clinton = William Clinton
- Detecting and resolving data value conflicts
    - For the same real world entity, attribute values from different sources are different
    - Possible reasons: different representations, different scales, e.g., metric vs. British units

# Handling Redundancy in Data Integration

- Redundant data occur often when integration of multiple databases

  - *Object identification*:  The same attribute or object may have different names in different databases

  - *Derivable data:* One attribute may be a "derived" attribute in another table, e.g., annual revenue

- Redundant attributes may be able to be detected by *correlation analysis* and *covariance analysis*

- Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

# Correlation Analysis (Nominal Data)

- **X² (chi-square) test**

$$\chi^2 = \sum \frac{(Observed - Expected)^2}{Expected}$$

- The larger the X² value, the more likely the variables are not related

- The cells that contribute the most to the X² value are those whose actual count is very different from the expected count

- Correlation does not imply causality
  - # of hospitals and # of car-theft in a city are correlated
  - Both are causally linked to the third variable: population

# Chi-Square Calculation: An Example

| | Play chess | Not play chess | Sum (row) |
|---|---|---|---|
| Like science fiction | 250(90) | 200(360) | 450 |
| Not like science fiction | 50(210) | 1000(840) | 1050 |
| Sum(col.) | 300 | 1200 | 1500 |

- X$^2$ (chi-square) calculation (numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

- It shows that like_science_fiction and play_chess are not correlated in the group
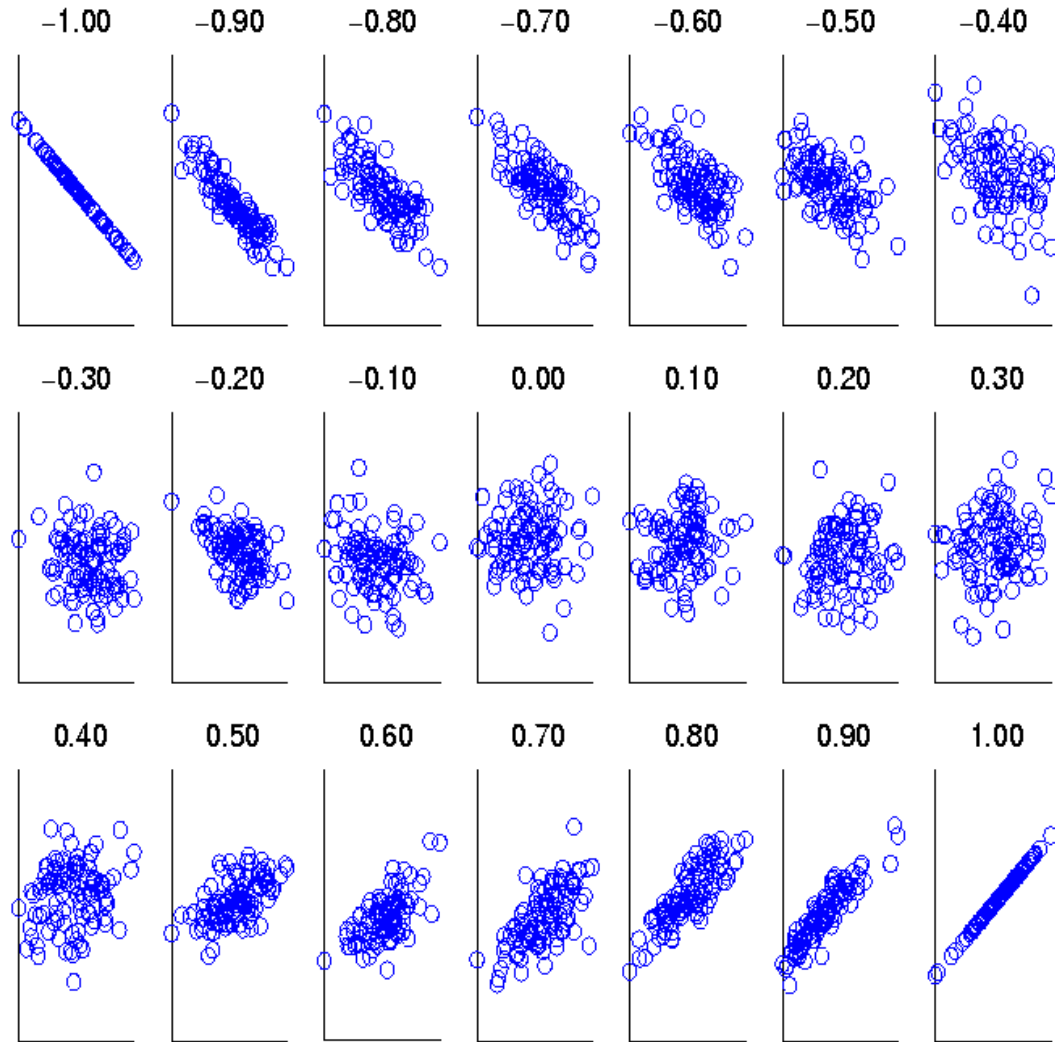
# Correlation Analysis (Numeric Data)

- Correlation coefficient (also called Pearson's product moment coefficient)

$$r_{A,B} = \frac{\sum_{i=1}^{n}(a_i - \overline{A})(b_i - \overline{B})}{(n-1)\sigma_A \sigma_B} = \frac{\sum_{i=1}^{n}(a_i b_i) - n\overline{A}\overline{B}}{(n-1)\sigma_A \sigma_B}$$

where n is the number of tuples, $\overline{A}$ and $\overline{B}$ are the respective means of A and B, $\sigma_A$ and $\sigma_B$ are the respective standard deviation of A and B, and $\Sigma(a_i b_i)$ is the sum of the AB cross-product.

- If $r_{A,B} > 0$, A and B are positively correlated (A's values increase as B's). The higher, the stronger correlation.

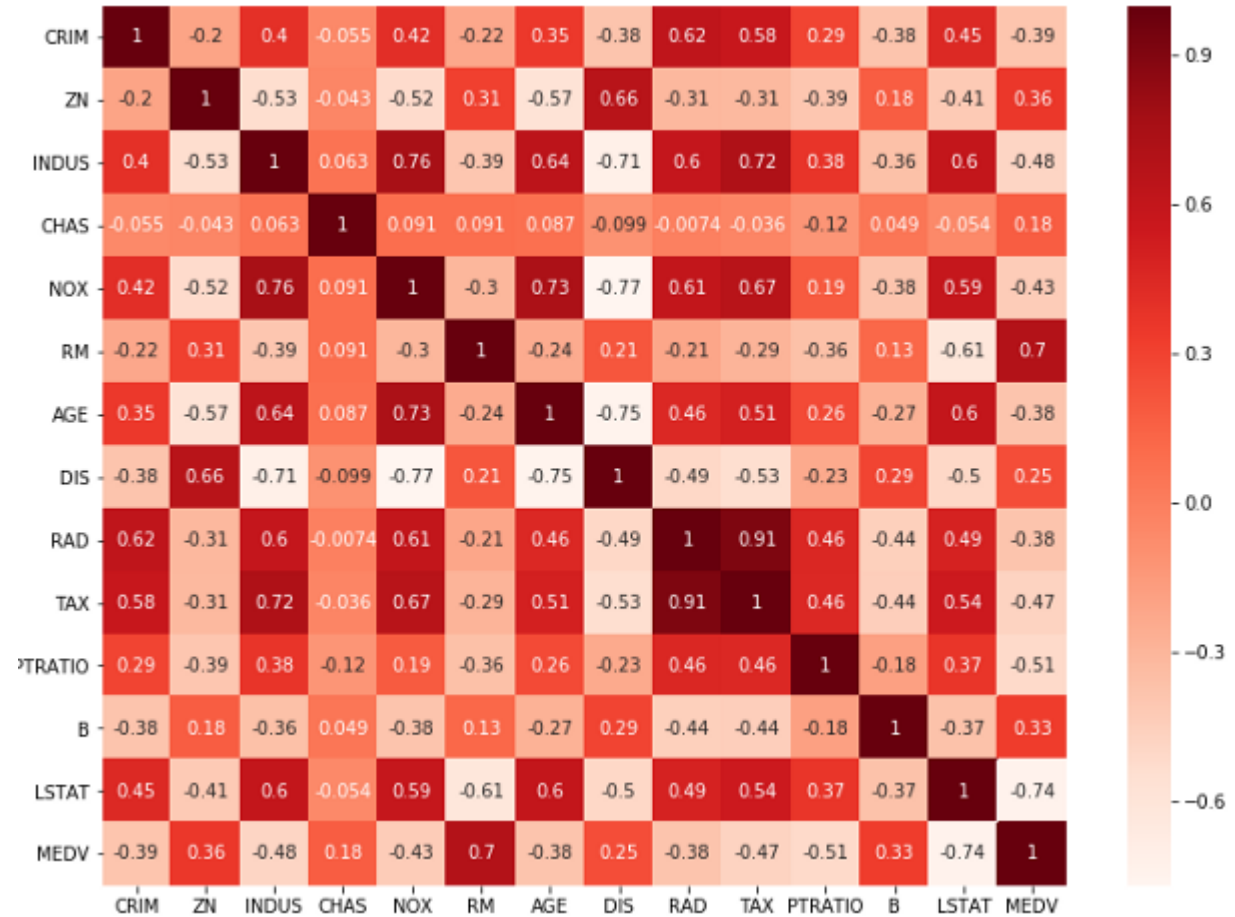- $r_{A,B} = 0$: independent;  $r_{AB} < 0$: negatively correlated

# Visually Evaluating Correlation



Scatter plots showing the similarity from –1 to 1.

# Using Pandas for Correlation

```
#Using Pearson Correlation
plt.figure(figsize=(12,10))
cor = df.corr()
sns.heatmap(cor, annot=True,
cmap=plt.cm.Reds)
plt.show()
```



As we can see, only the features RM, PTRATIO and LSTAT are highly correlated with the output variable MEDV. Hence we will drop all other features apart from these.

# Correlation Approaches

| Feature\Response | Continuous | Categorical |
|---|---|---|
| Continuous | Pearson's Correlation | LDA |
| Categorical | Anova | Chi-Square |

# Data Transformation


What a transformation

# 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

  - Attribute/feature construction

    - New attributes constructed from the given ones

  - 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 \$73,000 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):

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

  - Ex. Let μ = 54,000, σ = 16,000. Then $\frac{73,600 - 54,000}{16,000} = 1.225$

- **Normalization by decimal scaling**

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

# Miscellaneous

# Binarizer

```
In [1]:  import numpy as np
         from sklearn import preprocessing
```

```
In [2]:  data = np.array([[5.1, -2.9, 3.3],
                          [-1.2, 7.8, -6.1],
                          [3.9, 0.4, 2.1],
                          [7.3, -9.9, -4.5]])
```

```
In [6]:  binarizedData = preprocessing.Binarizer(threshold=2.1).transform(data)
         print("Binarized Data",binarizedData)

         Binarized Data [[1. 0. 1.]
          [0. 1. 0.]
          [1. 0. 0.]
          [1. 0. 0.]]
```

# One Hot Encoding

## Label Encoding

| Food Name | Categorical # | Calories |
|-----------|---------------|----------|
| Apple | 1 | 95 |
| Chicken | 2 | 231 |
| Broccoli | 3 | 50 |

$\rightarrow$

## One Hot Encoding

| Apple | Chicken | Broccoli | Calories |
|-------|---------|----------|----------|
| 1 | 0 | 0 | 95 |
| 0 | 1 | 0 | 231 |
| 0 | 0 | 1 | 50 |

# Other Types of Data

# Image, Video, Audio and Text

- Image/Video Features

  - Pixels, Corners, Edges, Keypoints, Color, Segments, CNN features

- Audio Features

  - Volume, power, air pressure, Frequency, RNN Features

- Text Features

  - Words, POS tags, Linguistic Features, TFIDF, word embedding

- How the previously mentioned preprocessing approaches can be applied for image, video, audio and text data