Data Mining

# Tour 2
# Data Preprocessing

# Data Preprocessing

## Terminology

| | |
|---|---|
| **Data Preprocessing** | • وهى العمليات اللى بتم قبل ما ابدأ فى التحليل لان البيانات بتكون غير مهيئة للتحليل بسبب المشاكل اللى فى البيانات<br>• Datasets are highly susceptible to noisy, missing, and inconsistent data.<br>• Low-quality data will lead to low-quality mining results.<br>• Data quality factors includes:<br>accuracy,<br>completeness,<br>consistency,<br>timeliness,<br>believability<br>interpretability. |
| **Reasons of Low Data Quality** | • Inaccurate data<br>   • Having incorrect attribute values (e.g., by choosing the default value "January 1" displayed for birthday)<br>   • عدى يا عمى التاريخ بسرعة اختاره اى حاجة ,, المفروض يبقى فيه اسلوب تحقق زى ويب فاليديشن<br>• **Incomplete data**<br>   • **Missing data (may not always be available or of interest)**<br>• **Inconsistent data**<br>   • **different assessments of the quality depend on the intended use of the data**<br>• Timeliness<br>   • (e.g. month-end data are not updated in a timely fashion has a negative impact on the data quality. )<br>• **Believability**<br>   • **reflects how much the data are trusted by users**<br>• Interpretability<br>   • reflects how easy the data are **understood** (e.g. sales codes)<br><br><table><tr><td>Tid</td><td>Refund</td><td>Marital Status</td><td>Taxable Income</td><td>Cheat</td></tr><tr><td>1</td><td>Yes</td><td>Single</td><td>125K</td><td>No</td></tr><tr><td>2</td><td>No</td><td>Married</td><td>100K</td><td>No</td></tr><tr><td>3</td><td>No</td><td>Single</td><td>70K</td><td>No</td></tr><tr><td>4</td><td>Yes</td><td>Married</td><td>120K</td><td>No</td></tr><tr><td>5</td><td>No</td><td>Divorced</td><td>10000K</td><td>Yes</td></tr><tr><td>6</td><td>No</td><td>NULL</td><td>60K</td><td>No</td></tr><tr><td>7</td><td>Yes</td><td>Divorced</td><td>220K</td><td>NULL</td></tr><tr><td>8</td><td>No</td><td>Single</td><td>85K</td><td>Yes</td></tr><tr><td>9</td><td>No</td><td>Married</td><td>90K</td><td>No</td></tr><tr><td>9</td><td>No</td><td>Single</td><td>90K</td><td>No</td></tr></table> |
| **Preprocessing Tasks That Improve Data Quality** | **Data cleaning**: missing values, noisy data, outliers<br><br>**Data integration**: data from multiple sources<br><br>**Data reduction**: reduced representation of the data set<br><br>**Data transformation:** data scaled to fall within a smaller range like 0.0 to 1.0 |
| **Data Cleaning** | is about:<br>• filling in **missing values**<br>   ○ Missing Values: {Nan, Null, Na, ""}<br>• Smooth out **noise**<br>   ○ Noisy: Contain **Errors** => Salary = -1000<br>• Identifying or removing **outliers**<br>• Removing **inconsistencies**<br>   ○ (e.g. rating was "1, 2, 3", now rating "A, B, C") |

| | |
|---|---|
| | • **Intentional** manipulating (e.g., by choosing the default value "January 1" displayed for birthday) |
| **Missing Values Methods** | • **Ignore** the Tuple<br>   o **Effective when Class label is missing, and the Required task is Classification**<br>   o Not Effective unless the tuple has varied of missing attributes<br>   o **Ignoring make no use of other attributes of the tuple that can be useful**<br>• **Filling** in missing values manually<br>   o Time Consuming<br>   o May not be feasible in a large data set<br>• **Using** a **global constant** to fill the missing value<br>   o **Simple**<br>   o **Not foolproof in datamining tasks** => مش مضمون لانه منكم يأثر على عمليات الماينينج<br>• Use a measure of **central tendency** for the attribute to fill in the missing value<br>   o **Symmetric** data distribution => Use **Mean**<br>   o **Skewed** data distribution => Use **Median**<br><br>• **Use** the attribute **mean or median** for all samples belonging to the same class as the given tuple<br>   o زى اللى فوقها بالظبط بس على مستوى الداتا اللى من نفس الكلاس<br>• **Use** the **most probable value** to fill in the missing value<br>   o Can be determined using **Regression, Decision Tree Induction, Inference Based tools**<br>   o Uses the most information of present data to predict the missing value<br><br>A missing value may not imply an error in the data!<br>⇨ Forms should allow respondents to specify values such as "not applicable." Software routines may also be used to uncover other null values (e.g., "don't know," "?" or "none") |
| **Noisy Data** | • Noise is a **random error** or **variance** in a measured variable.<br>• "**smooth**" out the data to **remove** the noise.<br>• Smooth Data Techniques:<br>   • Binning<br>   • Regression<br>   • Outlier Analysis |
| **Binning** | o The original data values are divided into "**buckets**" known as "**bins**" and then they are replaced by a **general value** calculated for that bin.<br>o In the context of **image processing**, binning is the procedure of combining a **cluster of pixels into a single pixel**.<br>o 2 Steps for Binning<br>   o Partitioning<br>   o Smoothing<br><br>o "Partition" sorted data by 2 methods:<br>   o **Equal depth (frequency) bins**: each bin has same number of values<br>   o **Equal width bins**: interval range of values per bin is equal<br><br>o "Smooth" each bin by:<br>   o **bin means**: each bin value is replaced by the bin **mean**<br>   o **bin medians**: each bin value is replaced by the bin **median**<br>   o **bin boundaries**: each bin value is replaced by the **closest boundary value** (min & max in a bin are bin boundaries) |

# Example

- Use smoothing to smooth the Age data, Partition them into three bins by
  - equal-frequency and
  - equal-width Partitioning
- Age: 23, 23, 27, 27, 39, 41, 47, 49, 50, 52, 54, 54, 56, 56, 57, 58, 60, 61

**Partition**
- Age data are first sorted,
- 23, 23, 27, 27, 39, 41, 47, 49, 50, 52, 54, 54, 56, 56, 57, 58, 60, 61
- Partitioned into three equal frequency bins of size 6
- Or Partitioned into 3-equal interval bins

Partition into (equal-frequency bins:

Bin 1: 23, 23, 27, 27, 39, 41

Bin 2: 47, 49, 50, 52, 54, 54

Bin 3: 56, 56, 57, 58, 60, 61

Partition into (equal-width) bins:

Bin 1: 23, 23, 27, 27

Bin 2: 39, 41, 47, 49, 50
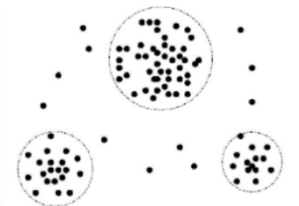
Bin 3: 52, 54, 54, 56, 56, 57, 58, 60, 61

| | |
|---|---|
| **Regression** | o <u>Linear</u> regression involves finding the "best" line to **fit two attributes** (or variables) so that **one attribute can be used to predict the other.**<br>o Replace **noisy or missing values** by **predicted** values. |
| **Outlier Analysis** | o May be Detected by **Clustering**<br>o The data outside the cluster {circle} may be analyzed as Outlier<br>o باذن الله مشروحة في جزء ال<br>o Clustering |
| **Data Integration** | o Merging of data from multiple **data stores**.<br>o Be Careful when integration because of Redundancy and inconsistencies<br>o فيه مشاكل هتقابلك و انت بتعمل تجميع للبيانات،،، تعالى نبدأ نتعرف عليهم!<br>o **Problems:**<br>o **Entity Identification Problem**<br>   o How can equivalent real-world entities from multiple data sources be matched up?<br>   o For example, how can the data analyst or the computer be sure that customer id in one database and cust_number in another refer to the same attribute?<br>   o **metadata** can be used to help avoid errors in schema integration.<br>   o You Can See the Metadata about Iris dataset<br>o **Redundancy**<br>   o An attribute may be redundant if it can be "*derived*" from another attribute or set of attributes. => Age, Date of Birth, annual revenue, for instance<br>   o So some Redundancies Can be detected using **Correlation analysis**<br>   o **Correlation analysis**, given two attributes, such analysis can measure how strongly one attribute implies the other.<br>   o There are 2 test:<br>      ▪ Chi-Squre for Nominal Data<br>      ▪ Covariance for Numeric<br>o تعالى نفتح الموضوع ده في سكشن جديد 🙂 |
| **Correlation Test for Nominal Data** | o a correlation relationship between two attributes, A and B, can be discovered by a $X^2$ (chi-square) test<br><br>$$\chi^2 = \sum_{i=1}^{c}\sum_{j=1}^{r}\frac{(o_{ij}-e_{ij})^2}{e_{ij}},$$ |

Datasets Used:

| Dataset | Attributes MetaInfo |
|---|---|
| Iris | 1. sepal length in cm<br>2. sepal width in cm<br>3. petal length in cm<br>4. petal width in cm<br>5. class:<br>-- Iris Setosa<br>-- Iris Versicolour<br>-- Iris Virginica |

- Where $o_{ij}$ is the observed frequency (i.e., actual count) of the joint event $(A_i, B_j)$ and $e_{ij}$ is the expected frequency of $(A_i, B_j)$ which can be computed as

$$e_{ij} = \frac{count(A = a_i) \times count(B = b_j)}{n},$$

- tests the hypothesis that A and B are independent, that is, there is no correlation between them.[Null vs Research Hypothesis]
- وده اللى كان موضحه دكتور صيام في مادة سنة ثالثة مادة بحث علمى
- دلوقتى تعالى نحل المثال ده لاحظ ان الألفا بيكون مرجعى يعنى مش بيتحسب بس احنا بنستخدمه على حسب المجال
- ولاحظ ان ال كاى سكوير بتثبت ان لاوجود ترابط و معنى كده ان هي بتثبت ال
- Null Hypothesis

Example :- Determine if is a relationship Between Gender and getting in Trouble

| | Get in Trouble | not — | Total |
|---|---|---|---|
| Boys | 46 | 71 | 117 |
| girls | 37 | 83 | 120 |
| Total | 83 | 154 | 237 |

$$C_{11} = \frac{117 \times 83}{237} = 40.97$$

$$C_{12} = \frac{117 \times 154}{237} = 76.02$$

$$C_{21} = \frac{120 \times 83}{237} = 42.03$$

$$C_{22} = \frac{120 \times 154}{237} = 77.97$$

$$X^2 = \frac{(46 - 40.97)^2}{40.97} + \frac{(71 - 76.02)^2}{76.02} + \frac{(37 - 42.03)^2}{42.03}$$

$$+ \frac{(83 - 77.97)^2}{77.97} = 1.87$$

(degree of freedom) ⇐ Critical Statistics معتاحينمحالا نحن على عمود

↪ DF = ( # rows -1) ( # col -1) = (2-1) ( 2-1) = 1

reference

منوت الجدول الى تحتت بحدد مابكن DF . 1 و ال (Alpha) يعنى:0.05

حساب Critical ⇐ لو قتح مقارنت مابين Critical J عامل
وال حساب 3.941 > 1.87
لو تكبر بيتقبل ⇐ Null H مالها و بالتالي مافيض علاقة مابيهم

- Correlation doesn't imply causality
- يعنى في المثال حتى لو فيه ترابط بين الشغب و النوع ده مش بيثبت ان النوع هو اللى بيسبب الشغب

<table>
<tr><td valign="top">

**Correlation Coefficient for Numeric Data**

</td><td>

- The correlation coefficient between two attributes, A and B, is
  - If $r_{A,B}$ is *greater* than 0, then *A* and *B* are *positively* correlated, The higher the value, the stronger the correlation
  - If $r_{A,B}$ = 0, then *A* and *B* are *independent*

$$r_{A,B} = \frac{Cov(A,B)}{\sigma_A \sigma_B}$$

- Note that Correlation does not imply causality! That is, if A and B are correlated, this does not necessarily imply that A causes B or that B causes A
- Covarience:

  $$Cov(A,B) = E(A.B) - \bar{A}\bar{B}$$

  - Measures how two things change together .
  - *Covariance is +ve → A & B change together, and if $A > \bar{A}$ then $B > \bar{B}$*
  - *Covariance is -ve → one is above its mean and one is below*
  - If *A* and *B* are independent → *Covariance* = 0
  - لو كانوا موجب يبقى الأول لو زاد التانى هيزداد و العكس علاقة طردية
  - Example



</td></tr>
</table>

| Data Reduction | o More is not always better. |
|---|---|
| | o Obtain a **reduced representation** of the data set that is much smaller in volume, yet closely maintains the **integrity** of the original data. |
| | o Data Reduction Startegies |
| |     o Dimensionality Reduction |
| |     o Numerosity Reduction |
| |     o Data Compression |
| **Dimensionality Reduction** | ▪ Reduce the number of attributes under consideration <br> ▪ Methods include: <br>     o wavelet transforms <br>     o principal components analysis (PCA) <br>     o Attribute subset selection  |
| **Numerosity Reduction Techniques** | ▪ Data are replaced or estimated by alternative. <br>     o parametric methods, a model is used to estimate the data (PCA) <br>     o Nonparametric methods histograms, clustering, sampling, and data cube aggregation |
| **Data Compression** | ▪ Reducing the amount of capacity required to store data. <br>     o **lossless** : No loss of information (e.g. Text ) <br>     o **Lossy**: the size of the file is reduced by eliminating data in the file (e.g. Image) |
| **Attribute Subset Selection** | ▪ How can we find a '**good**' subset of the original attributes? <br> ▪ Rmove the redundent or irrelevent attributes <br> ▪ For $n$ attributes, there are $2^n$ possible subsets!!! <br> ▪ Solution: **Heuristic (Greedy) methods** <br>     o while searching for attribute subsets, they always make what looks to be the best choice at the time. <br> ▪ **Heuristic : Stepwise forward selection {empty} => {Reduced set}** <br>     o The **best** of the attributes is determined and added to the reduced set. <br>     o "**best**" is determined by some predetermined criteria <br> ▪ **Heuristic : Stepwise backward selection {Fill} => {Reduced Set}** <br>     o start with the full set of attributes. <br>     o At each step, remove the worst attribute remaining in the set <br> ▪ **Heuristic : Compination of Stepwise Forward and Backward** <br> ▪ **Heuristic: Decision tree induction** <br>     o **classification** و دى هتندرس في شابتر ال <br><br>  |

| Regression | <ul><li>$y = wx + b$</li><li>y (response variable), can be modeled as a linear function of x (predictor variable)</li><li>W (slope) and b (intercept) could be optimized to get the best fitting</li></ul> |
|---|---|

| X | Y |
|---|---|
| 1.00 | 1.00 |
| 2.00 | 2.00 |
| 3.00 | 1.30 |
| 4.00 | 3.75 |
| 5.00 | 2.25 |



$y = 0.425x + 0.785$

---

**Histograms (binning)**

- The following data are a list of AllElectronics prices for commonly sold items (rounded to the nearest dollar). The numbers have been sorted:
  - 1, 1, 5, 5, 5,5, 5, 8, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18,18, 18, 18, 20, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30,30, 30.



---

**Sampling**

- Obtain (smaller) subsets of the dataset called data sample.

- Simple random sample without replacement (SRSWOR) of size s: all tuples are equally likely to be sampled.

- Simple random sample with replacement (SRSWR) of size s: similar to SRSWOR, but a tuple is drawn recorded then placed back so it may be drawn again

- **Cluster sample** : non overlapping

- **Stratified sample** : if the tuples are divided into strata (overlapping)



---

**Data Transformation**

- Data are transformed into forms **appropriate** for mining
- **Transformation** Strategies
  - **Smoothing**
  - **Attribute Selection**
  - **Aggregation** For example, the daily sales data may be aggregated so as to compute monthly and annual total amounts.
  - **Normalization**: scaling values
  - **Discretization**: (e.g., age) are replaced by interval labels (e.g., 0–10, 11–20, etc.)
  - **Concept Hierarchy**: street can be generalized to higher-level concepts, like city or country

| | |
|---|---|
| **Transformation by Normalization** | ▪ To help avoid dependence on the choice of measurement units<br>▪ Normalizing the data attempts to give all attributes an equal weight<br>    ○ **Min-max normalization**<br><br>• $v = \frac{v-min}{max-min}(new_{max} - new_{min}) + new_{min}$<br><br>• Suppose that the minimum and maximum values for the attribute age are 13 and 70, respectively. We would like to map age to the range [0.0, 1.0].<br><br>• By min-max normalization, a value of 35 for age is transformed to $map(35) = \frac{35-13}{70-13}(1-0) + 0 = 0.39$<br><br>    ○ **Z-score normalization**<br><br>• Normalized based on the mean and standard deviation .<br><br>• $v = \frac{v-mean}{standard\ deviation}$<br><br>• Useful when the actual minimum and maximum of attribute A are unknown, or<br><br>• when there are outliers that dominate the min-max normalization |
| **Concept Hierarchy Generation** | ▪ • It Recursively reduce data by replacing low level concepts (e.g. age values) by higher level concepts (e.g. age groups: youth, adult, or senior).<br>▪ explicitly specified by domain experts<br>▪ formed for both **numeric** and **nominal** data |