

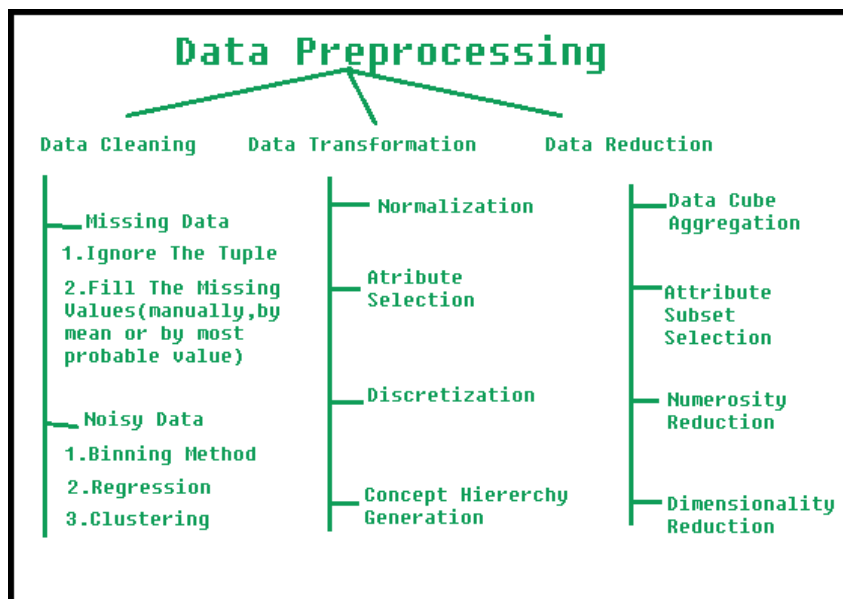


Data Preprocessing in Data Mining

Difficulty Level : Medium • Last Updated : 29 Jun, 2021

Preprocessing in Data Mining:

Data preprocessing is a data mining technique which is used to transform the raw data in a useful and efficient format.



Steps Involved in Data Preprocessing:

1. Data Cleaning:

The data can have many irrelevant and missing parts. To handle this part, data cleaning is done. It involves handling of missing data, noisy data etc.

- **(a). Missing Data:**

This situation arises when some data is missing in the data. It can be handled in various ways.

Some of them are:

1. **Ignore the tuples:**

This approach is suitable only when the dataset we have is quite large and multiple values are missing within a tuple.

2. **Fill the Missing values:**

There are various ways to do this task. You can choose to fill the missing values manually, by attribute mean or the most probable value.

- **(b). Noisy Data:**

Noisy data is a meaningless data that can't be interpreted by machines. It can be generated due to faulty data collection, data entry errors etc. It can be handled in following ways :

1. **Binning Method:**

This method works on sorted data in order to smooth it. The whole data is divided into segments of equal size and then various methods are performed to complete

 **Start Your Coding Journey Now!**

[Login](#)[Register](#)

segment by its mean or boundary values can be used to complete the task.

2. Regression:

Here data can be made smooth by fitting it to a regression function. The regression used may be linear (having one independent variable) or multiple (having multiple independent variables).

3. Clustering:

This approach groups the similar data in a cluster. The outliers may be undetected or it will fall outside the clusters.

2. Data Transformation:

This step is taken in order to transform the data in appropriate forms suitable for mining process. This involves following ways:

1. Normalization:

It is done in order to scale the data values in a specified range (-1.0 to 1.0 or 0.0 to 1.0)

2. Attribute Selection:

In this strategy, new attributes are constructed from the given set of attributes to help the mining process.

3. Discretization:

This is done to replace the raw values of numeric attribute by interval levels or conceptual levels.

4. Concept Hierarchy Generation:

Here attributes are converted from lower level to higher level in hierarchy. For Example-The attribute "city" can be converted to "country".

3. Data Reduction:

Since data mining is a technique that is used to handle huge amount of data. While



Start Your Coding Journey Now!

Login

Register

rid of this, we use data reduction technique. It aims to increase the storage efficiency and reduce data storage and analysis costs.

The various steps to data reduction are:

1. Data Cube Aggregation:

Aggregation operation is applied to data for the construction of the data cube.

2. Attribute Subset Selection:

The highly relevant attributes should be used, rest all can be discarded. For performing attribute selection, one can use level of significance and p-value of the attribute. The attribute having p-value greater than significance level can be discarded.

3. Numerosity Reduction:

This enables to store the model of data instead of whole data, for example: Regression Models.

4. Dimensionality Reduction:

This reduces the size of data by encoding mechanisms. It can be lossy or lossless. If after reconstruction from compressed data, original data can be retrieved, such reduction is called lossless reduction else it is called lossy reduction. The two effective methods of dimensionality reduction are: Wavelet transforms and PCA (Principal Component Analysis).



Like 89

 Start Your Coding Journey Now!

Login

Register

RECOMMENDED ARTICLES

Page : 1 2 3

- | | | | |
|----|--|----|--|
| 01 | Difference Between Data Mining and Text Mining
12, Apr 20 | 05 | Types of Sources of Data in Data Mining
11, Jun 18 |
| 02 | Difference Between Data Mining and Web Mining
10, Apr 20 | 06 | Data Normalization in Data Mining
13, Jun 19 |
| 03 | Text Mining in Data Mining
29, May 21 | 07 | Data Mining: Data Attributes and Quality
17, Jan 20 |
| 04 | Generalized Sequential Pattern (GSP) Mining in Data Mining
20, Mar 22 | 08 | Data Reduction in Data Mining
27, Jan 20 |

Article Contributed By :



deepak_jain

@deepak_jain

Vote for difficulty

Current difficulty : [Medium](#)

Easy

Normal

Medium

Hard

Expert

 Start Your Coding Journey Now!

Login

Register

Improved By : [deepak_jain](#), [abhishekolympics](#), [mukuljain1092](#)

Article Tags : [data mining](#), [DBMS](#)

Practice Tags : [Data Mining](#), [DBMS](#)

[Improve Article](#)[Report Issue](#)

Writing code in comment? Please use ide.geeksforgeeks.org, generate link and share the link here.

[Load Comments](#)

A-143, 9th Floor, Sovereign Corporate Tower,
Sector-136, Noida, Uttar Pradesh - 201305

feedback@geeksforgeeks.org

Company

[About Us](#)
[Careers](#)
[In Media](#)
[Contact Us](#)
[Privacy Policy](#)
[Copyright Policy](#)

Learn

[Algorithms](#)
[Data Structures](#)
[SDE Cheat Sheet](#)
[Machine learning](#)
[CS Subjects](#)
[Video Tutorials](#)
[Courses](#)

 **Start Your Coding Journey Now!**

[Login](#)[Register](#)

Top News	Python
Technology	Java
Work & Career	CPP
Business	Golang
Finance	C#
Lifestyle	SQL
Knowledge	Kotlin

Web Development

[Web Tutorials](#)
[Django Tutorial](#)
[HTML](#)
[JavaScript](#)
[Bootstrap](#)
[ReactJS](#)
[NodeJS](#)

Contribute

[Write an Article](#)
[Improve an Article](#)
[Pick Topics to Write](#)
[Write Interview Experience](#)
[Internships](#)
[Video Internship](#)

@geeksforgeeks , Some rights reserved