

```
cells": [
    {
        "cell_type": "markdown",
        "metadata": {},
        "source": [
            "# Module 4: Data Preprocessing\n",
            "\n",
            "This notebook contains Python examples for data preprocessing. You should refer to the `Data` chapter of the `Introduction to Data Mining` book (slides are available at https://www-users.cs.umn.edu/~kumar01/dmbook/index.php) to understand some of the concepts introduced in this tutorial. The notebook can be downloaded from http://www.cse.msu.edu/~ptan/dmbook/tutorials/tutorial4/tutorial4.ipynb.\n",
            "\n",
            "Data preprocessing consists of a broad set of techniques for cleaning, selecting, and transforming data to improve data mining analysis. Read the step-by-step instructions below carefully. To execute the code, click on the corresponding cell and press the SHIFT-ENTER keys simultaneously."
        ],
        "execution_count": null,
        "outputs": [],
        "source": [
            "# Poor data quality can have an adverse effect on data mining. Among the common data quality issues include noise, outliers, missing values, and duplicate data. This section presents examples of Python code to alleviate some of these data quality problems. We begin with an example dataset from the UCI machine learning repository containing information about breast cancer patients. We will first download the dataset using Pandas read_csv() function and display its first 5 data points.\n",
            "\n",
            "***font color='red'>Code:</font>***\n",
            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"import pandas as pd\\n\", \"data = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.data', header=None)\\n\", \"data.columns = ['sample_code', 'Clump Thickness', 'Uniformity of Cell Size', 'Uniformity of Cell Shape',\\n\", \"Marginal Adhesion', 'Single Epithelial Cell Size', 'Bare Nuclei', 'Bland Chromatin',\\n\", 'Normal Nucleoli', 'Mitoses', 'Class']\\n\", \"data = data.drop(['Sample code'],axis=1)\\n\", \"print('Number of instances = %d' % (data.shape[0]))\\n\", \"print('Number of attributes = %d' % (data.shape[1]))\\n\", \"data.head()\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"## 4.1.1 Missing Values\\n\", \"It is not unusual for an object to be missing one or more attribute values. In some cases, the information was not collected; while in other cases, some attributes are inapplicable to the data instances. This section presents examples on the different approaches for handling missing values. \\n\", \"According to the description of the data (https://archive.ics.uci.edu/ml/datasets/breast-cancer-wisconsin+original), the missing values are encoded as '?' in the original data. Our first task is to convert the missing values to NaNs. We can then count the number of missing values in each column of the data.\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"import numpy as np\\n\", \"data = data.replace('?',np.NaN)\\n\", \"\\n\", \"print('Number of instances = %d' % (data.shape[0]))\\n\", \"print('Number of attributes = %d' % (data.shape[1]))\\n\", \"\\n\", \"print('Number of missing values:')\\n\", \"for col in data.columns:\\n\", \"    print('\\t%s: %d' % (col,data[col].isna().sum()))\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"The boxplots in the table that contain outliers. Note that the values in all columns (except for 'Bare Nuclei') are originally stored as 'int64' whereas the median value of that column. The values before and after replacement are shown for a subset of the data points.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"data2 = data['Bare Nuclei']\\n\", \"\\n\", \"print('Before replacing missing values:')\\n\", \"print(data2[20:25])\\n\", \"data2 = data2.fillna(data2.median())\\n\", \"\\n\", \"print('\\\\nAfter replacing missing values:')\\n\", \"print(data2[20:25])\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"Instead of replacing the missing values, another common approach is to discard the data points that contain missing values. This can be easily accomplished by applying the dropna() function to the data frame.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"data2 = data.dropna()\\n\", \"print('Number of rows in original data = %d' % (data.shape[0]))\\n\", \"\\n\", \"data2 = data.dropna()\\n\", \"print('Number of rows after discarding missing values = %d' % (data2.shape[0]))\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"## 4.1.2 Outliers\\n\", \"Outliers are data instances with characteristics that are considerably different from the rest of the dataset. In the example code below, we will draw a boxplot to identify the outliers in the table that contain outliers. Note that the values in all columns (except for 'Bare Nuclei') are originally stored as 'int64' whereas the median value of that column. The values before and after replacement are shown for a subset of the data points.\\n\", \"\\n\", \"We must convert the column into numeric values first before creating the boxplot. Otherwise, the column will not be displayed when drawing the boxplot.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"\\n\", \"data2 = data.drop(['Class'],axis=1)\\n\", \"data2['Bare Nuclei'] = pd.to_numeric(data2['Bare Nuclei'])\\n\", \"data2.boxplot(figsize=(20,3))\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"The boxplots suggest that only 5 of the columns (Marginal Adhesion, Single Epithelial Cell Size, Bland Chromatin, Normal Nucleoli, and Mitoses) contain abnormally high values. To discard the outliers, we can compute the Z-score for each attribute and remove those instances containing attributes with abnormally high or low Z-scores (e.g., if Z > 3 or Z <= -3). \\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"\\n\", \"The following code shows the results of standardizing the columns of the data. Note that missing values (NaN) are not affected by the standardization process.\"\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"z = (data2-data2.mean())/data2.std()\\n\", \"z[20:25]\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"## 4.1.3 Duplicate Data\\n\", \"Some datasets, especially those obtained by merging multiple data sources, may contain duplicates or near duplicate instances. The term deduplication is often used to refer to the process of dealing with duplicate data issues. \\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"\\n\", \"In the following example, we first check for duplicate instances in the breast cancer dataset.\"\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"dups = data.duplicated()\\n\", \"print('Number of duplicate rows = %d' % (dups.sum()))\\n\", \"data.loc[[11,28]]\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"The duplicated() function will return a Boolean array that indicates whether each row is a duplicate of a previous row in the table. The results suggest there are 236 duplicate rows in the breast cancer dataset. For example, the instance with row index 11 has identical attribute values as the instance with row index 28. Although such duplicate rows may correspond to samples for different individuals, in this hypothetical example, we assume that the duplicates are samples taken from the same individual and illustrate below how to remove the duplicated rows.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"\\n\", \"data2 = data.drop_duplicates()\\n\", \"print('Number of rows after discarding duplicates = %d' % (data2.shape[0]))\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"## 4.2 Aggregation\\n\", \"Data aggregation is a preprocessing task where the values of two or more objects are combined into a single object. The motivation for aggregation includes (1) reducing the size of data to be processed, (2) changing the granularity of analysis (from fine-scale to coarser-scale), and (3) improving the stability of the data.\\n\", \"\\n\", \"In the example below, we will use the daily precipitation time series data for a weather station located at Detroit Metro Airport. The raw data was obtained from the Climate Data Online website (https://www.ncdc.noaa.gov/cdo-web/). The daily precipitation time series will be compared against its monthly values.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"\\n\", \"The code below will load the precipitation time series data and draw a line plot of its daily time series.\"\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"daily = pd.read_csv('DTM_prec.csv', header='infer')\\n\", \"daily.index = pd.to_datetime(daily['DATE'])\\n\", \"daily = daily['PRCP']\\n\", \"ax = daily.plot(kind='line',figsize=(15,3))\\n\", \"ax.set_title('Daily Precipitation (variance = %.4f)' % (daily.var()))\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"Observe that the daily time series appear to be quite chaotic and varies significantly from one time step to another. The time series can be grouped and aggregated by month to obtain the total monthly precipitation values. The resulting time series appears to vary more smoothly compared to the daily time series.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"monthly = daily.groupby(pd.Grouper(freq='M')).sum()\\n\", \"ax = monthly.plot(kind='line',figsize=(15,3))\\n\", \"ax.set_title('Monthly Precipitation (variance = %.4f)' % (monthly.var()))\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"In the example below, the daily precipitation time series are grouped and aggregated by year to obtain the annual precipitation values. \\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"annual = daily.groupby(pd.Grouper(freq='Y')).sum()\\n\", \"ax = annual.plot(kind='line',figsize=(15,3))\\n\", \"ax.set_title('Annual Precipitation (variance = %.4f)' % (annual.var()))\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"## 4.3 Sampling\\n\", \"Sampling is an approach commonly used to facilitate (1) data reduction for exploratory data analysis and scaling up algorithms to big data applications and (2) quantifying uncertainties due to varying data distributions. There are various methods available for data sampling, such as sampling without replacement, where each selected instance is removed from the dataset, and sampling with replacement, where each selected instance is not removed, thus allowing it to be selected more than once in the sample.\\n\", \"\\n\", \"In the example below, we will apply sampling with replacement and without replacement to the breast cancer dataset obtained from the UCI machine learning repository.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"\\n\", \"We initially display the first five records of the table.\"\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"\\n\", \"In the following code, a sample of size 3 is randomly selected (without replacement) from the original data.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"sample = data.sample(n=3)\\n\", \"sample\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"For the equal frequency method, we can apply the qcut() function to discretize the attribute into 4 bins of similar interval widths. The value_counts() function can be used to determine the number of instances in each bin.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"bins = pd.qcut(data['Clump Thickness'],4)\\n\", \"bins.value_counts(sort=False)\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"Finally, we perform a sampling with replacement to create a sample whose size is equal to 1% of the entire data. You should be able to observe duplicate instances in the sample by increasing the sample size.\\n\", \"\\n\", \"***font color='red'>Code:</font>***\n                ]\n            },\n            {\n                \"cell_type\": \"code\", \"execution_count\": null, \"metadata\": {}, \"outputs\": [], \"source\": [\n                    \"sample = data.sample(frac=0.01, random_state=1)\\n\", \"sample\"\n                ]\n            },\n            {\n                \"cell_type\": \"markdown\", \"metadata\": {}, \"source\": [\n                    \"## 4.4 Discretization\\n\", \"Discretization is a data preprocessing step that is often used to transform a continuous-valued attribute to a categorical attribute. The example below illustrates two simple but widely-used unsupervised discretization methods (equal width and equal depth) applied to the 'Clump
```