



Photo by Vilmos Heim on Unsplash

## 6 Different Ways to Compensate for Missing Values In a Dataset (Data Imputation with examples)

Popular strategies to statistically impute missing values in a dataset.



Will Badr

Follow

Jan 5, 2019 · 5 min read ★

Many real-world datasets may contain missing values for various reasons. They are often encoded as NaNs, blanks or any other placeholders. Training a model with a dataset that has a lot of missing values can drastically impact the machine learning model's quality.

Some algorithms such as *scikit-learn estimators* assume that all values are numerical and have and hold meaningful value.

One way to handle this problem is to get rid of the observations that have missing data. However, you will risk losing data points with valuable information. A better strategy would be to impute the missing values. In other words, we need to infer those missing values from the existing part of the data. There are three main types of missing data:

- Missing completely at random (MCAR)
- Missing at random (MAR)
- Not missing at random (NMAR)

However, in this article, I will focus on 6 popular ways for data imputation for cross-sectional datasets ( Time-series dataset is a different story ).

## 1- Do Nothing:

That's an easy one. You just let the algorithm handle the missing data. Some algorithms can factor in the missing values and learn the best imputation values for the missing data based on the training loss reduction (ie. XGBoost). Some others have the option to just ignore them (ie. LightGBM — *use\_missing=false*). However, other algorithms will panic and throw an error complaining about the missing values (ie. Scikit learn — LinearRegression). In that case, you will need to handle the missing data and clean it before feeding it to the algorithm.

Let's see some other ways to impute the missing values before training:

---

*Note: All the examples below use the California Housing Dataset from Scikit-learn.*

---

## 2- Imputation Using (Mean/Median) Values:

This works by calculating the mean/median of the non-missing values in a column and then replacing the missing values within each column separately and independently from the others. It can only be used with numeric data.

col1	col2	col3	col4	col5
------	------	------	------	------

col1	col2	col3	col4	col5
------	------	------	------	------

<b>0</b>	2	5.0	3.0	6	NaN	<b>mean()</b>	<b>0</b>	2.0	5.0	3.0	6.0	7.0
<b>1</b>	9	NaN	9.0	0	7.0	→	<b>1</b>	9.0	11.0	9.0	0.0	7.0
<b>2</b>	19	17.0	NaN	9	NaN		<b>2</b>	19.0	17.0	6.0	9.0	7.0

Mean Imputation

**Pros:**

- Easy and fast.
- Works well with small numerical datasets.

**Cons:**

- Doesn't factor the correlations between features. It only works on the column level.
- Will give poor results on encoded categorical features (do NOT use it on categorical features).
- Not very accurate.
- Doesn't account for the uncertainty in the imputations.

```

1  from sklearn.datasets import fetch_california_housing
2  from sklearn.linear_model import LinearRegression
3  from sklearn.model_selection import StratifiedKFold
4  from sklearn.metrics import mean_squared_error
5  from math import sqrt
6  import random
7  import numpy as np
8  random.seed(0)
9
10 #Fetching the dataset
11 import pandas as pd
12 dataset = fetch_california_housing()
13 train, target = pd.DataFrame(dataset.data), pd.DataFrame(dataset.target)
14 train.columns = ['0','1','2','3','4','5','6','7']
15 train.insert(loc=len(train.columns), column='target', value=target)
16
17 #Randomly replace 40% of the first column with NaN values

```

```
18 column = train['0']
19 print(column.size)
20 missing_pct = int(column.size * 0.4)
21 i = [random.choice(range(column.shape[0])) for _ in range(missing_pct)]
22 column[i] = np.NaN
23 print(column.shape[0])
24
25 #Impute the values using scikit-learn SimpleImpute Class
26 from sklearn.impute import SimpleImputer
27 imp_mean = SimpleImputer( strategy='mean') #for median imputation replace 'mean' with 'median'
28 imp_mean.fit(train)
29 imputed_train_df = imp_mean.transform(train)
```

mean impute.py hosted with ❤ by GitHub

[view raw](#)

mean/median imputation

### 3- Imputation Using (Most Frequent) or (Zero/Constant) Values:

**Most Frequent** is another statistical strategy to impute missing values and YES!! It works with categorical features (strings or numerical representations) by replacing missing data with the most frequent values within each column.

#### Pros:

- Works well with categorical features.

#### Cons:

- It also doesn't factor the correlations between features.
- It can introduce bias in the data.

```
1 #Impute the values using scikit-learn SimpleImpute Class
2
3 from sklearn.impute import SimpleImputer
4 imp_mean = SimpleImputer( strategy='most_frequent')
5 imp_mean.fit(train)
6 imputed_train_df = imp_mean.transform(train)
```

impute\_most\_frequent.py hosted with ❤ by GitHub

[view raw](#)

Most Frequent Imputation

**Zero or Constant** imputation — as the name suggests — it replaces the missing values with either zero or any constant value you specify



## 4- Imputation Using k-NN:

The  $k$  nearest neighbours is an algorithm that is used for simple classification. The algorithm uses ‘**feature similarity**’ to predict the values of any new data points. This means that the new point is assigned a value based on how closely it resembles the points in the training set. This can be very useful in making predictions about the missing values by finding the  $k$ ’s closest neighbours to the observation with missing data and then imputing them based on the non-missing values in the neighbourhood. Let’s see some example code using `Impute` library which provides a simple and easy way to use KNN for imputation:

```
1  import sys
2  from impute.imputation.cs import fast_knn
3  sys.setrecursionlimit(100000) #Increase the recursion limit of the OS
4
5  # start the KNN training
6  imputed_training=fast_knn(train.values, k=30)
```

impute\_knn.py hosted with ♥ by GitHub

[view raw](#)

KNN Imputation for California Housing Dataset

## How does it work?

It creates a basic mean impute then uses the resulting complete list to construct a KDTree. Then, it uses the resulting KDTree to compute nearest neighbours (NN). After it finds the  $k$ -NNs, it takes the weighted average of them.

**Pros:**

- Can be much more accurate than the mean, median or most frequent imputation methods (It depends on the dataset).

**Cons:**

- Computationally expensive. KNN works by storing the whole training dataset in memory.
- K-NN is quite sensitive to outliers in the data (**unlike SVM**)

## 5- Imputation Using Multivariate Imputation by Chained Equation (MICE)

### Main steps used in multiple imputations [1]

This type of imputation works by filling the missing data multiple times. Multiple Imputations (MIs) are much better than a single imputation as it measures the uncertainty of the missing values in a better way. The chained equations approach is also very flexible and can handle different variables of different data types (ie., continuous or binary) as well as complexities such as bounds or survey skip patterns. For more information on the algorithm mechanics, you can refer to the Research Paper

```
1 from impute.imputation.cs import mice
2
3 # start the MICE training
4 imputed_training=mice(train.values)
```

impute\_mice.py hosted with ♥ by GitHub

[view raw](#)

MICE imputation using impyute

## 6- Imputation Using Deep Learning (Datawig):

This method works very well with categorical and non-numerical features. It is a library that learns Machine Learning models using Deep Neural Networks to impute missing values in a dataframe. It also supports both CPU and GPU for training.

```
1  import datawig
2
3  df_train, df_test = datawig.utils.random_split(train)
4
5  #Initialize a SimpleImputer model
6  imputer = datawig.SimpleImputer(
7      input_columns=['1','2','3','4','5','6','7', 'target'], # column(s) containing information ab
8      output_column= '0', # the column we'd like to impute values for
9      output_path = 'imputer_model' # stores model data and metrics
10 )
11
12 #Fit an imputer model on the train data
13 imputer.fit(train_df=df_train, num_epochs=50)
14
15 #Impute missing values and return original dataframe with predictions
16 imputed = imputer.predict(df_test)
```

impute-datawig.py hosted with ❤ by GitHub

[view raw](#)

Imputation using Datawig

### Pros:

- Quite accurate compared to other methods.
- It has some functions that can handle categorical data (Feature Encoder).
- It supports CPUs and GPUs.

### Cons:

- Single Column imputation.
- Can be quite slow with large datasets.
- You have to specify the columns that contain information about the target column that will be imputed.



## Other Imputation Methods:

### Stochastic regression imputation:

It is quite similar to regression imputation which tries to predict the missing values by regressing it from other related variables in the same dataset plus some random residual value.

### Extrapolation and Interpolation:

It tries to estimate values from other observations within the range of a discrete set of known data points.

### Hot-Deck imputation:

Works by randomly choosing the missing value from a set of related and similar variables.

In conclusion, there is no perfect way to compensate for the missing values in a dataset. Each strategy can perform better for certain datasets and missing data types but may perform much worse on other types of datasets. There are some set rules to decide which strategy to use for particular types of missing values, but beyond that, you should experiment and check which model works best for your dataset.

## References:

- [1] Buuren, S. V., & Groothuis-Oudshoorn, K. (2011). Mice: Multivariate Imputation by Chained Equations in R. Journal of Statistical Software
- <https://impyute.readthedocs.io/en/master/index.html>