

Python Feature Engineering Cookbook

Stephen CUI¹

2022 年 12 月 21 日

¹cuixuanStephen@gmail.com

目录

Chapter 1

Imputing Missing Data

Missing data, that is, the absence of values for certain observations, is an unavoidable problem in most data sources.

The act of replacing missing data with statistical estimates of missing values is called imputation. The goal of any imputation technique is to produce a complete dataset. There are multiple imputation methods that we can use, depending on whether the data is missing at random, the proportion of missing values, and the machine learning model we intend to use.

1.1 Technical requirements

We will use the Credit Approval dataset from the UCI Machine Learning Repository (<https://archive.ics.uci.edu/>). To prepare the dataset, follow these steps:

1. Visit <https://archive.ics.uci.edu/ml/machine-learning-databases/credit-screening/>.
2. Click on **crx.data** to download the data.

```
1 import random
2 import numpy as np
3 import pandas as pd
4 data = pd.read_csv('data/crx.data', header=None)
5 varname = [f"A{s}" for s in range(1, 17)]
6 data.columns = varname
7 data = data.replace("?", np.nan)
8 data["A2"] = data["A2"].astype("float")
9 data["A14"] = data["A14"].astype("float")
10 data["A16"] = data["A16"].map({"+": 1, "-": 0})
11 data.rename(columns={"A16": "target"}, inplace=True)
```

```
12 random.seed(37)
13 values = list(set([random.randint(0, len(data)) for p in range(0, 100)]))
14 data.loc[values, ["A3", "A8", "A9", "A10"]] = np.nan
15 data.to_csv('data/credit_approval_uci.csv', index=False)
```

Tip

Make sure you store the dataset in the same folder from which you will execute the code in the recipes.

处理后的数据集有一个target字段以及15个（A1-A15）特征字段组成。其中A11, A12, A13, A15, Target不包含缺失值（??）。

1.2 Removing observations with missing data

Complete Case Analysis (CCA), also called list-wise deletion of cases, consists of discarding observations with missing data. CCA can be applied to both categorical and numerical variables. With CCA, we preserve the distribution of the variables after the imputation, provided the data is missing at random and only in a small proportion of observations. However, if data is missing for many variables, CCA may lead to the removal of a large portion of the dataset.

1.2.1 How to do it...

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 data = pd.read_csv('data/credit_approval_uci.csv')
4 with plt.style.context('ggplot'):
5     data.isnull().mean().sort_values(ascending=True).plot.bar(rot=45)
6     plt.ylabel("Proportion of missing data")
7     plt.title("Proportion of missing data per variable")
```

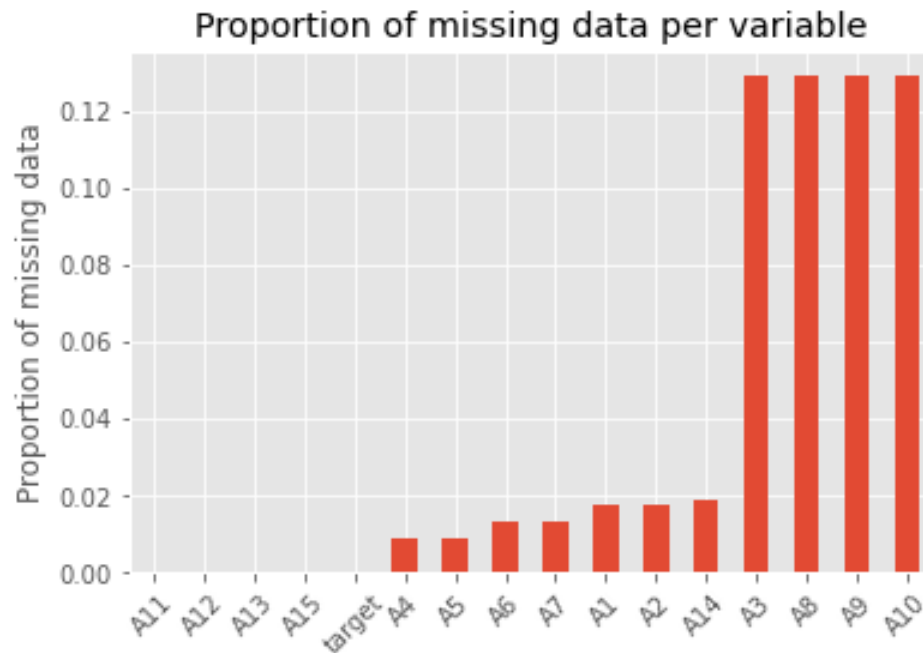


图 1.1: Proportion of missing data per variable

```

1 data_cca = data.dropna()
2 print(f"Total number of observations: {len(data)}")
3 print(f"Number of observations without missing data: {len(data_cca)}")
4 from feature_engine.imputation import DropMissingData
5 cca = DropMissingData(variables=None, missing_only=True)
6 cca.fit(data)
7 print(cca.variables_)
8 data_cca = cca.transform(data)

```

Tip

The `dropna()` method drops observations with any missing value by default. We can remove observations with missing data in a subset of variables like this: `data.dropna(subset=["A3", "A4"])`.

Tip

To remove observations with missing data in a subset of variables, use `DropMissingData(variables=['A3', 'A4'])`. To remove observations with missing values in at least 5% of the variables, use `DropMissingData(threshold=0.95)`.

1.3 Performing mean or median imputation

Mean or median imputation consists of replacing missing values with the mean or median variable. The mean or median is calculated using a train set, and these values are used to impute missing data in train and test sets, as well as in all future data we intend to use with the machine learning model. Scikit-learn and feature-engine transformers learn the mean or median from the train set and store these parameters for future use out of the box.

Tip

Use mean imputation if variables are normally distributed and median imputation otherwise. Mean and median imputation may distort the distribution of the original variables if there is a high percentage of missing data.

1.3.1 How to do it...

```

1  import pandas as pd
2  from sklearn.model_selection import train_test_split
3  from sklearn.impute import SimpleImputer
4  from sklearn.compose import ColumnTransformer
5  from feature_engine.imputation import MeanMedianImputer
6
7  data = pd.read_csv('data/credit_approval_uci.csv')
8
9  X_train, X_test, y_train, y_test = train_test_split(data.drop('target',
10                                                         axis='columns'),
11                                                         data['target'],
12                                                         test_size=.3,
13                                                         random_state=39)
14  numeric_vars = X_train.select_dtypes(exclude='O').columns.to_list()
15  numeric_vars
16  median_values = X_train[numeric_vars].median().to_dict()
17  median_values
18  for df in [X_train, X_test]:
19      # unable to modify X_train and X_test
20      # because df pointer to different address
21      #     df = df.fillna(value=median_values)
22      df.fillna(value=median_values, inplace=True)
23      # equivalent to following codes
24      # X_train = X_train.fillna(value=median_values)
25      # X_test = X_test.fillna(value=median_values)

```

```

26
27 X_train, X_test, y_train, y_test = train_test_split(data.drop('target',
28                                                         axis='columns'),
29                                                         data['target'],
30                                                         test_size=.3,
31                                                         random_state=39)
32 remaining_vars = [var for var in X_test.columns if var not in numeric_vars]
33 remaining_vars
34 imputer = SimpleImputer(strategy='median')
35 ct = ColumnTransformer([('imputer', imputer, numeric_vars)],
36                         remainder='passthrough')
37 ct.fit(X_train)
38 ct.named_transformers_.imputer.statistics_
39 # warning: unable to perform median imputation
40 # for df in [X_train, X_test]:
41 #     df = ct.transform(df)
42 # return NumPy arrays
43 X_train = ct.transform(X_train)
44 X_test = ct.transform(X_test)
45 X_train = pd.DataFrame(X_train, columns=numeric_vars + remaining_vars)
46 X_train
47 X_test = pd.DataFrame(X_test, columns=numeric_vars + remaining_vars)
48 X_test
49
50 X_train, X_test, y_train, y_test = train_test_split(data.drop('target',
51                                                         axis='columns'),
52                                                         data['target'],
53                                                         test_size=.3,
54                                                         random_state=39)
55 imputer = MeanMedianImputer(imputation_method='median', variables=numeric_vars)
56 imputer.fit(X_train)
57 imputer.imputer_dict_
58 X_train = imputer.transform(X_train)
59 X_test = imputer.transform(X_test)
60 print(X_train[numeric_vars].isnull().any().any())
61 print(X_test[numeric_vars].isnull().any().any())

```

The mean or median values should be learned from the train set variables. Thus, we divided the dataset into train and test sets using `scikit-learn`'s `train_test_split()` function. The function takes the predictor variables, the target, the fraction of observations to retain in the test set, and a `random_state` value for reproducibility as arguments. We obtained a train set with 70% of the original observations and a test set with 30% of the original observations. The 70:30 split was done at random.

Tip

`SimpleImputer()` operates on the entire `DataFrame` and returns NumPy arrays. To impute just a subset of variables, we need to use `ColumnTransformer()`. In contrast, `MeanMedianImputer()` can take an entire `DataFrame` and yet it will only impute the specified variables, returning a pandas `DataFrame`.

Categorical variables usually contain strings as values, instead of numbers. We replace missing data in categorical variables with the most frequent category, or with a different string. Frequent categories are estimated using the train set and then used to impute values in the train, test, and future datasets. Thus, we need to learn and store these values, which we can do using `scikit-learn` and `feature-engine`’s out-of-the-box transformers.

[illegible]


```

13                                     random_state=37)
14 # using pandas
15 categorical_vars = X_train.select_dtypes(include='O').columns.to_list()
16 # important: iloc[0]
17 frequent_values = X_train[categorical_vars].mode().iloc[0].to_dict()
18 print(frequent_values)
19
20 X_train = X_train.fillna(value=frequent_values)
21 X_test = X_test.fillna(value=frequent_values)
22 print(X_train[categorical_vars].isnull().any().any())
23 print(X_test[categorical_vars].isnull().any().any())
24
25 # Imputation with a string
26 X_train, X_test, y_train, y_test = train_test_split(data.drop('target',
27                                                         axis='columns'),
28                                                         data['target'],
29                                                         test_size=.3,
30                                                         random_state=37)
31 imputation_dict = {var: 'no_data' for var in categorical_vars}
32 X_train = X_train.fillna(value=imputation_dict)
33 X_test = X_test.fillna(value=imputation_dict)
34 print(X_train['A1'].value_counts())
35
36 # using scikit-learn
37 X_train, X_test, y_train, y_test = train_test_split(data.drop('target',
38                                                         axis='columns'),
39                                                         data['target'],
40                                                         test_size=.3,
41                                                         random_state=37)
42
43 remaining_vars = [
44     var for var in X_train.columns if var not in categorical_vars
45 ]
46
47 imputation_dict = {var: 'not_data' for var in categorical_vars}
48 imputer = SimpleImputer(strategy='most_frequent')
49
50 ct = ColumnTransformer([('Imputer', imputer, categorical_vars)],

```

```

51         remainder='passthrough')
52
53     ct.fit(X_train)
54     # this string "Imputer" here must be identical to
55     # the string 'Imputer' in ColumnTransformer
56     print(ct.named_transformers_.Imputer.statistics_)
57
58     X_train = ct.transform(X_train)
59     X_test = ct.transform(X_test)
60
61     X_train = pd.DataFrame(X_train, columns=categorical_vars + remaining_vars)
62     print(X_train[categorical_vars].isnull().any().any())
63
64     # using feature-engine
65     X_train, X_test, y_train, y_test = train_test_split(data.drop('target',
66                                                                axis='columns'),
67                                                         data['target'],
68                                                         test_size=.3,
69                                                         random_state=37)
70
71     imputer = CategoricalImputer(imputation_method='frequent', variables=categorical_vars)
72
73     imputer.fit(X_train)
74
75     print(imputer.imputer_dict_)
76
77     X_train = imputer.transform(X_train)
78     X_test = imputer.transform(X_test)
79     print(X_train[categorical_vars].isnull().any().any())
80     print(X_test[categorical_vars].isnull().any().any())

```

1.4.2 How it works...

`SimpleImputer()` returns a NumPy array by default. We converted this array into a pandas DataFrame. We had to pass the variable names in the correct order: the imputed variables are located first in the array, followed by the remaining variables. 这个需要特别的注意。

Tip

Note that, unlike `SimpleImputer()`, `CategoricalImputer()` will only impute categorical variables, unless specifically told not to do so by setting the `ignore_format` parameter to `True`.

1.4.3 Replacing missing values with an arbitrary number

When replacing missing values with arbitrary numbers, we need to be careful not to select a value close to the mean, the median, or any other common value of the distribution.

Tip

Arbitrary number imputation can be used when data is not missing at random, when we are building non-linear models, and when the percentage of missing data is high. This imputation technique distorts the original variable distribution.

1.4.4 How to do it...

```
1  # Replacing missing values with an arbitrary number
2
3  import pandas as pd
4  from sklearn.model_selection import train_test_split
5  from sklearn.impute import SimpleImputer
6  from feature_engine.imputation import ArbitraryNumberImputer
7
8  data = pd.read_csv('../data/credit_approval_uci.csv')
9
10 X_train, X_test, y_train, y_test = train_test_split(data.drop('target', axis='columns'),
11                                                    data['target'],
12                                                    test_size=.3,
13                                                    random_state=37)
14
15 arbitrary_cols = ['A2', 'A3', 'A8', 'A11']
16 print(X_train[arbitrary_cols].max().max())
17
18 X_train[arbitrary_cols] = X_train[arbitrary_cols].fillna(99)
19 X_test[arbitrary_cols] = X_test[arbitrary_cols].fillna(99)
20 print(X_train[arbitrary_cols].isnull().any().any())
21 print(X_test[arbitrary_cols].isnull().any().any())
22
```

```

23  # using scikit-learn
24  X_train, X_test, y_train, y_test = train_test_split(data.drop('target', axis='columns'),
25                                                    data['target'],
26                                                    test_size=.3,
27                                                    random_state=37)
28  imputer = SimpleImputer(strategy='constant', fill_value=99)
29  imputer.fit(X_train[arbitrary_cols])
30  X_train[arbitrary_cols] = imputer.transform(X_train[arbitrary_cols])
31  X_test[arbitrary_cols] = imputer.transform(X_test[arbitrary_cols])
32  print(X_train[arbitrary_cols].isnull().any().any())
33  print(X_test[arbitrary_cols].isnull().any().any())
34
35
36  # using feature-engine
37  X_train, X_test, y_train, y_test = train_test_split(data.drop('target', axis='columns'),
38                                                    data['target'],
39                                                    test_size=.3,
40                                                    random_state=37)
41  imputer = ArbitraryNumberImputer(arbitrary_number=99, variables=arbitrary_cols)
42
43  X_train = imputer.fit_transform(X_train)
44  X_test = imputer.transform(X_test)
45  print(X_train[arbitrary_cols].isnull().any().any())
46  print(X_test[arbitrary_cols].isnull().any().any())

```

第16行，首先要找到需要处理缺失值字段的最大值，然后将填充的任意值设置为比这个最大值要大。

第18-19行，也可以使用`imputation_dict`为每列单独设置填充值，将字典的key和value分别设置为字段名和想要填充的任意值。

第28行，如果你的数据集存在有缺失值的分类变量，那么`SimpleImputer()`也会用value将其填充。

On line 41, `ArbitraryNumberImputer()` can automatically select all numerical variables in the train set if we set the `variables` parameter to `None`.

1.5 Finding extreme values for imputation

Replacing missing values with a value at the end of the variable distribution (extreme values) is equivalent to replacing them with an arbitrary value, but instead of identifying the arbitrary values manually, these values are automatically selected as those at the very end of the variable distribution. Missing data can be replaced with a value that is greater or smaller than the remaining values in the variable. To select a value that is greater, we

can use the mean plus a factor of the standard deviation, or the 75th quantile + (IQR * 1.5), where IQR is the IQR given by the 75th quantile - the 25th quantile. To replace missing data with values that are smaller than the remaining values, we can use the mean minus a factor of the standard deviation, or the 25th quantile – (IQR * 1.5).

Tip

End-of-tail imputation may distort the distribution of the original variables, so it may not be suitable for linear models.

1.5.1 How to do it...

```
1  import pandas as pd
2  from sklearn.model_selection import train_test_split
3  from feature_engine.imputation import EndTailImputer
4
5  data = pd.read_csv('../data/credit_approval_uci.csv')
6
7  numeric_vars = [
8      var for var in data.select_dtypes(exclude='O').columns.to_list() if var != 'target'
9  ]
10 print(numeric_vars)
11
12 X_train, X_test, y_train, y_test = train_test_split(
13     data[numeric_vars],
14     data['target'],
15     test_size=.3,
16     random_state=0)
17
18 IQR = X_train.quantile(0.75) - X_train.quantile(0.25)
19 print(IQR)
20
21 imputation_dict = (X_train.quantile(.75) + 1.5 * IQR).to_dict()
22 print(imputation_dict)
23
24 X_train = X_train.fillna(value=imputation_dict)
25 X_test = X_test.fillna(value=imputation_dict)
26
27 print(X_train.isnull().any().any())
```

```

28 print(X_test.isnull().any().any())
29
30 # We can also replace missing data with values at the left tail of the distribution
31 X_train, X_test, y_train, y_test = train_test_split(
32     data[numeric_vars],
33     data['target'],
34     test_size=.3,
35     random_state=0)
36
37 imputer = EndTailImputer(
38     imputation_method='iqr',
39     tail='right',
40     fold=3,
41     variables=None)
42 imputer.fit(X_train)
43
44 print(imputer.imputer_dict_)
45
46 X_train = imputer.transform(X_train)
47 X_test = imputer.transform(X_test)
48
49 print(X_train.isnull().any().any())
50 print(X_test.isnull().any().any())

```

1.5.2 How it works...

On line 21, If we want to use the Gaussian approximation instead of the IQR proximity rule, we can calculate the value to replace missing data using `imputation_dict = (X_train.mean() + 3 * X_train.std()).to_dict()`.

On line 37-41, To use the mean and standard deviation to calculate the replacement values, we need to set `imputation_method="Gaussian"`. We can use 'left' or 'right' in the tail argument to specify the side of the distribution where we'll place the missing values.

1.6 Marking imputed values

A missing indicator is a binary variable that takes the value 1 or True to indicate whether a value was missing, or 0 or False otherwise. It is common practice to replace missing observations with the mean, median, or most frequent category while simultaneously marking those missing observations with missing indicators.

在数据库中，使用null查询会降低查询的性能，因此常见的操作是将缺失值转化为0来进行处理。

1.6.1 How to do it...

```
1  import pandas as pd
2  import numpy as np
3  from sklearn.model_selection import train_test_split
4  from sklearn.impute import SimpleImputer
5  from sklearn.compose import ColumnTransformer
6  from sklearn.pipeline import Pipeline
7  from feature_engine.imputation import (
8      AddMissingIndicator,
9      CategoricalImputer,
10     MeanMedianImputer,
11 )
12
13 data = pd.read_csv('../data/credit_approval_uci.csv')
14
15 X_train, X_test, y_train, y_test = train_test_split(
16     data.drop('target', axis='columns'), data['target'],
17     test_size=.3, random_state=37
18 )
19
20 varnames = ['A1', 'A3', 'A4', 'A5', 'A6', 'A7', 'A8']
21
22 indicators = [f'{var}_na' for var in varnames]
23
24 X_train[indicators] = X_train[varnames].isnull().astype(int)
25 X_test[indicators] = X_test[varnames].isnull().astype(int)
26
27 X_train.sample(5)
28
29 # using feature-engine
30 X_train, X_test, y_train, y_test = train_test_split(
31     data.drop('target', axis='columns'), data['target'],
32     test_size=.3, random_state=37
33 )
34
35 imputer = AddMissingIndicator(
36     variables=None,
37     missing_only=True,
```

```
38 )
39 imputer.fit(X_train)
40 imputer.variables_
41
42 X_train = imputer.transform(X_train)
43 X_test = imputer.transform(X_test)
44
45 pipe = Pipeline([
46     ('ind', AddMissingIndicator(missing_only=True)),
47     ('cat', CategoricalImputer(imputation_method='frequent')),
48     ('num', MeanMedianImputer(imputation_method='mean'))
49 ])
50
51 X_train = pipe.fit_transform(X_train)
52 X_test = pipe.fit_transform(X_test)
53
54 # using scikit-learn
55 X_train, X_test, y_train, y_test = train_test_split(
56     data.drop('target', axis='columns'), data['target'],
57     test_size=.3, random_state=37
58 )
59
60 num_vars = X_train.select_dtypes(exclude='O').columns.to_list()
61 cat_vars = X_train.select_dtypes(include='O').columns.to_list()
62
63 pipe = ColumnTransformer(
64     [
65         ('num_imputer',
66          SimpleImputer(strategy='mean',
67                        add_indicator=True),
68          num_vars),
69         ('cat_imputer',
70          SimpleImputer(strategy='most_frequent',
71                        add_indicator=True),
72          cat_vars)]
73 )
74 X_train = pipe.fit_transform(X_train)
75 X_test = pipe.fit_transform(X_test)
```


这段代码中，仔细看Pipeline的操作。

On line 45, 此处的Pipeline首先在DataFrame后面添加了存在missing value的字段，字段名设置为原字段名_na，然后将分类型变量设置用众数填充，数值型变量用均值填充。Feature-engine imputers automatically identify all numerical or categorical variables, modifying only the appropriate variables. So there is no need to slice the data or pass the variable names as arguments to the transformers in this case

1.7 Performing multivariate imputation by chained equations (MICE)

Multivariate imputation methods, as opposed to univariate imputation, use multiple variables to estimate the missing values. In other words, the missing values of a variable are modeled based on the other variables in the dataset. **Multivariate Imputation by Chained Equations (MICE)** models each variable with missing values as a function of the remaining variables and uses that estimate for imputation.

The following steps are required to perform MICE:

1. A simple univariate imputation is performed for every variable with missing data, for example, median imputation.
2. One specific variable is selected, say, var_1, and the missing values are set back to missing.
3. A model is trained to predict var_1 using the remaining variables as input features.
4. The missing values of var_1 are replaced with the new estimates.
5. Steps 2 to 4 are repeated for each of the remaining variables.

这里有个MICE的[例子](#)，更详细的内容见[Paper](#)。

Once all the variables have been modeled based on the rest, a cycle of imputation is concluded. Multiple imputation cycles are carried out, typically 10. The idea is that by the end of the cycles, the distribution of the imputation parameters should have converged, which means that we should have found the best estimates for the missing data.

1.7.1 How to do it...

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import BayesianRidge
4 from sklearn.impute import IterativeImputer
5
6 variables = ['A2', 'A3', 'A8', 'A11', 'A14', 'A15', 'target']
7
8 data = pd.read_csv('../data/credit_approval_uci.csv',
9                    usecols=variables)
```

```

10
11 X_train, X_test, y_train, y_test = train_test_split(
12     data.drop('target', axis='columns'), data['target'],
13     test_size=.3, random_state=37
14 )
15
16 imputer = IterativeImputer(
17     estimator=BayesianRidge(),
18     max_iter=10, random_state=37
19 )
20
21 imputer.fit(X_train)
22 X_train = imputer.transform(X_train)
23 X_test = imputer.transform(X_test)
24 # Remember that scikit-learn returns NumPy arrays and not DataFrames.
25 pd.DataFrame(X_train).isnull().any().any()
26 pd.DataFrame(X_test).isnull().any().any()

```

1.7.2 See also

- [A Multivariate Technique for Multiply Imputing Missing Values Using a Sequence of Regression Models](#)
- [mice: Multivariate Imputation by Chained Equations in R](#)

1.8 Estimating missing data with nearest neighbors

In imputation with **K-Nearest Neighbors (KNN)**, missing values are replaced with the mean value from their k closest neighbors. The neighbors of each observation are found utilizing distances like the Euclidean distance, and the replacement value can be estimated as the mean or weighted mean of the neighbor's value, where further neighbors have less influence on the replacement value.

1.8.1 How to do it...

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.impute import KNNImputer
4
5 variables = ['A2', 'A3', 'A8', 'A11', 'A14', 'A15', 'target']
6 data = pd.read_csv('../data/credit_approval_uci.csv',

```

```
7         usecols=variables)
8
9 X_train, X_test, y_train, y_test = train_test_split(
10     data.drop('target', axis='columns'),
11     data['target'], test_size=.3, random_state=37
12 )
13
14 imputer = KNNImputer(n_neighbors=5, weights='distance')
15 imputer.fit(X_train)
16 X_train = imputer.transform(X_train)
17 X_test = imputer.transform(X_test)
18 pd.DataFrame(X_train).isnull().any().any()
19 pd.DataFrame(X_test).isnull().any().any()
```

Chapter 2

Encoding Categorical Variables

Categorical variables are those whose values are selected from a group of categories or labels. In some categorical variables, the labels have an intrinsic order. These are called ordinal categorical variables. Variables in which the categories do not have an intrinsic order are called nominal categorical variables.

The values of categorical variables are often encoded as strings. To train mathematical or machine learning models, we need to transform those strings into numbers. The act of replacing strings with numbers is called categorical encoding.

补充： 数据的四个等级

1. nominal level (定类等级)
2. ordinal level (定序等级)
3. interval level (定距等级)
4. ratio level (定比等级)

The values of categorical variables are often encoded as strings. To train mathematical or machine learning models, we need to transform those strings into numbers. The act of replacing strings with numbers is called categorical encoding.

2.1 Technical requirements

We will also use the open-source [Category Encoders](#) Python library, which can be installed using pip:

```
1 pip install category_encoders
```

We will also use the Credit Approval dataset(?? ??).

```
1 import random
2 import numpy as np
3 import pandas as pd
4
5 data = pd.read_csv('../data/credit_approval_uci.csv')
6
7 cat_cols = [
8     c for c in data.columns if data[c].dtypes == 'O'
9 ]
10 num_cols = [
11     c for c in data.columns if data[c].dtypes != 'O'
12 ]
13
14 data[num_cols] = data[num_cols].fillna(0)
15 data[cat_cols] = data[cat_cols].fillna('Missing')
16
17 if(not data.isnull().any().any()):
18     print('not exist missing value')
```

2.2 Creating binary variables through one-hot encoding

In one-hot encoding, we represent a categorical variable as a group of binary variables, where each binary variable represents one category. The binary variable takes a value of 1 if the category is present in an observation, or 0 otherwise.

A categorical variable with k unique categories can be encoded using $k - 1$ binary variables. For the Color variable, which has three categories ($k = 3$; red, blue, and green), we need to create two ($k - 1 = 2$) binary variables to capture all the information so that the following occurs:

- If the observation is red, it will be captured by the red variable (red = 1, blue = 0).
- If the observation is blue, it will be captured by the blue variable (red = 0, blue = 1).
- If the observation is green, it will be captured by the combination of red and blue (red = 0, blue = 0).

Encoding into $k - 1$ binary variables is well-suited for linear models. **There are a few occasions in which we may prefer to encode the categorical variables with k binary variables:**

- When training decision trees since they do not evaluate the entire feature space at the same time.
- When selecting features recursively.
- When determining the importance of each category within a variable.

2.2.1 How to do it...

In this recipe, we will compare the one-hot encoding implementations of pandas, scikit-learn, Feature-engine, and Category Encoders.

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3
4 data = pd.read_csv('../data/credit_approval_uci.csv')
5
6
7 # Let's do one-hot encoding using pandas
8 X_train, X_test, y_train, y_test = train_test_split(
9     data.drop(labels=['target'], axis='columns'), data['target'],
10     test_size=.3,
11     random_state=0
12 )
13
14 X_train['A4'].unique()
15
16 dummies = pd.get_dummies(X_train['A4'], drop_first=True)
17 dummies.head()
18
19
20 X_train_enc = pd.get_dummies(X_train, drop_first=True)
21 X_test_enc = pd.get_dummies(X_test, drop_first=True)
22 X_train_enc.head()
23
24
25 X_train_enc.columns
26
27
28 X_test_enc = pd.concat([X_test, X_test_enc], axis='columns')
29
30
31 X_test_enc.drop(
32     labels=X_test_enc.select_dtypes(include='O').columns,
33     axis='columns',
34     inplace=True
35 )
```

```
1 # Let's do one-hot encoding using scikit-learn
2 from sklearn.preprocessing import OneHotEncoder
3
```

```
4 X_train, X_test, y_train, y_test = train_test_split(
5     data.drop(labels=['target'], axis='columns'), data['target'],
6     test_size=.3,
7     random_state=0
8 )
9
10 # handle_unknown='ignore'
11 encoder = OneHotEncoder(drop='first', sparse=False)
12
13 vars_categorical = X_train.select_dtypes(include='O').columns.to_list()
14
15 encoder.fit(X_train[vars_categorical])
16
17 encoder.categories_
18
19 X_train_enc = encoder.transform(
20     X_train[vars_categorical]
21 )
22
23 X_test_enc = encoder.transform(
24     X_test[vars_categorical]
25 )
26
27 encoder.get_feature_names_out()
28
29 X_test_enc = pd.DataFrame(X_test_enc)
30 X_test_enc.columns = encoder.get_feature_names_out()
31
32 X_test_enc.index = X_test.index
33
34 X_test_enc = pd.concat([X_test, X_test_enc], axis='columns')
35
36 X_test_enc.drop(
37     labels=X_test_enc.select_dtypes(include='O').columns,
38     axis='columns',
39     inplace=True
40 )
```

```
1  # Let's perform one-hot encoding with Feature-engine
2  from feature_engine.encoding import OneHotEncoder
3
4  X_train, X_test, y_train, y_test = train_test_split(
5      data.drop(labels=['target'], axis='columns'), data['target'],
6      test_size=.3,
7      random_state=0
8  )
9
10 ohe_enc = OneHotEncoder(drop_last=True)
11 ohe_enc.fit(X_train)
12
13 ohe_enc.variables_
14
15 ohe_enc.encoder_dict_
16
17 X_train_enc = ohe_enc.transform(X_train)
18 X_test_enc = ohe_enc.transform(X_test)
19
20 ohe_enc.get_feature_names_out()
```

2.3 Performing one-hot encoding of frequent categories

One-hot encoding represents each variable's category with a binary variable. Hence, one-hot encoding of highly cardinal variables or datasets with multiple categorical features can expand the feature space dramatically. This, in turn, may increase the computational cost of using machine learning models or deteriorate their performance. To reduce the number of binary variables, we can perform one-hot encoding of the most frequent categories. One-hot encoding the top categories is equivalent to treating the remaining, less frequent categories as a single, unique category.

2.3.1 How to do it...

```
1  import pandas as pd
2  import numpy as np
3  from sklearn.model_selection import train_test_split
4  from feature_engine.encoding import OneHotEncoder
```



```

1 data = pd.read_csv('../data/credit_approval_uci.csv')
2 X_train, X_test, y_train, y_test = train_test_split(
3     data.drop(labels=['target'], axis='columns'), data['target'],
4     test_size=.3,
5     random_state=0
6 )

```

数据集必须在统计最大频数之前切割，避免出现数据泄露。

The most frequent categories need to be determined in the train set. This is to avoid data leakage.

```

1 X_train['A6'].unique()
2 X_train['A6'].value_counts().sort_values(ascending=False).head(5)
3
4 top_5 = [
5     x for x in X_train['A6'].value_counts().sort_values(
6         ascending=False).head(5).index
7 ]
8 top_5
9
10 for label in top_5:
11     X_train[f'A6_{label}'] = np.where(X_train['A6'] == label, 1, 0)
12     X_test[f'A6_{label}'] = np.where(X_test['A6'] == label, 1, 0)
13
14 X_train[['A6']+ [f'A6_{label}' for label in top_5]].sample(10)

```

在下面代码执行之前，需要运行数据集切割的代码。

```

1 ohe_enc = OneHotEncoder(top_categories=5, variables=['A6', 'A7'])
2 ohe_enc.fit(X_train)
3 X_train_enc = ohe_enc.transform(X_train)
4 X_test_enc = ohe_enc.transform(X_test)

```

Feature-engine's `OneHotEncoder()` will encode all categorical variables in the dataset by default unless we specify the variables to encode.

You can view the new binary variables in the DataFrame by executing `X_train_enc.head()`. You can also find the top five categories learned by the encoder by executing `ohe_enc.encoder_dict_`.

Feature-engine replaces the original variable with the binary ones returned by one-hot encoding, leaving the dataset ready to use in machine learning.

2.3.2 There's more...

This recipe is based on the winning solution of the KDD 2009 cup, [Winning the KDD Cup Orange Challenge with Ensemble Selection](#), where the authors limited one-hot encoding to the 10 most frequent categories of each variable.

2.4 Replacing categories with counts or the frequency of observations

In count or frequency encoding, we replace the categories with the count or the fraction of observations showing that category. These encoding methods, which capture the representation of each label in a dataset, are very popular in data science competitions. **The assumption is that the number of observations per category is somewhat predictive of the target.**

Tip

Note that if two different categories are present in the same number of observations, they will be replaced by the same value, which leads to information loss.

如果两个不同的特征具有相同的频数或者counts，那么这两个特征将会具有相同的值，这将导致信息的减少。