

```
1 import pandas as pd
2 import numpy as np
3 import sys
4 from sklearn.impute import KNNImputer, SimpleImputer
5 from sklearn.preprocessing import StandardScaler, PowerTransformer
```

▼ Load data

```
1 if 'google.colab' in sys.modules:
2     from google.colab import drive
3     drive.mount('/content/drive')
4     DATA_PATH = '/content/drive/MyDrive/DL Project/data/'
5 else:
6     DATA_PATH = './Documents/Classes/AML/proj/archive/'

Mounted at /content/drive

1 train = pd.read_feather(DATA_PATH + 'train_data.ftr')
2 train.head()
3 # test = pd.read_feather(DATA_PATH + 'test_data.ftr')

1 CAT_COL = ['B_30', 'B_38', 'D_114', 'D_116', 'D_117', 'D_120', 'D_126', 'D_63', 'D_64', 'D_66', 'D_68']
```

▼ Preprocess

▼ Transform categorical data into numerical data (Laplace smoothing)

```
1 def cat_to_num(df):
2     y_bar = df['target'].mean()
3
4     def apply_laplace_smooth(df, feature_col, k=30):
5         G = feature_col.value_counts()
6         y_per_G = df.groupby(feature_col.name)['target'].sum()
7         x_hat = (k * y_bar + y_per_G) / (k + G)
8         return feature_col.map(x_hat)
9
10    df[CAT_COL] = df[CAT_COL].apply(lambda f: apply_laplace_smooth(df, f)).astype(np.float16)
11
12    return df

1 train = cat_to_num(train)
2 train.head()
```

	customer_ID	S_2	P_2	D_39	B_1	B_2	R_1	S_3	D_41	B_3	...	D_137	D_138
0	0000099d6bd597052cdcda90ffabf56573fe9d7c79be5f...	2017-03-09	0.938477	0.001734	0.008728	1.006836	0.009224	0.124023	0.008774	0.004707	...	NaN	NaN
1	0000099d6bd597052cdcda90ffabf56573fe9d7c79be5f...	2017-04-07	0.936523	0.005775	0.004925	1.000977	0.006153	0.126709	0.000798	0.002714	...	NaN	NaN
2	0000099d6bd597052cdcda90ffabf56573fe9d7c79be5f...	2017-05-28	0.954102	0.091492	0.021652	1.009766	0.006817	0.123962	0.007599	0.009422	...	NaN	NaN
3	0000099d6bd597052cdcda90ffabf56573fe9d7c79be5f...	2017-06-13	0.960449	0.002455	0.013687	1.002930	0.001372	0.117188	0.000685	0.005531	...	NaN	NaN
4	0000099d6bd597052cdcda90ffabf56573fe9d7c79be5f...	2017-07-16	0.947266	0.002483	0.015190	1.000977	0.007607	0.117310	0.004654	0.009308	...	NaN	NaN

5 rows x 191 columns

```
1 train.to_pickle(DATA_PATH + 'train_cat_transed.pkl')
```

▼ Handle missing values

```
1 def impute_mean(df):
2     imputer = SimpleImputer().set_output(transform='pandas')
```

```

3
4     trans_df = imputer.fit_transform(df.drop(['customer_ID', 'S_2'], axis=1))
5     imputed_df = pd.concat([df[['customer_ID', 'S_2']], trans_df], axis=1)
6
7     assert sum(imputed_df.isna().any()) == 0
8
9     return imputed_df

1 def impute_knn(df):
2     imputer = KNNImputer(n_neighbors=3).set_output(transform='pandas')
3     trans_df = imputer.fit_transform(df.drop(['customer_ID', 'S_2'], axis=1))
4     imputed_df = pd.concat([df[['customer_ID', 'S_2']], trans_df], axis=1)
5
6     assert sum(imputed_df.isna().any()) == 0
7
8     return imputed_df

1 def impute_drop(df):
2     # Get all columns that have NaN values
3     nan_cols = df.columns[df.isna().any()].tolist()
4
5     # Drop columns that have > 20% NaN values
6     df = df.dropna(thresh=df.shape[0]*0.8, axis=1)
7
8     df = df.dropna()
9     df = df.reset_index(drop=True)
10
11     # No NaN values at this point
12     assert sum(imputed_df.isna().any()) == 0
13
14     return df

1 def find_impute_val(df, feature_series):
2     feature_name = feature_series.name
3     bin_col_name = feature_name + "_bin"
4
5     # 1. Calculate percentiles (with 10 bins) for each feature that has missing values
6     df[bin_col_name] = pd.qcut(x=feature_series, q=10, labels=False, duplicates='drop')
7
8     # 2. Calculate target label rates for each bin of each feature
9     target_rate_per_bin = df.groupby(bin_col_name)['target'].mean()
10
11     # 3. Determine the bin that the imputed missing values fall inside and assign the target label rate
12     for bin_num in range(len(target_rate_per_bin)):
13         bool_mask = (df[bin_col_name]==bin_num) & (train[feature_name].isna())
14         feature_series.mask(bool_mask, target_rate_per_bin[bin_num])
15
16     df.drop(bin_col_name, axis=1, inplace=True)
17
18     return feature_series

1 def handle_missing_values(df, impute_fn):
2     filled_df = impute_fn(df)
3
4     nan_col = filled_df.columns[filled_df.isna().any()].tolist()
5     filled_df[nan_col] = filled_df[nan_col].apply(lambda f: find_impute_val(df, f))
6
7     assert sum(filled_df.isna().any()) == 0
8
9     return filled_df

1 train = handle_missing_values(train, impute_mean)
2 train.head()

```



```
/Users/roxyk/opt/anaconda3/lib/python3.9/site-packages/numpy/core/_methods.py:236: RuntimeWarning: overflow encountered in multiply
  x = um.multiply(x, x, out=x)
/Users/roxyk/opt/anaconda3/lib/python3.9/site-packages/numpy/core/_methods.py:247: RuntimeWarning: overflow encountered in reduce
  ret = umr_sum(x, axis, dtype, out, keepdims=keepdims, where=where)
```

	customer_ID	S_2	target	P_2	D_39	B_1	B_2	R_1	S_3	D_41	...	D_136
5	0000099d6bd597052cdcda90ffabf56573fe9d7c79be5f...	2017-08-04	0.0	1.167299	-0.753275	-0.596688	1.050376	-0.200328	-0.899846	0.779986	...	-1.110223e-16
		2017-										-1.110223e-
1 train.to_pickle(DATA_PATH + 'train_skewness_reduced.pkl')												
		2017-										-1.110223e-

▼ Normalization

```
1 # scaler = StandardScaler().set_output(transform='pandas')
2 # train = scaler.fit_transform(train)
```

▼ Save processed data

```
1 train.to_pickle(DATA_PATH + 'train.pkl')
```

1