```python
1    import pandas as pd
2    import numpy as np
3    from sklearn.preprocessing import StandardScaler, LabelEncoder
4    from sklearn.model_selection import train_test_split
5    import joblib
6    import os
7
8
9    def get_values(value):
10       return value.values.reshape(-1, 1)
11
12
13   def load_raw():
14       train = pd.read_csv('./data/train.csv')
15       test = pd.read_csv('./data/test.csv')
16
17       categorical_features = ['COMPONENT_ARBITRARY', 'YEAR']
18
19       train = train.fillna(0)
20       test = test.fillna(0)
21       additional_test = train[train["Y_LABEL"] == 1]
22       train = train[train["Y_LABEL"] == 0]
23
24       all_X = train.drop(['ID', 'Y_LABEL'], axis=1)
25       all_y = train['Y_LABEL']
26
27       test = test.drop(['ID'], axis=1)
28       additional_test = additional_test.drop(["ID"], axis=1)[test.columns]
29
30       train_X, val_X, train_y, val_y = train_test_split(all_X, all_y, test_size=0.2)
31
32       scaler = StandardScaler()
33       for col in train_X.columns:
34           if col not in categorical_features:
35               train_X[col] = scaler.fit_transform(get_values(train_X[col]))
36               val_X[col] = scaler.transform(get_values(val_X[col]))
37               if col in test.columns:
38                   test[col] = scaler.transform(get_values(test[col]))
39                   additional_test[col] = scaler.transform(get_values(additional_test[col]))
40       le = LabelEncoder()
41       for col in categorical_features:
42           train_X[col] = le.fit_transform(train_X[col])
43           val_X[col] = le.transform(val_X[col])
44           if col in test.columns:
45               test[col] = le.transform(test[col])
46               additional_test[col] = le.transform(additional_test[col])
47
48       # test = pd.concat([test, additional_test])
49       return train_X, val_X, train_y, val_y, test, additional_test
50
51
52   class CustomDataset:
53       def __init__(self, data_X: pd.DataFrame, data_y, distillation=False):
54           super(CustomDataset, self).__init__()
55           self.data_X = data_X
56           self.data_y = data_y
57           self.distillation = distillation
58           self.test_stage_features = ['COMPONENT_ARBITRARY', 'ANONYMOUS_1', 'YEAR',
59                                       'ANONYMOUS_2', 'AG', 'CO', 'CR', 'CU', 'FE', 'H20',
60                                       'MN', 'MO', 'NI', 'PQINDEX', 'TI', 'V', 'V40', 'ZN']
61
62       def __len__(self):
63           return len(self.data_X)
64
65       def __getitem__(self, index):
66           if self.distillation:
```

```python
                    # 지식 증류 학습 시
                    teacher_X = self.data_X.iloc[index].values
                    student_X = self.data_X[self.test_stage_features].iloc[index].values
                    y = self.data_y.values[index]
                    return teacher_X, student_X, y
            else:
                if self.data_y is None:
                    test_X = self.data_X.iloc[index].values
                    return test_X
                else:
                    teacher_X = self.data_X.iloc[index].values
                    y = self.data_y.values[index]
                    return teacher_X, y


class DataLoader:
    def __init__(self, dataset: CustomDataset, batch_size, shuffle=True):
        self.dataset = dataset
        self.batch_size = batch_size
        self.shuffle = shuffle
        self.on_epoch_end()

    def __len__(self):
        return len(self.dataset) // self.batch_size

    def __iter__(self):
        for item in (self[i] for i in range(len(self))):
            yield item

    def __getitem__(self, idx):
        indices = self.indices[int(idx*self.batch_size): int((idx+1) * self.batch_size)]
        batch_input = self.dataset[indices]
        return batch_input

    def on_epoch_end(self):
        self.indices = list(range(len(self.dataset)))
        if self.shuffle:
            np.random.shuffle(self.indices)
```