Double-click (or enter) to edit

Project: Machine Learning of Salary and Demographic Factors Name: Shaohua Feng Supervisor:

Double-click (or enter) to edit

```
1 from google.colab import drive
2
3 # Mount Google Drive
4 drive.mount('/content/drive')

  Mounted at /content/drive

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt

1 from sklearn.metrics import confusion_matrix
2 #from sklearn.metrics import plot_confusion_matrix

1 # from sklearn.metrics import plot_confusion_matrix doesn't work, so
2 !pip install plot_confusion_matrix

  Collecting plot_confusion_matrix
    Downloading plot_confusion_matrix-0.0.2-py3-none-any.whl (3.6 kB)
  Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from plot_confusion_matrix) (3.7.1)
  Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from plot_confusion_matrix) (1.23.5)
  Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (1.2.0)
  Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (0.12.1)
  Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (4.44.3)
  Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (1.4.5)
  Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (23.2)
  Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (9.4.0)
  Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (3.1.1)
  Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib->plot_confusion_matrix) (2.8.2)
  Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib->plot_confusion_matr
  Installing collected packages: plot_confusion_matrix
  Successfully installed plot_confusion_matrix-0.0.2
```

```
1 ####################################
2 ## Read data and data wrangling
3 ####################################
4 # read in data loaded in google drive
5 file_path_1 = '/content/drive/My Drive/adult.data'
6 adult_1= pd.read_csv(file_path_1,header=None)
7 file_path_2 = '/content/drive/My Drive/adult.test.txt'
8 adult_2= pd.read_csv(file_path_2,header=None)
9 adult=pd.concat([adult_1, adult_2], ignore_index=True)

1 # add column names
2 cols=['age','workclass','fnlwgt','education','education-num','marital-status','occupation','relationship','race','sex','capital-gain','capital-
3 adult.columns=cols
4
5 # add y column to data frame. target=1 for label '>50k' and y=0 for label '<=50k'
6 adult['target']=np.where(adult['label']=='>50K',1,0)
7 #adult['target'] = adult['target'].astype(bool)
8 #
9 print(adult.describe())
10 adult.dtypes
11 adult.info()
```

```
                age        fnlwgt  education-num  capital-gain  capital-loss  \
count  48842.000000  4.884200e+04   48842.000000  48842.000000  48842.000000
mean      38.643585  1.896641e+05      10.078089   1079.067626     87.502314
std       13.710510  1.056040e+05       2.570973   7452.019058    403.004552
min       17.000000  1.228500e+04       1.000000      0.000000      0.000000
25%       28.000000  1.175505e+05       9.000000      0.000000      0.000000
50%       37.000000  1.781445e+05      10.000000      0.000000      0.000000
75%       48.000000  2.376420e+05      12.000000      0.000000      0.000000
max       90.000000  1.490400e+06      16.000000  99999.000000   4356.000000

       hours-per-week        target
count    48842.000000  48842.000000
mean        40.422382      0.160538
std         12.391444      0.367108
min          1.000000      0.000000
25%         40.000000      0.000000
50%         40.000000      0.000000
75%         45.000000      0.000000
max         99.000000      1.000000
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48842 entries, 0 to 48841
Data columns (total 16 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   age             48842 non-null  int64
```

```
 1   workclass        48842 non-null   object
 2   fnlwgt           48842 non-null   int64
 3   education        48842 non-null   object
 4   education-num    48842 non-null   int64
 5   marital-status   48842 non-null   object
 6   occupation       48842 non-null   object
 7   relationship     48842 non-null   object
 8   race             48842 non-null   object
 9   sex              48842 non-null   object
10   capital-gain     48842 non-null   int64
11   capital-loss     48842 non-null   int64
12   hours-per-week   48842 non-null   int64
13   native-country   48842 non-null   object
14   label            48842 non-null   object
15   target           48842 non-null   int64
dtypes: int64(7), object(9)
memory usage: 6.0+ MB
```

```
1 # Data Manipulation: replace '?' with None
2 adult['workclass']=adult['workclass'].replace(' ?',None)
3 adult['occupation']=adult['occupation'].replace(' ?',None)
4 adult['native-country']=adult['native-country'].replace(' ?',None)
```

```
1 # charactegorical columns
2 cols_cat=['workclass','education','marital-status','occupation','relationship','race','sex','native-country']
3
4 for x in cols_cat:
5   adult[x] = adult[x].astype('category')
6   #print(x)
7
8 adult.dtypes
```

```
age               int64
workclass         category
fnlwgt            int64
education         category
education-num     int64
marital-status    category
occupation        category
relationship      category
race              category
sex               category
capital-gain      int64
capital-loss      int64
hours-per-week    int64
native-country    category
label             object
```

```
target    int64
dtype: object
```

```
1 # delete missing value
2 adult_cleaned=adult.dropna()
3 print(len(adult_cleaned))
```

```
45222
```

```
1 # drop original label
2 del adult_cleaned['label']
3 adult_cleaned.head(5)
```

| | age | workclass | fnlwgt | education | education-num | marital-status | occupation | relationship | race | sex | capital-gain | capital-loss |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 39 | State-gov | 77516 | Bachelors | 13 | Never-married | Adm-clerical | Not-in-family | White | Male | 2174 | 0 |
| 1 | 50 | Self-emp-not-inc | 83311 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband | White | Male | 0 | 0 |
| 2 | 38 | Private | 215646 | HS-grad | 9 | Divorced | Handlers-cleaners | Not-in-family | White | Male | 0 | 0 |
| 3 | 53 | Private | 234721 | 11th | 7 | Married-civ-spouse | Handlers-cleaners | Husband | Black | Male | 0 | 0 |

```
1 ####################################
2 # Functions used by the machine learning algorithms
3 #Normalizing numeric data
4 def normalize(x):
5   if x.dtype == 'int' or x.dtype == 'float':
6     return ((x - min(x)) / (max(x) - min(x)))
7   else:
8     return x
```

```
1 #Converting categorical data to dummy/one-hot variables
2 def dummy(x):
3   cat_col=['workclass','education','marital-status','occupation','relationship','race','sex','native-country']
4   x = pd.get_dummies(x, columns=cat_col, prefix = cat_col)
5   return x
6 #print the dataset
7 #adult_new.head(5)
8
```

```
1 # Print out Accuracy
2 def printAcc(y_test,y_pred):
3   from sklearn.metrics import accuracy_score
4
5   accuracy = accuracy_score(y_test, y_pred)
6   print(f"Accuracy: {accuracy}")
```

```
1 # print out confusion matrix
2 def printConfusion(y_test, y_pred):
3
4   from sklearn.metrics import confusion_matrix
5
6   cf=confusion_matrix(y_test, y_pred)
7   print(cf)
8   tn, fp, fn, tp=cf.ravel()
9   print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

```
1 #print precision, recall, and accuracy from the perspective of each of the class (0 and 1 for German dataset)
2 def printReport(y_test,y_pred):
3
4   from sklearn.metrics import classification_report
5   from sklearn import metrics
6
7   print(classification_report(y_test, y_pred))
```

```
1 ##################################################
2 # Decision tree
3 ##################################################
4
5 # Decision tree without oversampling and normalization
6 # I want to use it as baseline to proof that oversampling and normalization improves decision tree
7 from sklearn.tree import DecisionTreeClassifier
8 from sklearn.model_selection import train_test_split
9
10 # Create a deep copy of the data frame adult_cleaned and name it adult_new
11 # adult_new is not normalized or oversampled
12 # I will use adult_new as the base line
13 adult_new=adult_cleaned.copy(deep=True)
14
15 #X = list(set(list(adult_cleaned)) - set(['target']))
16 X = adult_new.drop('target', axis=1)
17 y = adult_new['target']
18
19 # dummy variable
20 X=dummy(X)
21
22 # split train and test. test size is 0.35
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.35, random_state=42)
24
25 # create decision tree classifier
26 dtc = DecisionTreeClassifier(random_state=52)
27 # Train the model on the training data
28 dtc.fit(X_train, y_train)
29 # Make predictions on the test data
30 y_pred = dtc.predict(X_test)
31
32 print("Decision Tree for data without normalization and oversampling")
33 # Print Accuracy
34 printAcc(y_test,y_pred)
35
36 # Print Confusion Matrix
37 print("")
38 print("Confustion Matrix")
39 printConfusion(y_test, y_pred)
40
41 # Print Diagnosis
42 print("")
43 printReport(y_test,y_pred)
44
45
```

```
Decision Tree for data without normalization and oversampling
Accuracy: 0.7998483699772555

Confustion Matrix
[[11507 1697]
 [ 1471 1153]]
TP: 1153 , FP: 1697 , TN: 11507 , FN: 1471

              precision    recall  f1-score   support

           0       0.89      0.87      0.88     13204
           1       0.40      0.44      0.42      2624

    accuracy                           0.80     15828
   macro avg       0.65      0.66      0.65     15828
weighted avg       0.81      0.80      0.80     15828
```

```python
1  # Decision tree with normalization but no oversampling
2
3  from sklearn.tree import DecisionTreeClassifier
4  from sklearn.model_selection import train_test_split
5
6  # Create a deep copy of the data frame adult_cleaned and name it adult_new
7  # adult_new is not normalized or oversampled
8  # I will use adult_new as the base line
9  adult_dt_norm = adult_cleaned.copy(deep=True)
10 X = adult_dt_norm.drop('target', axis=1)
11 y = adult_dt_norm['target']
12
13 # Normalization the numerical columns
14 num_cols = ['age', 'fnlwgt', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']
15 # why this code doesn't work?
16 #X[num_cols] = X[num_cols].apply(lambda x:normalize(x), axis=0)
17 X_num = X[num_cols]
18 X_num_normalized = X_num.apply(normalize, axis=0)
19 # combine the normalized numerical columns with the categorical columns
20 X = pd.concat([X_num_normalized, X.drop(num_cols, axis=1)], axis=1)
21
22 #print(X.head(5))
23 #print(X.tail(5))
24 # dummy variable
25 X=dummy(X)
26 print(X.head(5))
27 # split train and test. test size is 0.35
28 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.35, random_state=52)
29
30 # create decision tree classifier
31 dtc_norm = DecisionTreeClassifier(random_state=52)
32 # Train the model on the training data
33 dtc_norm.fit(X_train, y_train)
34 # Make predictions on the test data
35 y_pred = dtc_norm.predict(X_test)
36
37 print("Decision Tree for data with normalization but no oversampling")
38 # Print Accuracy
39 printAcc(y_test,y_pred)
40
41 # Print Confustion Matrix
42 print("")
43 print("Confustion Matrix")
44 printConfusion(y_test, y_pred)
45
46 # Print Diagnosis
```

```
47  print("")
48  printReport(y_test,y_pred)
49
50
51
52
53
```

```
         age     fnlwgt  education-num  capital-gain  capital-loss  \
0   0.301370   0.043350       0.800000      0.02174           0.0
1   0.452055   0.047274       0.800000      0.00000           0.0
2   0.287671   0.136877       0.533333      0.00000           0.0
3   0.493151   0.149792       0.400000      0.00000           0.0
4   0.150685   0.219998       0.800000      0.00000           0.0

   hours-per-week  workclass_ Federal-gov  workclass_ Local-gov  \
0        0.397959                       0                     0
1        0.122449                       0                     0
2        0.397959                       0                     0
3        0.397959                       0                     0
4        0.397959                       0                     0

   workclass_ Never-worked  workclass_ Private  ...  native-country_ Portugal  \
0                        0                   0  ...                         0
1                        0                   0  ...                         0
2                        0                   1  ...                         0
3                        0                   1  ...                         0
4                        0                   1  ...                         0

   native-country_ Puerto-Rico  native-country_ Scotland  \
0                            0                         0
1                            0                         0
2                            0                         0
3                            0                         0
4                            0                         0

   native-country_ South  native-country_ Taiwan  native-country_ Thailand  \
0                      0                       0                         0
1                      0                       0                         0
2                      0                       0                         0
3                      0                       0                         0
4                      0                       0                         0

   native-country_ Trinadad&Tobago  native-country_ United-States  \
0                                0                              1
1                                0                              1
2                                0                              1
3                                0                              1
4                                0                              0
```

```
   native-country_ Vietnam  native-country_ Yugoslavia
0                        0                           0
1                        0                           0
2                        0                           0
3                        0                           0
4                        0                           0

[5 rows x 105 columns]
Decision Tree for data with normalization but no oversampling
Accuracy: 0.8035759413697245

Confustion Matrix
[[11590  1636]
 [ 1473  1129]]
TP: 1129 , FP:  1636 , TN:  11590 , FN: 1473
```

```
1 # Decision tree with normalization and random oversampling
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.model_selection import train_test_split
4 from imblearn.over_sampling import RandomOverSampler
5
6 # Create a deep copy of the data frame adult_dt_norm_ros
7 # X is normalized and random oversampled
8 adult_dt_norm_ros = adult_cleaned.copy(deep=True)
9 X = adult_dt_norm_ros.drop('target', axis=1)
10 y = adult_dt_norm_ros['target']
11
12 # Normalization the numerical columns
13 num_cols = ['age', 'fnlwgt', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']
14 X_num = X[num_cols]
15 X_num_normalized = X_num.apply(normalize, axis=0)
16 # combine the normalized numerical columns with the categorical columns
17 X = pd.concat([X_num_normalized, X.drop(num_cols, axis=1)], axis=1)
18
19 # dummy variable
20 X=dummy(X)
21
22 # split train and test. test size is 0.35
23 X_train_norm, X_test_norm, y_train_norm, y_test_norm = train_test_split(X, y, test_size=0.35, random_state=52)
24
25 #Random Oversampling
26 ros = RandomOverSampler(sampling_strategy='auto', random_state=52)
27 X_train_ros, y_train_ros = ros.fit_resample(X_train_norm, y_train_norm)
28
29 # create decision tree classifier for normalized data with random oversampling
30 dtc_norm_ros = DecisionTreeClassifier(random_state=52)
31 # Train the model on the training data
32 dtc_norm.fit(X_train_ros, y_train_ros)
33 # Make predictions on the test data
34 y_pred_norm_ros = dtc_norm.predict(X_test_norm)
35
36 print("Decision Tree for data with normalization and random oversampling")
37 # Print Accuracy
38 printAcc(y_test_norm,y_pred_norm_ros)
39
40 # Print Confusion Matrix
41 print("")
42 print("Confustion Matrix")
43 printConfusion(y_test_norm, y_pred_norm_ros)
44
45 # Print Diagnosis
```

---

```
46 print("")
47 printReport(y_test_norm,y_pred_norm_ros)
```

```
Decision Tree for data with normalization and random oversampling
Accuracy: 0.8037022997220116

Confustion Matrix
[[11656  1570]
 [ 1537 1065]]
TP: 1065 , FP: 1570 , TN: 11656 , FN: 1537
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.88 | 0.88 | 0.88 | 13226 |
| 1 | 0.40 | 0.41 | 0.41 | 2602 |
| accuracy |  |  | 0.80 | 15828 |
| macro avg | 0.64 | 0.65 | 0.64 | 15828 |
| weighted avg | 0.80 | 0.80 | 0.80 | 15828 |

Double-click (or enter) to edit

```
1  # Cross validated decision tree normalized no over sampling
2  from sklearn.model_selection import cross_val_score, cross_val_predict, KFold
3  # StratifiedKFold
4  from sklearn.tree import DecisionTreeClassifier
5  from imblearn.pipeline import Pipeline
6  from sklearn.metrics import classification_report
7
8  # Create a deep copy of the data frame adult_cleaned and name it adult_c_dt_norm_ros
9  # X is normalized and random oversampled
10 adult_c_dt = adult_cleaned.copy(deep=True)
11 X = adult_c_dt.drop('target', axis=1)
12 y = adult_c_dt['target']
13
14 # Normalization the numerical columns
15 num_cols = ['age','fnlwgt','education-num','capital-gain','capital-loss','hours-per-week']
16 X_num = X[num_cols]
17 X_num_normalized = X_num.apply(normalize, axis=0)
18 # combine the normalized numerical columns with the categorical columns
19 X = pd.concat([X_num_normalized, X.drop(num_cols, axis=1)], axis=1)
20
21 # dummy variable
22 X=dummy(X)
23
24 # Create a Decision Tree classifier
25 c_dt= DecisionTreeClassifier()
26
27 # Create a pipeline with oversampling and the decision tree classifier
28 model = Pipeline([('ros', ros), ('dt', c_dt_rom)])
29
30 # Set up cross-validation using StratifiedKFold
31 cv = KFold(n_splits=10, shuffle=True, random_state=52)
32
33 # Perform cross-validation and obtain predicted labels
34 y_pred = cross_val_predict(c_dt,X,y, cv=cv)
35
36 # Calculate and print classification report
37 print("Classification Report:\n", classification_report(y, y_pred))
38
39 # Perform cross-validation
40 scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv)
41
42 # Print the mean accuracy across all folds
43 print("Mean Accuracy:", scores.mean())
44
```

```
Classification Report:
              precision    recall  f1-score   support

           0       0.89      0.88      0.88     37714
           1       0.41      0.43      0.42      7508

    accuracy                           0.80     45222
   macro avg       0.65      0.65      0.65     45222
weighted avg       0.81      0.80      0.81     45222

Mean Accuracy: 0.8046967325976435
```

```
1  # cross validated decision tree normalized random oversampling
2  # I use pipeline!
3  from sklearn.model_selection import cross_val_score, cross_val_predict, KFold
4  # StratifiedKFold
5  from sklearn.tree import DecisionTreeClassifier
6  from imblearn.over_sampling import RandomOverSampler
7  from imblearn.pipeline import Pipeline
8  from sklearn.metrics import classification_report
9  import pandas as pd
10
11 # Create a deep copy of the data frame adult_cleaned and name it adult_c_dt_norm_ros
12 # X is normalized and random oversampled
13 adult_c_dt_norm_ros = adult_cleaned.copy(deep=True)
14 X = adult_c_dt_norm_ros.drop('target', axis=1)
15 y = adult_c_dt_norm_ros['target']
16
17 # Normalization the numerical columns
18 num_cols = ['age','fnlwgt','education-num','capital-gain','capital-loss','hours-per-week']
19 X_num = X[num_cols]
20 X_num_normalized = X_num.apply(normalize, axis=0)
21 # combine the normalized numerical columns with the categorical columns
22 X = pd.concat([X_num_normalized, X.drop(num_cols, axis=1)], axis=1)
23
24 # dummy variable
25 X=dummy(X)
26
27 # Create a RandomOverSampler
28 ros = RandomOverSampler()
29
30 # Create a Decision Tree Classifier
31 c_dt_rom = DecisionTreeClassifier()
32
33 # Apply oversampling to X and y
34 X_resampled, y_resampled = ros.fit_resample(X, y)
```

```
35
36
37  # Create a pipeline with oversampling and the decision tree classifier
38  model = Pipeline([('ros', ros), ('dt', c_dt_rom)])
39
40  # Set up cross-validation using StratifiedKFold
41  cv = KFold(n_splits=10, shuffle=True, random_state=52)
42
43  # Perform cross-validation and obtain predicted labels
44  y_pred = cross_val_predict(model, X_resampled, y_resampled, cv=cv)
45
46  # Calculate and print classification report
47  print("Classification Report:\n", classification_report(y_resampled, y_pred))
48
49  # Perform cross-validation
50  scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv)
51
52  # Print the mean accuracy across all folds
53  print("Mean Accuracy:", scores.mean())
54
55  #################### comments: how to print out all those stuff?
```

```
Classification Report:
              precision    recall  f1-score   support

           0       1.00      0.88      0.93     37714
           1       0.89      1.00      0.94     37714

    accuracy                           0.94     75428
   macro avg       0.94      0.94      0.94     75428
weighted avg       0.94      0.94      0.94     75428

Mean Accuracy: 0.80710712156444292
```

```
1  ####################################
2  # Logistic regression
3  ####################################
4
5  # Logistic regression with cross validation and no oversampling
6  from sklearn.linear_model import LogisticRegression
7  from sklearn.model_selection import cross_val_score, KFold
8
9  # Create a deep copy of the data frame adult_cleaned and name it adult_c_dt_norm_ros
10 # X is normalized and random oversampled
11 adult_c_log = adult_cleaned.copy(deep=True)
12 X = adult_c_log.drop('target', axis=1)
13 y = adult_c_log['target']
14
15 # Normalization the numerical columns
16 num_cols = ['age', 'fnlwgt', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']
17 X_num = X[num_cols]
18 X_num_normalized = X_num.apply(normalize, axis=0)
19 # combine the normalized numerical columns with the categorical columns
20 X = pd.concat([X_num_normalized, X.drop(num_cols, axis=1)], axis=1)
21
22 # dummy variable
23 X=dummy(X)
24
25 # create a logistic model
26 model = LogisticRegression(max_iter=1000)   # Increase max_iter if needed for convergence
27
28 # Set up 10-fold cross-validation
29 kfold = KFold(n_splits=10, shuffle=True, random_state=52)
30
31 # Perform cross-validation and get the accuracy scores for each fold
32 scores = cross_val_score(model, X, y, cv=kfold)
33
34 # Print the accuracy for each fold and the mean accuracy
35 for i, score in enumerate(scores, 1):
36     print(f'Fold {i}: {score}')
37
38 print(f'Mean Accuracy: {scores.mean()}')
```

```
Fold 1: 0.8425823568428035
Fold 2: 0.8534158744419633
Fold 3: 0.8423264024459089
Fold 4: 0.8491817779743477
Fold 5: 0.8524988942945599
Fold 6: 0.8429898275099513
Fold 7: 0.8516143299425033
```

```
Fold 8: 0.8398938522777533
Fold 9: 0.8403361344537815
Fold 10: 0.8423264042459089
Mean Accuracy: 0.8457165856207152
```

```
 1 # Logistic regression with cross validation and random oversampling
 2
 3 from sklearn.linear_model import LogisticRegression
 4 from sklearn.model_selection import cross_val_score, KFold
 5 from imblearn.over_sampling import RandomOverSampler
 6
 7 # Create a deep copy of the data frame adult_cleaned and name it adult_c_dt_norm_ros
 8 # X is normalized and random oversampled
 9 adult_c_log_ros = adult_cleaned.copy(deep=True)
10 X = adult_c_log_ros.drop('target', axis=1)
11 y = adult_c_log_ros['target']
12
13 # Normalization the numerical columns

 1 ################################################################
 2 # KNN
 3 ################################################################
 4 # KNN with cross validation no oversampling
 5 from sklearn.neighbors import KNeighborsClassifier
 6 from sklearn.model_selection import cross_val_score, KFold
 7
 8 # Create a deep copy of the data frame adult_cleaned and name it adult_c_knn
 9 adult_c_knn = adult_cleaned.copy(deep=True)
10 X = adult_c_knn.drop('target', axis=1)
11 y = adult_c_knn['target']
12
13 # Normalization the numerical columns
14 num_cols = ['age','fnlwgt','education-num','capital-gain','capital-loss','hours-per-week']
15 X_num = X[num_cols]
16 X_num_normalized = X_num.apply(normalize, axis=0)
17 # combine the normalized numerical columns with the categorical columns
18 X = pd.concat([X_num_normalized, X.drop(num_cols, axis=1)], axis=1)
19
20 # dummy
21 X=dummy(X)
22
```