**question2.py**

```python
import numpy as np
import pandas as pd
import math
from sklearn.tree import DecisionTreeClassifier
from sklearn.preprocessing import OneHotEncoder
from sklearn.metrics import accuracy_score

# The CSV data has no headers.  These headers will be added after import.
column_headers = [
    "age",
    "class_of_worker",
    "detailed_industry_recode",
    "detailed_occupation_recode",
    "education",
    "wage_per_hour",
    "enroll_in_edu_inst_last_wk",
    "marital_stat",
    "major_industry_code",
    "major_occupation_code",
    "race",
    "hispanic_origin",
    "sex",
    "member_of_labor_union",
    "reason_for_unemployment",
    "full_or_part_time_employment_stat",
    "capital_gains",
    "capital_losses",
    "dividends_from_stocks",
    "tax_filer_stat",
    "region_of_previous_residence",
    "state_of_previous_residence",
    "detailed_household_and_family_stat",
    "detailed_household_summary_in_household",
    "unknown_value",
    "migration_code_change_in_msa",
    "migration_code_change_in_reg",
    "migration_code_move_within_reg",
    "live_in_this_house_1_year_ago",
    "migration_prev_res_in_sunbelt",
    "num_persons_worked_for_employer",
    "family_members_under_18",
    "country_of_birth_father",
    "country_of_birth_mother",
    "country_of_birth_self",
    "citizenship",
    "own_business_or_self_employed",
    "fill_inc_questionnaire_for_veterans_admin",
```

```
48          "veterans_benefits",
49          "weeks_worked",
50          "year",
51          "income_50k_plus"
52  ]
53
54  # The names of the files to import.
55  datafiles = ["data/census-income.data.csv", "data/census-income.test.csv"]
56
57  # Initialize an array to hold dataframes for the data and test sets.
58  dataframes = []
59
60  # Load the CSV files.
61  for f in datafiles:
62      # Read the file into a dataframe.  As noted, there are no headers.
63      df = pd.read_csv(f, header=None)
64      # Add the column headers.
65      df.columns = column_headers
66
67      # Data cleaning operations:
68
69      # Handle missing values using forward fill and backward fill.
70      df = df.ffill().bfill()
71
72      # Upon analyzying the data, the two values in the "income_50k_plus" column, after
73      # import, look like this: [' - 50000.', ' 50000+.'].  I am replacing these values using
74      # one-hot encoding so that less than 50k income is a 0, and 50k+ is a 1.
75
76      # Strip whitespace
77      df["income_50k_plus"] = df["income_50k_plus"].str.strip()
78      # Map the existing values to 0 or 1.
79      df["income_50k_plus"] = df["income_50k_plus"].map({"- 50000.": 0, "50000+.": 1})
80
81      # Upon analyzing the data, the the values in the "race" column contain leading
82      # and/or trailing spaces.  These will be removed.
83
84      # Strip whitespace
85      df["race"] = df["race"].str.strip()
86
87      # Add this dataframe to the array
88      dataframes.append(df)
89
90  # Assign the data frames as training or test data.
91  training_data = dataframes[0]
92  test_data = dataframes[1]
93
94  # Define features and target variable
95  target_col = "income_50k_plus"
96  feature_cols = [col for col in training_data.columns if col != target_col]
```

```python
 97
 98  # Identify categorical and numerical features
 99  categorical_cols = training_data.select_dtypes(include=["object"]).columns.tolist()
100  numerical_cols = training_data.select_dtypes(include=["int64", "float64"]).columns.tolist()
101
102  # Create a one-hot encoder for categorical features.
103  # Note: The income_50k_plus column was encoded earlier, becaue that code was carried
104  # over from Assignment 0.
105  encoder = OneHotEncoder(handle_unknown="ignore", sparse_output=False)
106
107  # Encode the categorical columns from the test data.
108  encoded_train = encoder.fit_transform(training_data[categorical_cols])
109
110  # Encode the categorical columns from the test data.
111  encoded_test = encoder.transform(test_data[categorical_cols])
112
113  # Convert the encoded columns to dataframes.
114  encoded_train_df = pd.DataFrame(encoded_train, columns=encoder.get_feature_names_ou↵
     t(categorical_cols))
115  encoded_test_df = pd.DataFrame(encoded_test, columns=encoder.get_feature_names_ou↵
     t(categorical_cols))
116
117  # Assemble the prepared training and test data from the numerical features and the one-hot
     encoded
118  # categorical features.  Remove the target column from numerical features before during this
     process.
119  X_train =
     pd.concat([training_data[numerical_cols].drop(columns=[target_col]).reset_index(drop=True),
     encoded_train_df], axis=1)
120  X_test =
     pd.concat([test_data[numerical_cols].drop(columns=[target_col]).reset_index(drop=True),
     encoded_test_df], axis=1)
121
122  # Error checking: The code will exit if the target column is in the training or test data.
123  assert target_col not in X_train.columns, "Target column should not be in X_train!"
124  assert target_col not in X_test.columns, "Target column should not be in X_test!"
125
126  # Get the target column from the original training and test datasets.
127  y_train = training_data[target_col]
128  y_test = test_data[target_col]
129
130  # Part A: Train decision trees for depths of 2 - 10 and store accuracies.
131  depth_accuracies = []
132  for depth in range(2, 11):
133      # The random_state value is arbitrary, as long as the same value is used
134      # each time to allow reproducibility.
135      model = DecisionTreeClassifier(max_depth=depth, random_state=42)
136      model.fit(X_train, y_train)
137      train_acc = accuracy_score(y_train, model.predict(X_train))
138      depth_accuracies.append({"Depth": depth, "Training Accuracy": train_acc})
```

```
139
140  print(depth_accuracies)
141  # Find the maximum accuracy and its associated depth.
142  max_accuracy = max([entry["Training Accuracy"] for entry in depth_accuracies])
143  # The optimal depth is the depth associated with max_accuracy.
144  optimal_depth = next(item["Depth"] for item in depth_accuracies if item["Training Accuracy"]
     == max_accuracy)
145  print(f"max_accuracy = {max_accuracy}")
146  print(f"optimal_depth = {optimal_depth}")
147
148  # Part B: use the optimal depth found above to classify the test data.
149  # Re-train the model using the optimal depth found above.
150  model = DecisionTreeClassifier(max_depth=optimal_depth, random_state=42)
151  model.fit(X_train, y_train)
152  # Get the accuracy score for the training data (should be the same as before).
153  train_acc = accuracy_score(y_train, model.predict(X_train))
154  # Now get the accuracy score for the test data.
155  test_acc = accuracy_score(y_test, model.predict(X_test))
156  print(f"Training accuracy for depth=10: {train_acc}")
157  print(f"Test accuracy for depth=10: {test_acc}")
158
159
```