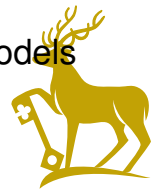


## Worksheet 4 (Week 4)

MATM063: Preparing data for ML & Regression models

Werner Bauer

Department of Mathematics



UNIVERSITY OF  
SURREY

**Key Learning:** Cleaning and scaling data, converting categories, preprocessing pipeline, regression models

### 1 USUAL WORK FLOW IN DATA SCIENCE (PART 2)

Here, in short, part 2 of the steps to be done to prepare your data for machine learning algorithms.

1. Split training data into predictors and labels
2. Clean data: remove missing values or fill in values (e.g. SimpleImputer)
3. Convert categorical to numerical values (e.g. OneHotEncoder)
4. Scale features (e.g. MinMaxScaler, StandardScaler)
5. Combine these steps into a transformer pipeline
6. Train and evaluate your ML model

### 2 REGRESSION MODELS

For the provided dataset, we want to use a regression model. Besides those regression models discussed in the lectures, another one is included that you will need further below in Q2.

- Linear Regression:

```
1 from sklearn.linear_model import LinearRegression
2
3 lin_reg = LinearRegression() # create instance
4 lin_reg.fit(housing_prepared, house_labels) # train the model
```

- Decision Tree Regressor:

```
1 from sklearn.tree import DecisionTreeRegressor
2 tree_reg = DecisionTreeRegressor()
3 tree_reg.fit(housing_prepared, housing_labels)
```

- In Q2, you should train a Random Forest Regressor:

```
1 from sklearn.ensemble import RandomForestRegressor
2
3 forest_reg = RandomForestRegressor()
4 forest_reg.fit(housing_prepared, housing_labels)
```

### 3 EVALUATION OF THE MODELS

Here are two ways of evaluating the performance of regression model:

- Root Mean Square Error (RMSE):

```
1 from sklearn.metrics import mean_squared_error
2 housing_predictions = lin_reg.predict(housing_prepared)
3 lin_rmse = np.sqrt(mean_squared_error(housing_labels, housing_predictions))
4 print(lin_rmse)
```

- $k$ -fold Cross Validation (CV):

```
1 from sklearn.model_selection import cross_val_score
2
3 lin_scores = cross_val_score(lin_reg, housing_prepared, housing_labels,
```

```

4         scoring="neg_mean_squared_error", cv=10)
5 lin_rmse_scores = np.sqrt(-lin_scores)
6
7 # Use output function:
8 def display_scores(scores):
9     print ("Scores:", scores)
10    print ("Mean:", scores.mean())
11    print ("Standard deviation:", scores.std())
12 display_scores(lin_rmse_scores)

```

## 4 CUSTOMER TRANSFORMER TO ADD COMBINED ATTRIBUTES

The following customer transformer is used below in the transformer pipeline.

```

1 rooms_ix, bedrooms_ix, population_ix, households_ix = 3, 4, 5, 6
2 class CombinedAttributesAdder (BaseEstimator, TransformerMixin):
3     def __init__ (self, add_bedrooms_per_room=True):
4         self.add_bedrooms_per_room =
5         add_bedrooms_per_room
6     def fit(self, X, y=None):
7         return self # nothing else to do
8     def transform(self, X):
9         rooms_per_household = X[:, rooms_ix] /
10                                X[:, households_ix]
11         population_per_household = X[:, population_ix] /
12                                    X[:, households_ix]
13         if self.add_bedrooms_per_room:
14             bedrooms_per_room = X[:, bedrooms_ix] /
15                                 X[:, rooms_ix]
16             return np.c_[X, rooms_per_household,
17                           population_per_household,
18                           bedrooms_per_room]
19         else:
20             return np.c_[X, rooms_per_household,
21                           population_per_household]

```

Listing 1. combined attribute adder

## 5 DATA PREPROCESSING CLASSES AND PIPELINES

Here is a summary of the packages we can use to do data preprocessing.

- Data cleaning:

```

1 from sklearn.impute import SimpleImputer
2 imputer = SimpleImputer(strategy='median') # create instance and define strategy
3 X = imputer.fit_transform(housing_num) #numpy array
4 housing_tr = pd.DataFrame(X, columns=housing_num.columns, index=housing_num.index) #DataFrame

```

- Converting text (categories) to numbers:

```

1 # define DataFrame with categories
2 housing_cat = housing[['ocean_proximity']]
3
4 from sklearn.preprocessing import OrdinalEncoder
5 ordinal_encoder = OrdinalEncoder()
6 housing_cat_encoded = ordinal_encoder.fit_transform(housing_cat)
7
8 # Output categories with:
9 cat_encoder.categories_
10 # Output:
11 [array(['<1H OCEAN', 'INLAND', 'ISLAND', 'NEAR BAY', 'NEAR OCEAN'], dtype=object)]
12
13 # similarly for OneHotEncoder:
14 from sklearn.preprocessing import OneHotEncoder
15 cat_encoder = OneHotEncoder()
16 housing_cat_1hot = cat_encoder.fit_transform(housing_cat)

```

- Feature scaling:

```

1 # normalising data:
2 from sklearn.preprocessing import MinMaxScaler
3 minmax_scale = MinMaxScaler()
4 housing_minmax = minmax_scale.fit_transform(housing_num)
5
6 # standardising data:
7 from sklearn.preprocessing import StandardScaler
8 std_scale = StandardScaler()
9 housing_std = std_scale.fit_transform(housing_num)

```

- Setting up a transformation pipeline:

```

1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3
4 # define numerical and categorical attribute lists
5 num_attribs = list(housing_num)
6 cat_attribs = ['ocean_proximity']
7
8 num_pipeline = Pipeline([
9     ('imputer', SimpleImputer(strategy='median')),
10    ('attribs_adder', CombinedAttributesAdder(add_bedrooms_per_room=True)),
11    ('std_scaler', StandardScaler()), ])
12
13 full_pipeline = ColumnTransformer([('num', num_pipeline, num_attribs),
14    ('cat', OneHotEncoder(), cat_attribs), ])

```

- Automatisation: easily transform data by applying `full_pipeline` directly to original data set `housing`.

```

1 housing_prepared = full_pipeline.fit_transform(housing)

```

## 6 DEFINING CELLS IN SPYDER

Sometimes it is useful not to execute the entire Python script. This could be the case when the training of a model takes longer (e.g. you will see that the Random Forest Regressor takes some time to be trained). Then, you can split your code into cells that can be executed separately (e.g. with CTRL-Return shortcut you can execute a selected cell). To begin a new cell in your code, simply use the following command:

```

1 ### HERE YOU CAN ADD SOME COMMENT
2
3 import ...
4 ...

```

**Listing 2.** organizing code parts into cells in Spyder

## 7 EXERCISES

### Question 1:

Here, you should combine the material we discussed in the lectures into **one running Python script** that reads in the housing file, creates a test and training set, transforms the data such that it can be used in a machine learning algorithm, and then trains a linear regression model.

Please proceed as follows to create such script:

1. Start your script with the part that reads in the housing data file from the internet <https://raw.githubusercontent.com/ageron/handson-ml2/master/datasets/housing/housing.tgz>.
2. Load the `housing.csv` file using the `fetch_housing_data()` function.
3. Split your data into test and training data by using Stratified Shuffle Split (20% test data) while assuring that `median_income` is proportionally represented in the test data.
4. Separate the target `median_house_value` from the predictors (similarly to `housing` and `housing_labels` as done in the lectures). Then, define numerical attributes and categorical attributes (in the lectures we called them `num_attribs` and `cat_attribs`).
5. Write a full pipeline that fits and transforms `housing` to `housing_prepared` that can be used to train ML models. Make sure that you can change the option `add_bedrooms_per_room` from `True` to `False` as we will need this at the next step.
6. Define and train a linear regression model (as done in the lectures). Output both, the model's RMSE for the training dataset and also the cross validation scores (10-folds). Compare two realization of these models, a) one where you set `add_bedrooms_per_room = True` (you should have very similar results as in the lecture notes); and b) one with `add_bedrooms_per_room = False`

7. Test also the performance of these two models (for the default customer transformer setting) when comparing min-max scaling vs standardisation.
8. Finally test (for the default customer transformer setting and standardisation) what happens when you use the Ordinal Encoder rather than the OneHotEncoder.

### **Question 2:**

This question requires you that you have successfully finished Q1.

1. Train a Random Forest Regressor on the prepared housing dataset (Steps 1-5 of Q1).
2. Make a 'no-CV' prediction (no cross-validation) with this model for the entire training set and output the RMSE.
3. Now, make a prediction using a 5-fold cross-validation and output: Scores, Mean, Standard deviation.
4. Compare the RMSE of the no-CV with the CV predictions, and deduce from these values if your model has underfit or overfit the training data?

### **Question 3:**

This question is related to the code you did produce for Question 1 and Question 2.

1. Update the customer transformer (Listing 1) such that it has two hyper-parameters: `add_bedrooms_per_room = True` as well as `add_rooms_per_household = True`.
2. Update the pipeline from Q1 such that it uses this new transformer.
3. Try out different combinations of adding or omitting these combined attributes to explore if their presents might improve the performance of the linear regression model.

### **Question 4:**

For this question, you have to write a Python script that makes predictions for the target value `median_income` (Note: in Q1 the target value is `median_house_values`). The values in `median_income` are in units of 10.000\$. (Hint: rather than writing the whole script from scratch, you could also copy the one from Q1 and change the corresponding entries).

1. Load the housing file.
2. Split your data into test and training data using Stratified Shuffle Split (15% test data) while assuring that `median_house_values` is proportionally represented in the test data.
3. Use the training set to proceed and separate `median_income` from the predictors.
4. Write a transformer pipeline (as in Q1) that summarizes all the necessary data preprocessing steps into one function (use default settings)!
5. Train a Linear Regression, a Decision Tree Regressor, and a Random Forest Regressor, evaluate their RMSEs on the training set and when applying cross-validation. Are the models capable of making reasonable predictions of `median_income`?