# Data Science: Predictive Analytics, Übung 3

Henry Haustein

## Modeling & Evaluation

```python
1  import pandas as pd
2  from sklearn.linear_model import LinearRegression
3  from sklearn.tree import DecisionTreeRegressor
4  from sklearn.model_selection import train_test_split
5  from keras.models import Sequential
6  from keras.layers import Dense
7  import numpy as np
8
9  data = pd.read_csv("data_ex3.csv", index_col = 0)
10 print(data.info())
11 print(data.describe())
12
13 # Task 1
14 X = data.drop(columns = ["Rented_Bike_Count", "Date"])
15 y = data["Rented_Bike_Count"]
16 x_train, x_test, y_train, y_test = train_test_split(X, y,
       test_size = 0.3)
17 model = LinearRegression()
18 model.fit(x_train, y_train)
19 print(model.score(x_test, y_test)) # 0.55 for full model
20
21 # using R's step() function, I get the following better model
22 X2 = data.drop(columns = ["Rented_Bike_Count", "Date", "
       Functioning_Day", "Month"])
23 y2 = data["Rented_Bike_Count"]
24 x2_train, x2_test, y2_train, y2_test = train_test_split(X2, y2,
       test_size = 0.3)
25 model_reduced = LinearRegression()
26 model_reduced.fit(x2_train, y2_train)
27 print(model_reduced.score(x2_test, y2_test)) # 0.47 for better (?)
       model
28
29 model2 = DecisionTreeRegressor()
30 model2.fit(x_train, y_train)
31 print(model2.score(x_test, y_test)) # 0.04 for full model
32
33 model2_reduced = DecisionTreeRegressor()
34 model2_reduced.fit(x2_train, y2_train)
```

```
35  print(model2_reduced.score(x2_test, y2_test)) # 0.02 for better
        (?) model
36
37  newData = pd.DataFrame({"Hour": 9, "Temperature": -2, "Humidity":
        42, "Wind_speed": 6, "Visibility": 2000, "Solar_Radiation": 0,
         "Rainfall": 0, "Snowfall": 0, "Seasons": 1, "Holiday": 0, "
        Functioning_Day": True, "Weekday": 1, "Day_Night": 0, "Month":
         3}, index = [0])
38  print(model.predict(newData)) # 453.36
39  print(model2.predict(newData)) # 206 (but huge spread depending on
        training data)
40
41  # Task 2
42  print(model.coef_[1]) # Temperature coefficient = 28.25
43  print(model.coef_[6]) # Rainfall ciefficient = -56.82
44  print(model_reduced.coef_[1]) # Temperature coefficient = 29.19
45  print(model_reduced.coef_[6]) # Rainfall ciefficient = -63.89
46
47  # Task 3
48  data["Functioning_Day"] = np.asarray(data["Functioning_Day"]).
        astype(np.float32) # boolean values don't work
49  X3 = data.drop(columns = ["Date"]).copy()
50  y3 = X3.pop("Rented_Bike_Count")
51  X3 = np.array(X3)
52  y3 = np.array(y3)
53
54  x3_train, x3_test, y3_train, y3_test = train_test_split(X3, y3,
        test_size = 0.3)
55
56  model3 = Sequential()
57  model3.add(Dense(32, input_dim = 14))
58  model3.add(Dense(128))
59  model3.add(Dense(64))
60  model3.add(Dense(64))
61  model3.add(Dense(64))
62  model3.add(Dense(32))
63  model3.add(Dense(1))
64  model3.compile(
65    optimizer = 'adam',
66    loss = 'mean_absolute_error',
67    metrics = ['mean_absolute_error']
68  )
69  model3.fit(x3_train, y3_train)
70  print(model3.evaluate(x3_test, y3_test)) # [404.8241882324219,
        404.8241882324219]
71  model3.fit(x3_train, y3_train, epochs = 10)
72  print(model3.evaluate(x3_test, y3_test)) # [386.04815673828125,
        386.04815673828125]
73
74  newData = pd.DataFrame({"Hour": 9, "Temperature": -2, "Humidity":
        42, "Wind_speed": 6, "Visibility": 2000, "Solar_Radiation": 0,
```

```
     "Rainfall": 0, "Snowfall": 0, "Seasons": 1, "Holiday": 0, "
     Functioning_Day": 0, "Weekday": 1, "Day_Night": 0, "Month":
     3}, index = [0])
75 print(model3.predict(newData)) # 197.68306 (but huge spread
     depending on training data)
```