

Stanford University

Final Project Report



Stanford
University

Thomas Sounack

Professor: Eric Darve

ME 343

Contents

| | | |
|---|--|---|
| 1 | Description of the project | 2 |
| 2 | Technical problem of interest and its difficulties | 2 |
| 3 | Algorithm description | 2 |
| 4 | Published literature and available approaches | 3 |
| 5 | Algorithm implementation | 3 |
| 6 | Algorithm verification | 3 |
| 7 | Datasets, problem settings and parameters | 4 |
| 8 | Results | 4 |
| 9 | Conclusion and limitations | 5 |

1 Description of the project

The goal of this project is to design a model that can predict wind speeds at different altitudes in the context of implementing offshore wind turbines. When building offshore wind turbines, on-site studies have to be made during several months to several years in order to determine which specific spots would provide the best wind speeds for the turbines.

Currently, lidar buoys are used to estimate wind speeds at various altitudes. This solution is rather costly (several thousand dollars per week for one buoy), which is why the objective is to replace the lidar buoys with a model that would compute wind speeds at various altitudes using simple sensors at sea level and satellite data. This is possible as several physical equations link conditions at sea level with wind speeds (such as the difference in temperature between the sea and the air, the rugosity - charnock - of the sea, etc...).

Using data from a lidar buoy and the ERA5 satellite, the project will consist in designing a multi-output regression machine learning model that can predict wind speeds at different heights.

This project is the continuation of an internship I did last summer in an offshore wind company. During this internship, I used XGBoost to accomplish a similar task without a good understanding of ML models and hyperparameter optimization. Through this project, I hope to be able to obtain similar or better performance with a model I built and optimized by myself, using the content I learned during this course.

2 Technical problem of interest and its difficulties

The technical problem of interest in this project is to design a machine learning model that can predict wind speeds at different altitudes in the context of implementing offshore wind turbines. The main difficulty lies in the fact that there are several physical equations that link conditions at sea level with wind speeds, and it can be challenging to accurately capture all the relevant variables and their interactions in a model.

Furthermore, predicting wind speeds at multiple altitudes simultaneously is needed as offshore wind turbines can have very different sizes, which adds another layer of complexity to the problem. Therefore, a multi-output regression machine learning model will be used to tackle this problem. This approach allows for predicting wind speeds at multiple altitudes simultaneously, taking into account the interdependence between the different altitudes.

3 Algorithm description

Since this problem requires multi-output regression and is complex to model, designing a deep learning architecture that has native support for multiple outputs appears to be an interesting approach, and has the added benefit of applying concepts taught in class during the quarter.

The model therefore relies on a sequential neural network whose structure is detailed in section 5. We arbitrarily set the model as site-dependent to limit the complexity of the problem.

4 Published literature and available approaches

During my internship, I used three papers to understand the problem at play and have a good understanding of the physical equations behind it:

- A paper on Offshore vertical wind shear by the DTU (1)
- A paper on Southern New England's vertical wind shear conditions (2)
- A study on the northern french coast wind conditions (3)

In addition to these resources, two articles from the website *Machine Learning Mastery* -(4) and (5)- have helped me understand how to perform multi output regression with pytorch.

5 Algorithm implementation

To implement this algorithm, we use a neural network model which consists of two {linear, activation, batch normalization} blocks and a block composed only of an activation layer.

To identify the best parameters we perform a random search, as we have seen in lectures that it is more effective than a grid search. The parameters that are being searched are the activation type, the batch size, the learning rate and the number of neurons in the first and second layer (the others are imposed by the input and output size). The results are as follows:

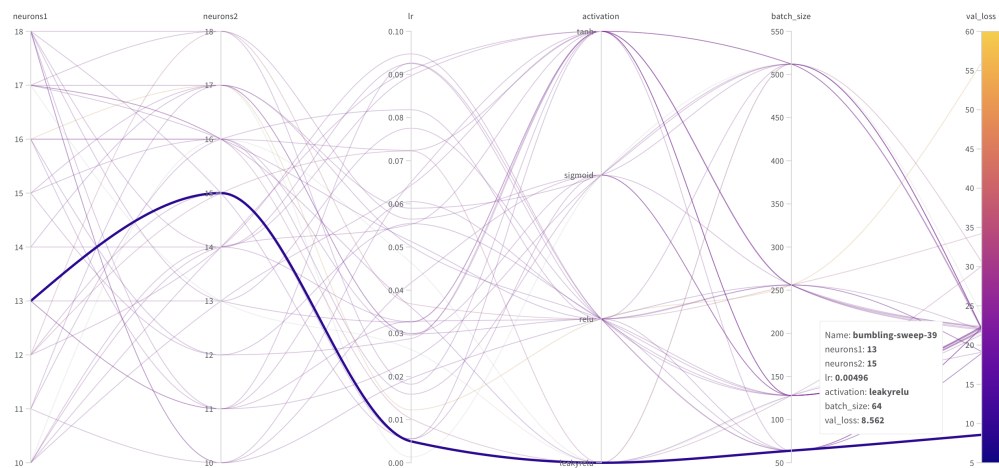


Figure 1: Random search over the parameter space, obtained through wandb

This allows us to identify the best configuration for our network, which will be used to evaluate the model. The code to achieve this is adapted from homework 6, and can be found in section 9.

6 Algorithm verification

To ensure that the code is working correctly, several steps have been implemented:

- The data has been imported in an auxiliary python script and tested before moving on to implementing the model. This ensures that the model receives correct data.

- Wandb is used to ensure that the model is learning properly, providing a feedback on the validation and training loss at each step.
- Finally, the model’s performance is tested after the learning process (see section 8).

7 Datasets, problem settings and parameters

The dataset used comes from several sources: the buoy data (lidar + anemometer) is the one I used during my internship and belongs to the company, while the satellite data (ERA5) is open-source and available online at the desired location. Both the location and the buoy data cannot be shared in compliance with the confidentiality agreement I signed.

To ensure that the algorithm is working properly, a good sanity check is to plot the correlation matrix between the inputs and the outputs of the model.

| | WS | Wdir | Delta T | p140209 | msl | t2m | d2m | blh | sst | chnk | ssr | month |
|---------|----------|----------|-----------|-----------|-----------|-----------|-----------|----------|-----------|----------|-----------|-----------|
| WS 40m | 0.962580 | 0.056202 | -0.037145 | 0.022448 | -0.206344 | -0.172457 | -0.161730 | 0.410736 | -0.189806 | 0.530855 | -0.108940 | -0.051696 |
| WS 57m | 0.951707 | 0.054832 | -0.020029 | 0.014973 | -0.209202 | -0.166127 | -0.154084 | 0.406570 | -0.192516 | 0.530396 | -0.108120 | -0.060123 |
| WS 77m | 0.940810 | 0.054939 | -0.006381 | 0.008025 | -0.211938 | -0.160547 | -0.147685 | 0.404326 | -0.195268 | 0.530346 | -0.106775 | -0.066966 |
| WS 97m | 0.932843 | 0.055677 | 0.001790 | 0.002735 | -0.214116 | -0.156652 | -0.143227 | 0.404060 | -0.197461 | 0.531089 | -0.105499 | -0.071202 |
| WS 117m | 0.926523 | 0.056475 | 0.008007 | -0.001786 | -0.215435 | -0.153157 | -0.139339 | 0.404111 | -0.198920 | 0.532066 | -0.104665 | -0.073964 |
| WS 137m | 0.920922 | 0.056712 | 0.013130 | -0.005586 | -0.216315 | -0.150151 | -0.135964 | 0.403874 | -0.200065 | 0.532769 | -0.103706 | -0.075780 |
| WS 157m | 0.916063 | 0.056748 | 0.017214 | -0.008947 | -0.217058 | -0.147485 | -0.133095 | 0.404001 | -0.200589 | 0.533665 | -0.103337 | -0.076571 |
| WS 177m | 0.911754 | 0.057688 | 0.020273 | -0.012046 | -0.218251 | -0.145436 | -0.130836 | 0.405020 | -0.200790 | 0.535067 | -0.103945 | -0.076427 |
| WS 197m | 0.907898 | 0.058215 | 0.022791 | -0.014380 | -0.219248 | -0.144229 | -0.129290 | 0.405963 | -0.201181 | 0.536236 | -0.104781 | -0.076545 |

Figure 2: *Input - Output correlation matrix*

As expected, the wind speeds are highly correlated with WS (the wind speed at sea level). Other parameters like the charnok (rugosity of the sea) are also important.

We can note that (as expected) the correlation is different for the various heights. Indeed, it makes sense that trying to infer wind speeds using data at sea level would produce different results based on the desired height. This observation implies that the benchmark needs to evaluate the different heights separately.

To perform this benchmark on the validation set, we can use correlation graphs at each height along with an R^2 value for the $y = x$ line. This is mostly a visual feedback, but is useful to see if our model has a strong tendency to overestimate or underestimate wind speeds (which was a recurrent problem during my internship because of outliers). For more quantitative data, we can use a table that displays at each height the R^2 value; the average error and the average absolute error to quantify the error in m/s ; the standard deviation to determine how spread-out the data is; and the average relative error to have a unitless metric of the model’s performance at this height.

8 Results

Using the benchmark described above, we get the following graph:

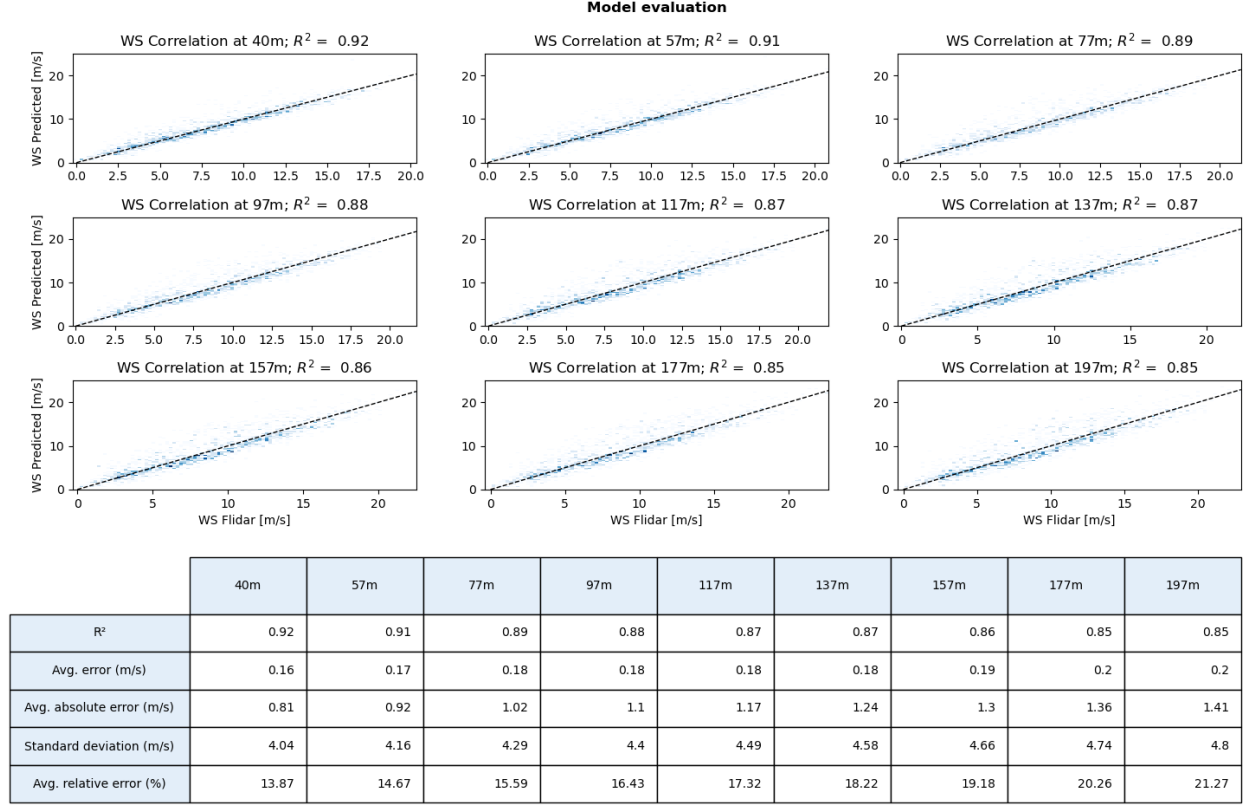


Figure 3: *Evaluating the performance of the model*

Visually, we can see that the data seems to match the correlation line well. This is confirmed by relatively high R^2 values and low error scores.

As expected, the accuracy decreases with the height. However, even at the maximal height of 197m, we have a relatively good score (21.27%) compared to the results obtained with XGBoost during my internship (29.3%). This means that the random search was effective in finding a good optimum for the hyperparameters of the neural network.

9 Conclusion and limitations

By using a random search to build our neural network, we were able to achieve satisfactory wind speed predictions. We can also note that the lower accuracy scores for the last heights is in part due to the lidar being less precise for higher measurements.

However, we can note that our predictions are still off by more than 20% for the higher wind speeds. In the context of predicting power outputs of potential wind turbines, this model could hardly replace lidar measurements. Since we have optimized the model, further improving the accuracy would require identifying and correcting the noise in our data (especially since we are using data from different sources). A good start would be to look into the precision of the lidar during rainy days (which is known to be lower than usual) as well as the precision of the ERA5 satellite.

References

- [1] A. Pena Diaz, T. Mikkelsen, S.-E. Gryning, C. Hasager, A. Hahmann, M. Badger, I. Karagali, and M. Courtney, *Offshore vertical wind shear: Final report on NORSEWIND's work task 3.1*, ser. DTU Wind Energy E. Denmark: DTU Wind Energy, 2012, no. 0005.
- [2] D. Borvarán, A. Peña, and R. Gandoin, "Characterization of offshore vertical wind shear conditions in southern new england," *Wind Energy*, vol. 24, no. 5, pp. 465–480, 2021. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/we.2583>
- [3] R. Husson, L. De Montera, H. Berger, P. Appelghem, "Etude de potentiel de vent à partir du sarao4 normandie 2021 -ddec 05," 2021.
- [4] J. Brownlee, "Deep learning models for multi-output regression," Available at <https://machinelearningmastery.com/deep-learning-models-for-multi-output-regression/> (2020/08/18).
- [5] M. A. I. Khan, "Multi-target predictions with multilinear regression in pytorch," Available at <https://machinelearningmastery.com/multi-target-predictions-with-multilinear-regression-in-pytorch/> (2022/12/14).

Appendix - Code

```
1 # %%
2 import torch
3 import torch.nn as nn
4 import torch.optim as optim
5 import pprint
6 import numpy as np
7 import os
8 import wandb
9 import pandas as pd
10 import matplotlib.pyplot as plt
11 import scipy.stats as st
12 from matplotlib.gridspec import GridSpec
13 device = torch.device("cpu")
14 from extract_data import *
15
16 # %%
17 X_data, y_data = create_dataframes()
18
19 # %%
20 df_corr = pd.DataFrame()
21 for column in y_data:
22     corrM = X_data.join(y_data[column]).corr()
23     corrM.drop(corrM.tail(1).index, inplace=True)
24     df_corr[column] = corrM[column]
25
26 df_corr = df_corr.T
27 df_corr
28
29 # %%
30 # Hyperparameters
31 num_epochs = 250
32 log_freq = 15
```

```

33 n_train = int(0.8 * len(X_data))
34
35 # Data loader
36 X_train = torch.from_numpy(X_data[:n_train].values).float().to(device)
37 y_train = torch.from_numpy(y_data[:n_train].values).float().to(device)
38 X_val = torch.from_numpy(X_data[n_train:].values).float().to(device)
39 y_val = torch.from_numpy(y_data[n_train:].values).float().to(device)
40 train_dataset = torch.utils.data.TensorDataset(X_train, y_train)
41 test_dataset = torch.utils.data.TensorDataset(X_val, y_val)
42
43 # %%
44 #wandb hyperparameter dictionary
45 sweep_configuration = {
46     "method": "random",
47     "name": "random_search",
48     "metric": {"goal": "minimize", "name": "val_loss"},
49     "parameters":
50     {
51         "lr": {"min": 0.0001, "max": 0.1},
52         "batch_size": {"values": [64, 128, 256, 512]},
53         "neurons1": {"values": [10, 11, 12, 13, 14, 15, 16, 17, 18]},
54         "neurons2": {"values": [10, 11, 12, 13, 14, 15, 16, 17, 18]},
55         "activation": {"values": ["tanh", "sigmoid", "relu", "leakyrelu"]}
56     },
57     "run_cap": 150
58 }
59 pprint.pprint(sweep_configuration)
60 project_name = "cme216_final_project"
61 group_name = "randomsearch"
62 sweep_id = wandb.sweep(sweep_configuration, project=project_name)
63
64 # %%
65 class Network(nn.Module):
66     def __init__(self, neurons1, neurons2, activation):
67         super().__init__()
68         self.network = nn.Sequential(
69             nn.Linear(12, neurons1),
70             activation,
71             nn.BatchNorm1d(neurons1),
72             nn.Linear(neurons1, neurons2),
73             activation,
74             nn.BatchNorm1d(neurons2),
75             nn.Linear(neurons2, 9)
76         )
77
78     def forward(self, x):
79         y_pred = self.network(x)
80         return y_pred
81
82 # %%
83 import time
84 t1 = time.time()
85 # =====
86 # Training
87 # =====
88 # Train the model
89 def train(config=None):
90     # Initialize the new wandb run
91     wandb.init(config=config, project=project_name, group=group_name)

```



```

92     config = wandb.config
93     train_loader = torch.utils.data.DataLoader(dataset=train_dataset,
94                                                batch_size=config.batch_size,
95                                                shuffle=True)
96     total_step = len(train_loader)
97     loss_list = []
98
99     activation = None
100     if config.activation == "tanh":
101         activation = nn.Tanh()
102     elif config.activation == "sigmoid":
103         activation = nn.Sigmoid()
104     elif config.activation == "relu":
105         activation = nn.ReLU()
106     else:
107         activation = nn.LeakyReLU()
108
109     model = Network(config.neurons1, config.neurons2, activation)
110     criterion = nn.MSELoss()
111     optimizer = torch.optim.Adam(model.parameters(), lr=config.lr)
112     for epoch in range(num_epochs):
113         for i, (train_x, train_y) in enumerate(train_loader):
114             # Run the forward pass
115             model.train()
116             output = model(train_x)
117             loss = criterion(output, train_y)
118             loss_list.append(loss.item())
119             # Backprop and perform Adam optimisation
120             optimizer.zero_grad()
121             loss.backward()
122             optimizer.step()
123
124             if (epoch+1) % log_freq == 0:
125                 # Calculate the validation loss
126                 model.eval()
127                 with torch.no_grad():
128                     X_val_pred = model(X_val)
129                     val_loss = criterion(X_val_pred, y_val)
130
131                 # diff_ = (X_val_pred - y_val.unsqueeze(1)).detach().cpu().numpy().
squeeze()
132                 # diff_vec = np.reshape(diff_, (diff_.shape[0], -1))
133                 # val_l2_pt_error = np.mean(np.linalg.norm(diff_vec, axis=1) / np.
linalg.norm(np.reshape(y_val.detach().cpu().numpy(), (y_val.shape[0], -1)),
axis=1), axis=0) * 100
134                 # rel_error = 100 * np.linalg.norm(diff_vec, axis=1) / np.linalg.norm(
np.reshape(y_val.detach().cpu().numpy(), (y_val.shape[0], -1)), axis=1)
135
136                 wandb.log({"val_loss": val_loss.item(), "train_loss": loss.item(), "
epoch": epoch})
137                 print (f"Epoch [{epoch+1}/{num_epochs}], Step [{i+1}/{total_step}], \
138                     Training Loss: {loss.item():.4f}, Validation Loss: {val_loss.
item():.4f}")
139
140             # Save the model checkpoint (optional)
141             save_path = os.path.join(wandb.run.dir, "model.ckpt")
142             torch.save(model.state_dict(), save_path)
143
144 wandb.agent(sweep_id, train)

```

```

145 t2 = time.time()
146 print(f"Total time taken: {t2-t1}")
147 wandb.finish()
148
149 # %% [markdown]
150 # Best performance obtained with:
151 # - activation: "leakyrelu"
152 # - batch_size: 64
153 # - lr: 0.0049603811225372675
154 # - neurons1: 13
155 # - neurons2: 15
156
157 # %%
158 train_loader = torch.utils.data.DataLoader(dataset=train_dataset,
159                                             batch_size=64,
160                                             shuffle=True)
161 total_step = len(train_loader)
162 model = Network(13, 15, nn.LeakyReLU())
163 criterion = nn.MSELoss()
164 optimizer = torch.optim.Adam(model.parameters(), lr=0.00496)
165 for epoch in range(num_epochs):
166     for i, (train_x, train_y) in enumerate(train_loader):
167         # Run the forward pass
168         model.train()
169         output = model(train_x)
170         loss = criterion(output, train_y)
171         # Backprop and perform Adam optimisation
172         optimizer.zero_grad()
173         loss.backward()
174         optimizer.step()
175
176 # %%
177 heights = [40,57,77,97,117,137,157,177,197]
178
179 df_X = pd.DataFrame(model(X_val).detach().numpy())
180 df_y = pd.DataFrame(y_val.numpy())
181 df_X.columns = [f"WS {z}m" for z in heights]
182 df_y.columns = [f"WS {z}m" for z in heights]
183
184
185 fig = plt.figure(figsize=(15, 10))
186 fig.tight_layout(pad=3)
187 spec = GridSpec(int(np.ceil(len(y_val[0])/3))+2,3, figure=fig)
188 column_headers = []
189
190 row_headers = ["R ", "Avg. error (m/s)", "Avg. absolute error (m/s)", "Standard
191               deviation (m/s)", "Avg. relative error (%)"]
192
193 cell_text = np.zeros((len(row_headers), len(heights)))
194
195 r = 0
196 c = 0
197
198 for i, z in enumerate(heights):
199     if len(heights)%3 == 1:
200
201         if c > 2 and i != len(heights)-1:
202             c = 0

```

```

203         r += 1
204
205         if i == 9:
206             c = 1
207             r += 1
208
209     if len(heights)%3 != 1:
210
211         if c > 2:
212             c = 0
213             r += 1
214
215     column_headers += [str(z)+"m"]
216
217
218     # R value
219     lr = st.linregress(df_X[f"WS {z}m"], df_y[f"WS {z}m"]) #scipy stats built-in
linear regression function
220
221     cell_text[0][i] = '{: 0.2f}'.format(lr.rvalue**2)
222
223
224     # Mean error
225     df_subtract = df_y[f"WS {z}m"].subtract(df_X[f"WS {z}m"], axis = 0)
226     subtract_avg = df_subtract.mean()
227
228     cell_text[1][i] = '{: 0.2f}'.format(subtract_avg)
229
230
231     # Mean absolute error
232     df_subtract_abs = df_subtract.abs()
233     subtract_abs_avg = df_subtract_abs.mean()
234
235     cell_text[2][i] = '{: 0.2f}'.format(subtract_abs_avg)
236
237
238     # Standard deviation
239     cell_text[3][i] = '{: 0.2f}'.format(np.std(df_y[f"WS {z}m"]))
240
241
242     # Average relative error
243     df_relative_error = 100 * df_subtract_abs.divide(df_X[f"WS {z}m"], axis = 0)
244     relative_error_avg = df_relative_error.mean()
245     cell_text[4][i] = '{: 0.2f}'.format(relative_error_avg)
246
247
248
249     ## Histogram
250
251     ax = fig.add_subplot(spec[r, c])
252
253     ax.hist2d(df_X[f"WS {z}m"], df_y[f"WS {z}m"], 100, cmin=0.0001, cmap='Blues')
254
255     if c == 0:
256         ax.set_ylabel("WS Predicted [m/s]")
257
258     if r == int(np.ceil(len(heights)/3))-1:
259         ax.set_xlabel("WS Flidar [m/s]")
260

```

```

261
262     ax.set_ylim([0, 25])
263
264     xpoints = ypoints = ax.get_xlim()
265     ax.axline((xpoints[0], ypoints[0]), (xpoints[1], ypoints[1]), linestyle='--',
266               color='k', lw=1)
267
268     # on the title, I write the height of the plot and the R2 value
269     ax.set_title(f"WS Correlation at {z}m;  $R^2$  = {lr.rvalue**2: 0.2f}")
270
271     c += 1
272
273 rcolors = plt.cm.Blues(np.full(len(row_headers), 0.1))
274 ccolors = plt.cm.Blues(np.full(len(column_headers), 0.1))
275
276 ax = fig.add_subplot(spec[int(np.ceil(len(heights)/3)):int(np.ceil(len(heights)/3)
277                        )+2, :])
278
279 ytable = ax.table(cellText=cell_text,
280                  rowLabels=row_headers,
281                  rowColours=rcolors,
282                  rowLoc='center',
283                  colColours=ccolors,
284                  colLabels=column_headers,
285                  loc='center',
286                  bbox=[0.1, 0.05, 0.9, 0.9])
287
288 ax.axis('tight')
289 ax.axis('off')
290
291 fig.suptitle("Model evaluation", fontweight="bold")
292
293 cellDict = ytable.get_celld()
294 for i in range(0, len(cell_text[0])):
295     cellDict[(0,i)].set_height(.15)
296     for j in range(1, len(cell_text)+1):
297         cellDict[(j,i)].set_height(.1)
298
299 for i in range(1, len(cell_text)+1):
300     cellDict[(i,-1)].set_height(.1)
301
302 fig.subplots_adjust(left=0.05, bottom=0.04, right=0.95, top=0.92, hspace=0.5,
303                    wspace=0.2)
304
305 plt.show()

```

Listing 1: main.ipynb

```

1 import pandas as pd
2
3 def create_dataframes():
4     # Extract Satellite Data
5     sat_data = pd.read_csv("A_ERA.csv", sep=",", index_col=0, parse_dates=True)
6     sat_data.index = pd.to_datetime(sat_data.index)
7     sat_data.drop(["longitude", "latitude"], axis=1, inplace=True)
8     sat_data.dropna(inplace=True)
9

```

```

10 # Extract Buoy Data
11 buoy_data = pd.read_csv("A_MeteoStation_Data.csv", sep=";", index_col=0,
12 parse_dates=True)
13 buoy_data = buoy_data.loc[:,["wind_speed [m/s]", "air_temperature [ C ]", "
14 wind_direction [ ]"]]
15 buoy_data.index = pd.to_datetime(buoy_data.index)
16 buoy_data = buoy_data.groupby(pd.Grouper(freq="1h")).mean()
17 buoy_data.dropna(inplace=True)
18
19 # Extract Anemometer Data
20 df_temp = buoy_data.merge(sat_data, right_index=True, left_index=True)
21 df_temp.dropna(inplace=True)
22 df_subtract = df_temp["air_temperature [ C ]"].subtract(df_temp["sst"]) +
23 273.15
24 anemo_data = pd.DataFrame([df_temp["wind_speed [m/s]"], df_temp["
25 wind_direction [ ]"], df_subtract]).transpose()
26 anemo_data.columns = ["WS", "Wdir", "Delta T"]
27
28 df_temp_lidar = pd.read_csv("A_Wind_Data.txt", sep="\t", parse_dates=True,
29 index_col=0)
30 df_temp_lidar = df_temp_lidar.loc[:,df_temp_lidar.columns.str.contains("
31 wind_speed ")]
32 df_temp_lidar = df_temp_lidar.groupby(pd.Grouper(freq="1h")).mean()
33 df_temp_lidar.dropna(inplace=True)
34
35 heights_A = [40,57,77,97,117,137,157,177,197]
36 regex = ""
37 for i in range(len(heights_A)):
38     z = heights_A[i]
39     regex += str(z)
40     if i!=len(heights_A)-1:
41         regex += "|"
42
43 lidar_data = df_temp_lidar.filter(regex=regex)
44 lidar_data.columns = [f"WS {z}m" for z in heights_A]
45
46 df_temp_X = anemo_data.merge(sat_data, right_index=True, left_index=True)
47 df_temp_X["month"] = df_temp_X.index.month
48 df_temp_X = df_temp_X[df_temp_X["WS"] != 0]
49 df_temp_X.dropna(inplace=True)
50
51 df_temp_y = lidar_data.groupby(pd.Grouper(freq="1h")).mean()
52 df_temp_y.dropna(inplace=True)
53
54 df_temp_Xy = df_temp_X.merge(df_temp_y, right_index=True, left_index=True)
55 df_temp_Xy.dropna(inplace=True)
56 df_temp_Xy = df_temp_Xy.sample(frac=1)
57
58 Xy_data = df_temp_Xy
59 X_data = df_temp_Xy.loc[:, ~df_temp_Xy.columns.str.contains("WS ")]
60 y_data = df_temp_Xy.loc[:, df_temp_Xy.columns.str.contains("WS ")]
61
62 return X_data, y_data

```

Listing 2: extract_data.py