

blatt09_Guth_Venker_Jaekel

June 29, 2021

1 Abgabe SMD Blatt 09

1.0.1 von Nico Guth, David Venker, Jan Jäkel

```
[1]: class PDF(object):  
    '''  
    Displays PDF files in Jupyter Notebooks  
    '''  
    def __init__(self, pdf, size=(200,200)):  
        self.pdf = pdf  
        self.size = size  
  
    def _repr_html_(self):  
        return '<iframe src={0} width={1[0]} height={1[1]}></iframe>'.  
        ↪format(self.pdf, self.size)  
  
    def _repr_latex_(self):  
        return r'\includegraphics[width=1.0\textwidth]{{{0}}}'.format(self.pdf)
```

2 Aufgabe 18 Projektaufgabe: Feature Engineering

2.1 a) Optimieren Sie den „roc auc score“, indem Sie neue Features generieren.

Der Code ist weiter unten im PDF eingebunden.

2.2 b) Wie hoch ist der größte „roc auc score“ den Sie erreichen können?

roc auc score = 0.870

2.3 c) Was würde ein „roc auc score“ von 0.5 bedeuten, was ein „roc auc score“ von 1?

- 0.5 : Die Klassifizierung ist genau so gut wie zufälliges Raten
- 1.0 : Die Klassifizierung ist perfekt auf den gegebenen Daten -> Overfitting

2.4 d) Wie ändern sich Genauigkeit und Sensitivität mit steigendem „prediction threshold“?

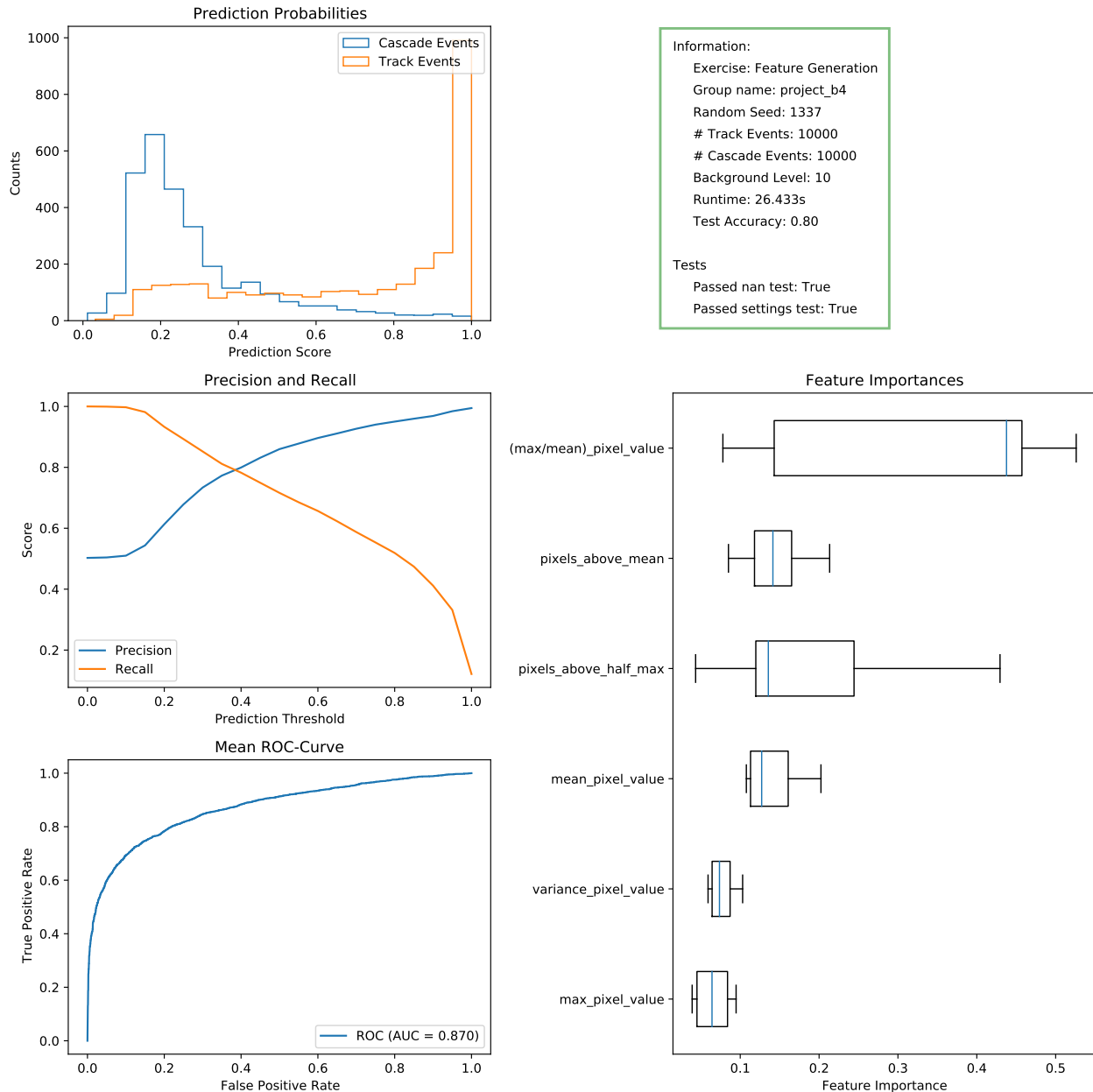
- Bei hohem Prediction Threshold steigt der Recall bzw. Sensitivität (TP/TP+FN)

- Bei niedrigem Prediction Threshold steigt die Precision bzw. Genauigkeit ($TP/TP+FP$)

2.5 e) Übersichts PDF

```
[2]: PDF('feature_generation.pdf', size=(500,500))
```

```
[2]:
```



3 Aufgabe 19 Projektaufgabe: Energieregression

3.1 a) Implementieren Sie den Random Forest Regressor b) Implementieren Sie 5-Fache Kreuzvalidierung

Der Code ist weiter unten im PDF eingebunden.

3.2 c) Regressor Confusion

3.2.1 Beschreiben Sie kurz, was eine Migrationsmatrix ist.

Wenn eine Migrationsmatrix eine Konfusionsmatrix ist:

In der Konfusionsmatrix kann man ablesen für welchen wahren Wert x wie viele Ereignisse als y vorhergesagt wurden.

Hier kann man also die Zusammenhänge der wahren Werte zu den vorhergesagten Werten einschätzen und so ein Modell bewerten.

3.2.2 Welche Eigenschaften Ihres Regressors können Sie von der Konfusionsmatrix ablesen?

Die Konfusionsmatrix ist auf der Diagonalen am größten und ungefähr um die Diagonale normalverteilt.

Also kann das Modell die Energie recht gut vorhersagen.

Aufgrund der log-log-Skala sollten die Abstände nach oben rechts kleiner werden. Allerdings ist oben rechts die Streuung recht groß. Also sieht man außerdem, dass große Energien nicht so gut vorhergesagt werden. Kleine Energien werden auch teilweise über eine Größenordnung falsch vorhergesagt.

3.3 d) Welche Eigenschaften Ihres Regressors können Sie an den Parametern „Bias“ und „Resolution“ ableiten?

Annahme: Mit Resolution ist Variance gemeint.

Bei Energien um 10^4 bis 10^5 GeV ist der Bias und die Variance niedrig.

Bei kleineren und größeren Energien wird sowohl Bias als auch Variance höher.

Also wird die Energie dort ungenauer aufgelöst (Variance) und im Mittel zu hoch bzw. zu niedrig eingeschätzt. (Bias)

Niedrige Energien werden eher zu hoch eingeschätzt und hohe Energien eher zu niedrig.

Am besten scheint der Regressor zwischen $2 \cdot 10^4$ und $2 \cdot 10^5$ GeV zu funktionieren.

3.4 e) Welche drei Features haben für die Energieschätzung in diesem Beispiel die größte Bedeutung?

Mit Abstand hat die Anzahl der Pixel mit einem Wert über dem Mittelwert die höchste Bedeutung. Gefolgt davon sind die Anzahl der Pixel mit einem Wert über der Hälfte des Maximums und der maximale Wert durch den Mittelwert.

In etwas anderer Reihenfolge sieht man auch in Aufgabe 18, dass diese Features am aussagekräftigsten sind.

3.5 f) Kann es hilfreich sein die zu schätzende Zielgröße (Teilchenenergie) zu skalieren?

Ja, denn hohe Werte können den Regressor teilweise daran hindern niedrige Werte gut einschätzen zu können.

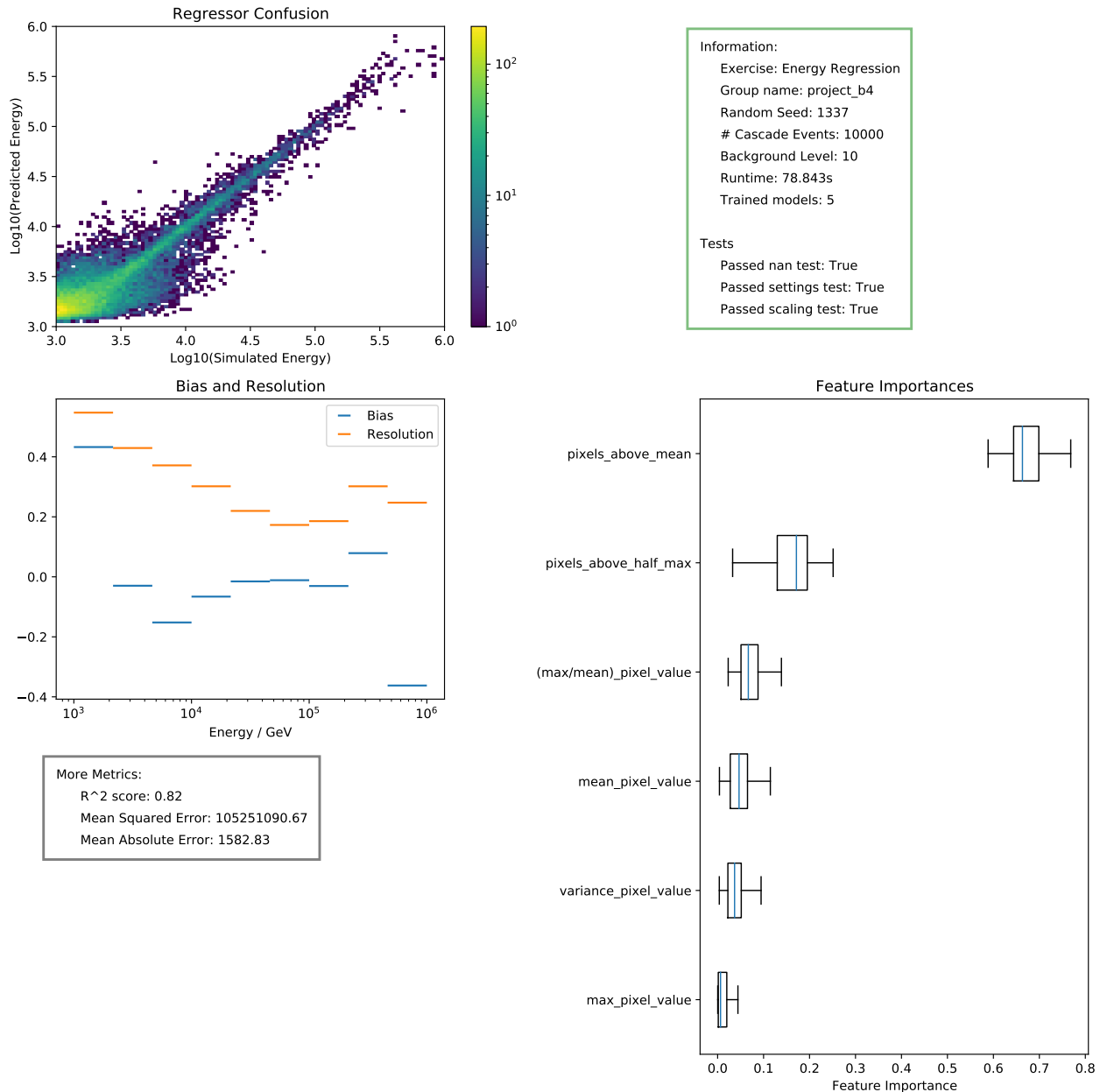
Z.B. wenn man als Loss-Function den Mean-Squared-Error benutzt, ist es ohne Skalierung deutlich

günstiger hohe Teilchenenergien richtig einzuschätzen als niedrige, obwohl die Problemstellung unter Umständen mehr an niedrigen Werten interessiert sein könnte.

3.6 g) Übersichts PDF

```
[3]: PDF('energy_regression.pdf', size=(500,500))
```

[3]:



[]:

project_b4/reconstruction/preprocessing/feature_generation.py
2021-06-25T15:49+02:00

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 class FeatureGenerator():
5     """This class generates features for an event
6
7     Based on the values measured in the detector, arbitrary features
8     can be generated. These can be used for machine learning exercises
9     such as classifying the particle type or estimating energy and
10    ↪ origin.
11    """
12
13     def __init__(self, detector):
14         """Sets up the feature generator
15
16         Parameters
17         -----
18         detector : Detector object
19             A detector object for which the feature generator will be
20    ↪ set up.
21         """
22         self.detector = detector
23
24     def analyse(self, event):
25         """Reconstruct an event measured in the assigned detector.
26
27         Parameters
28         -----
29         event : Event
30             The event which we want to reconstruct.
31
32         Returns
33         -----
34         Features: Dict
35             A dictionary of arbitrary features.
36             Values are to be scalar to be fed into a random forest
37         """
38         # -----
39         # Exercise 18 (Sheet 9):
40         # -----
```

```

39     # Feature Generation (exercises/feature_generation.py)
40
41     mean = np.mean(event.pixels)
42     max = np.max(event.pixels)
43
44     features = {
45         'max_pixel_value': max,
46         'variance_pixel_value': np.var(event.pixels),
47         'mean_pixel_value': mean,
48         '(max/mean)_pixel_value': max/mean,
49         'pixels_above_half_max': (event.pixels > 0.5*max).sum(),
50         'pixels_above_mean': (event.pixels > mean).sum()
51     }
52
53
54     return features

```

project_b4/reconstruction/machine_learning/energy_regression.py
2021-06-25T18:25+02:00

```
1 import numpy as np
2 from sklearn import (
3     ensemble, linear_model, neighbors, svm, tree, naive_bayes,
4     gaussian_process, neural_network, dummy)
5 from sklearn.model_selection import KFold
6 from sklearn.base import clone
7 from tqdm import tqdm
8
9
10 def define_model(seed):
11     """A helper function to retrieve a model for the energy regression
12     exercise. The parameter seed gets used to set the random state of
13     ↪ the model
14     and ensure reproducible results"""
15
16     # -----
17     # Exercise 19 (Sheet 9):
18     # -----
19     # Energy_regression (exercises/energy_regression.py)
20
21     model = ensemble.RandomForestRegressor(random_state=seed)
22
23     return model
24
25 def cross_validate_model(X, y, model, seed):
26     """This function implements a cross validation on a given model.
27
28     Required return values:
29     -----
30     predictions: np.array of the same shape as y
31     (These are the predictions on the test sets combined into one
32     ↪ array)
33     true_values: np.array of the same shape as y
34     (These are the y's chosen in the different cross validation
35     ↪ steps
36     combined into one array. This equals y but with different
37     ↪ order.)
38     models: list of (n_cross_validation_steps) models
39     These are used to calculate the feature importances later
```

```

37     """
38     # -----
39     # Exercise 19 (Sheet 9):
40     # -----
41     # Energy_regression (exercises/energy_regression.py)
42
43     kf = KFold(n_splits=5, shuffle=True, random_state=seed)
44
45     predictions = []
46     true_values = []
47     models = []
48
49     for train_index, test_index in kf.split(X):
50         X_train, y_train = X.iloc[train_index], y[train_index]
51         X_test, y_test = X.iloc[test_index], y[test_index]
52
53         model.fit(X_train, y_train)
54         y_pred = model.predict(X_test)
55
56         predictions.append(y_pred)
57         true_values.append(y_test)
58         models.append(model)
59
60     predictions = np.concatenate(predictions)
61     true_values = np.concatenate(true_values)
62
63     return predictions, true_values, models

```