

Lithology prediction with Python

Recently, machine learning has increased in popularity in many industries. In the Oil and Gas industry the implementation of data driven methods has increased, specially in the Upstream sector, where machine learning models are used to leverage the vast amount of data available for geoscientist although in some circumstances this data is relatively small when compared to the task that must be completed.

The main goal of this project is to build a machine learning model using popular python libraries and datasets available on the **FORCE 2020 ML competition** Github repository. The datasets available on the official Github repository are the training data, and test data, with the earlier having 29 columns and over one million observations, and the latter with just as many columns but with only around 122,000 examples. The final model will be build gradually, starting from a simple model trained on an imbalanced dataset with only the well logs and the depth as features, and increasing in complexity as we perform different data imputation techniques on the training data to improve its quality and boost learning, and tune the hyperparameters to improve the model performance.

Importing libraries

First, we will import commonly available python libraries. These libraries contain useful classes and functions that allow us to write a few lines of code in order to perform task that otherwise would require a lot of custom written code. Following is a short description of each library:

- Pandas:** A Python library build over numpy, useful for data analysis tasks. Stores functions that allow us to generate statistics for our data, manipulate the data, and visualize it through wrappers over other libraries such as Matplotlib.
- Numpy:** A scientific Python library that allows the generation of matrices and vectors, and supports complex operations between such data structures. Numpy also allows for vectorized operations which means that multiple operations are done on a single cup cycle.
- Matplotlib:** Maybe Python's the most commonly used data visualization library. It is used for visualizing data by generating plots.
- Seaborn:** A data visualization library that lets us further customize plots and generate more specific plots like violin and linear regression plots, as well as facilitating the generation of conditional plots and pairplots.
- Scikit-learn:** A popular machine learning library with tools for data processing, building data pipelines, and training, building and evaluating machine learning models.
- Boruta:** An open source Python library used for feature selection. Commonly used by kagglers on competitions.
- Umap-learn:** An open source library for Uniform Manifold Approximation and Projection a method for dimensionality reduction commonly preferred over PCA (Principal Component Analysis) and autoencoders because its results are easier to visualize on plots than the results of PCA and it is faster than autoencoders.

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, KFold, cross_val_score
from sklearn.metrics import accuracy_score
from sklearn.neural_network import MLPClassifier
from sklearn.svm import SVC, OneClassSVM
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.preprocessing import PowerTransformer
from sklearn.compose import ColumnTransformer
from boruta import BorutaPy
import umap
import matplotlib inline
```

For this project we will also use some custom functions, which are commonly stored in python scripts. We can generate folders that store scripts for functions of a similar topic and add such folders to our current python path by using the built-in `sys` python library. This allows us to import this custom functions as if we were importing built-in python available open source libraries.

```
In [2]: import sys
sys.path.append('formation_evaluation')
```

```
In [3]: import well_log_display
```

It should be noted that the `well_log_display` function, as well as other functions in the "formation_evaluation" folder, was originally coded by Yohanes Nuwara and is available on his own Github repository.

Loading the data

We will use the pandas library `pd.read_csv()` function to read our dataset, the default separator for csv files is the "." character, however our file uses the ";" character to separate values.

```
In [4]: logs_data = pd.read_csv('datasets1\log_train.csv', sep=';')
logs_data.head(10)
```

	WELL	DEPTH_MD	X_LOC	Y_LOC	Z_LOC	GROUP	FORMATION	CALI	RSHA	RMED	...	RDEP	DTS	GR	NPHI	PEF	DTC
0	15/9-13	494.528	437641.96875	6470972.5	-469.501831	NORDLAND	GP.	NaN	19.480835	NaN	1.611410	...	34.636410	NaN	Na	Na	Na
1	15/9-13	494.680	437641.96875	6470972.5	-469.653809	NORDLAND	GP.	NaN	19.468800	NaN	1.618070	...	34.636410	NaN	Na	Na	Na
2	15/9-13	494.832	437641.96875	6470972.5	-469.805786	NORDLAND	GP.	NaN	19.468800	NaN	1.626459	...	34.779556	NaN	Na	Na	Na
3	15/9-13	494.984	437641.96875	6470972.5	-469.957794	NORDLAND	GP.	NaN	19.459282	NaN	1.621594	...	39.965164	NaN	Na	Na	Na
4	15/9-13	495.136	437641.96875	6470972.5	-470.109772	NORDLAND	GP.	NaN	19.453100	NaN	1.602679	...	57.483765	NaN	Na	Na	Na
5	15/9-13	495.288	437641.96875	6470972.5	-470.261780	NORDLAND	GP.	NaN	19.453100	NaN	1.585567	...	75.281410	NaN	Na	Na	Na
6	15/9-13	495.440	437641.96875	6470972.5	-470.413788	NORDLAND	GP.	NaN	19.462496	NaN	1.576569	...	76.199951	NaN	Na	Na	Na
7	15/9-13	495.592	437641.96875	6470972.5	-470.565796	NORDLAND	GP.	NaN	19.468800	NaN	1.587011	...	76.199951	NaN	Na	Na	Na
8	15/9-13	495.744	437641.96875	6470972.5	-470.717773	NORDLAND	GP.	NaN	19.468800	NaN	1.613674	...	75.898796	NaN	Na	Na	Na
9	15/9-13	495.896	437641.96875	6470972.5	-470.869782	NORDLAND	GP.	NaN	19.468800	NaN	1.634622	...	68.121262	NaN	Na	Na	Na

10 rows × 29 columns

Sometimes our dataset is too big to be completely displayed. We can use the `pd.set_option()` function to specify the maximum number of rows or columns to be displayed.

```
In [5]: pd.set_option('display.max_columns', None)
logs_data.head(10)
```

•	PEF : Photo-electric factor
•	DC : Density tool compressional wave
•	SP : Spontaneous potential log
•	BS : Bit size
•	RQ : Rate of penetration
•	DTS : Density tool shear wave
•	DCAL : Differential caliper
•	DRHO : Bulk density correction
•	MUDWEIGHT : Density of drilling mud
•	RMIC : Microresistivity log
•	ROPAC : Electromagnetic propagation log (unsure if this is the correct log)
•	RXC : Resistivity of flushed zone

Now we will use the `DataFrame.describe()` method to display some basic summary statistics for our datasets. This lets us quickly scan the maximum and minimum values for the data in each column, along with other statistics to see if there is any unusual range or distribution. This is easier on smaller datasets. For datasets with many columns (like our current dataset), visual exploratory data analysis through plots may be preferable.

We can set the `include` parameter to display statistics for columns storing data in numerical or object type.

```
logs_data.describe(include=np.number)
```

	DEPTH_MD	X.LOC	Y.LOC	Z.LOC	CALI	RSHA	RMED	RDEP	RHOB
count	1.170511e+06	1.159738e+06	1.159738e+06	1.193736e+06	1.062834e+06	630650.000000	1.131518e+06	1.159496e+06	1.095242e+06
mean	2.184087e+03	4.856310e+05	6.681276e+05	-2.138527e+03	1.1318565e+01	10.694664	4.966978e+00	1.069101e+01	2.284987e+00
std	8.871931e+03	2.455641e+05	1.381263e+06	8.709436e+02	2.398903e+01	100.641507	6.457370e+01	1.138490e+01	2.623303e+00

Exploring the data

For quick basic data exploration we can use the `DataFrame.info()` method to display comprehensive information about our dataset like the number of data points(rows), the number of columns, and the name, number of non-null elements and dtype of each column.

```
In [6]: logs_data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1170511 entries, 0 to 1170510
Data columns (total 29 columns):
# Column                                Non-Null Count  Dtype
---  ---                                ---
0 WELL                                1170511 non-null  object
1 DEPTH_MD                           1170511 non-null  float64
2 X_LOC                              11979736 non-null float64
3 Y_LOC                              11979736 non-null float64
4 Z_LOC                              11979736 non-null float64
5 GROUP                              1169293 non-null  object
6 FORMATION                          1033917 non-null  object
7 CALI                               630650 non-null  float64
8 RSHA                               630650 non-null  float64
9 RMED                               113316 non-null  float64
10 RDEP                              1134986 non-null float64
11 RHOB                              1009242 non-null float64
12 GR                                1170511 non-null float64
13 NPHI                              69353 non-null   float64
14 PEF                               765409 non-null float64
15 DTS                               671692 non-null float64
16 DTC                               1089648 non-null float64
17 SP                               864247 non-null float64
18 DRHO                              682657 non-null float64
19 ROP                                535071 non-null float64
20 DTS                               174613 non-null float64
21 DCAL                              298933 non-null float64
22 DRHO                              987857 non-null float64
23 MUDWEIGHT                         316151 non-null float64
24 RMIC                              176160 non-null float64
25 ROPA                              192325 non-null float64
26 RXO                              327427 non-null float64
27 FORCE_2020_LITHOFACIES_LITHOLOGY  1170511 non-null int64
28 FORCE_2020_LITHOFACIES_CONFIDENCE 1170332 non-null float64
dtype: float64(23), int64(1), object(3)
memory usage: 259.0+ MB
```

Dtypes for all columns in our dataframe seem correct, however there is a noticeable amount of missing data in several columns.

Some column names are self explanatory, while others are commonly used log mnemonics. Following is a short list with the mnemonic and what it corresponds to.

- CALI:** Caliper
- RSHA:** Shallow resistivity
- RMED:** Medium resistivity
- RDEP:** Deep resistivity
- RHOB:** Bulk resistivity
- GR:** Gamma Ray
- SGR:** Spectral Gamma ray
- NPHI:** Neutron porosity
- PEF:** Photo-electric factor
- DTC:** Density tool compressional wave
- SP:** Spontaneous potential
- BS:** Bit size
- RDP:** Rate of penetrator
- DTS:** Density tool shear wave
- DCAL:** Differential caliper
- DRHO:** Bulk density correction
- MUDWEIGHT:** Density of drilling mud
- RMIC:** Microresistivity log
- ROPA:** Electromagnetic propagation log (unsure if this is the correct log)
- RXO:** Resistivity of flushed zone

Now we will use the `DataFrame.describe()` method to display some basic summary statistics for our datasets. This lets us quickly scan the maximum and minimum values of the data in each column, along with other statistics to see if there is any unusual range or distribution, this is easier on smaller datasets. For datasets with many columns (like our current dataset), visual exploratory data analysis through plots may be preferable.

```
We can set the include parameter to display statistics for columns storing data in numerical or object type.

In [7]: logs_data.describe(include=logc.number)
```

	DEPTH_MD	X_LOC	Y_LOC	Z_LOC	CALI	RSHA	RMED	RDEP	RHOB
count	1170511e+06	1.159736e+06	1.159736e+06	1.159736e+06	1.082634e+06	630650.000000	1.131518e+06	1.159496e+06	1.009242e+06
mean	1.218408e+03	4.856010e+05	6.681276e+06	-2.138527e+03	1.318568e+01	10.694664	4.986978e+00	1.069103e+01	2.284987e+00
std	9.97182e+02	3.455541e+04	1.281524e+05	9.709426e+00	3.798907e+00	100.642597	5.467269e+01	1.139480e+02	2.532835e+01
min	1.360860e+02	4.266886e+05	6.406641e+06	-5.395563e+03	2.344000e+00	0.000100	-8.418695e-03	9.102396e-02	2.092203e+00
25%	1.418597e+03	4.547996e+05	6.591327e+06	-2.811502e+03	9.429712e+00	0.854120	9.140862e-01	9.102396e-01	2.092203e+00
50%	2.076605e+03	4.976039e+05	6.737317e+06	-2.042785e+03	1.255575e+01	1.399920	1.443594e+00	1.439000e+00	3.321228e+00
75%	2.864393e+03	5.201532e+05	6.784886e+06	-1.391866e+03	1.671075e+01	3.099348	2.680930e+00	2.557220e+00	2.488500e+00
max	5.436632e+03	5.726238e+05	6.856661e+06	-1.110860e+02	2.827900e+01	2193.904541	1.988610e+03	1.999887e+03	3.457830e+00

```
In [8]: logs_data.describe(include=np.object)
```

	WELL	GROUP	FORMATION
count	1170511	1169293	1033517
unique	98	14	69
top	25/2-7	NORDLAND	GP.
freq	25131	293155	172636

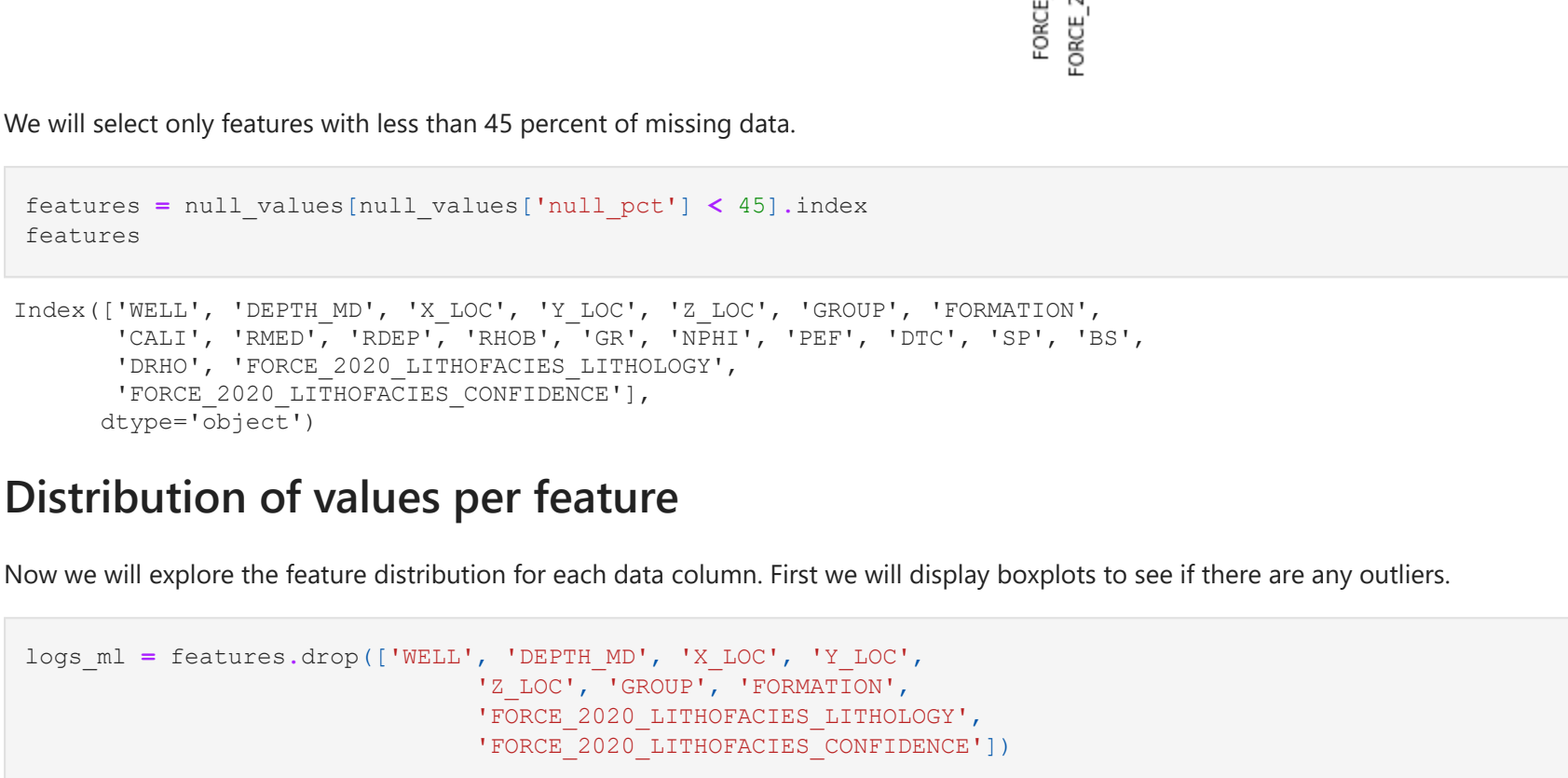
Visual exploratory data analysis

Before we explore each feature, lets first visualize the missing data in our dataset. This way we might focus our analysis in columns that have enough information, or where data can be imputed, to be useful for our machine learning algorithm.

```
In [9]: df = plot_null_matrix(df, figsize=(10, 8))
plt.figure(figsize=figsize)
```

```
df.null = df.isnull()
ax = sns.heatmap(-df.null, cbar=False)
plt.show()
```

```
In [10]: plot_null_matrix(logs_data)
```



Looking at the plot, it is obvious that there are some columns with a noticeable amount of missing data. We will generate a table with the count of null values and the percentage of the data missing for each column.

```
In [11]: null_count = logs_data.isnull().sum()
null_pct = (null_count / logs_data.shape[0]) * 100
null_values = pd.DataFrame({'null_count': null_count,
                           'null_pct': null_pct})
```

	null_count	null_pct
WELL	0	0.000000
DEPTH_MD	0	0.000000
X_LOC	10775	0.920538
Y_LOC	10775	0.920538
Z_LOC	10775	0.920538
GROUP	1278	0.109183
FORMATION	156994	11.703777
CALI	87877	7.507576
RSHA	539861	46.121822
RMED	38993	3.331280
RDEP	11015	0.941702
RHOB	161269	13.777658
GR	0	0.000000
NPHI	1101158	94.074981
SGR	405102	34.608987
PEF	498819	42.615490
DTC	80863	6.908350
SP	306264	26.164983
BS	487854	41.678720
ROP	635440	54.287401
DTS	955898	85.082327
DCAL	871678	74.469868
DRHO	182654	15.604368
MUDWEIGHT	854360	72.990344
RMIC	994351	84.950163
ROPA	978186	83.569142
RXO	843084	72.027004
FORCE_2020_LITHOFACIES_LITHOLOGY	0	0.000000
FORCE_2020_LITHOFACIES_CONFIDENCE	179	0.015292

Then we can plot these data to visually determine a threshold value for the percentage of missing data we will allow for the features to be used in our machine learning algorithm.

```
In [12]: null_values['null_pct'].plot(kind='bar', figsize=(10,8), title=
'Percentage of missing values per data column')
```

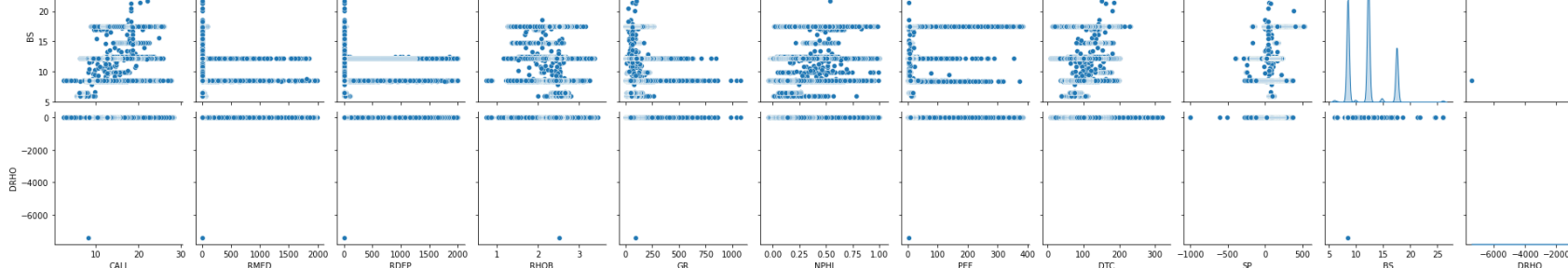


We will select only features with less than 45 percent of missing data.

```
In [13]: features = null_values[null_values['null_pct'] < 45].index
features
```

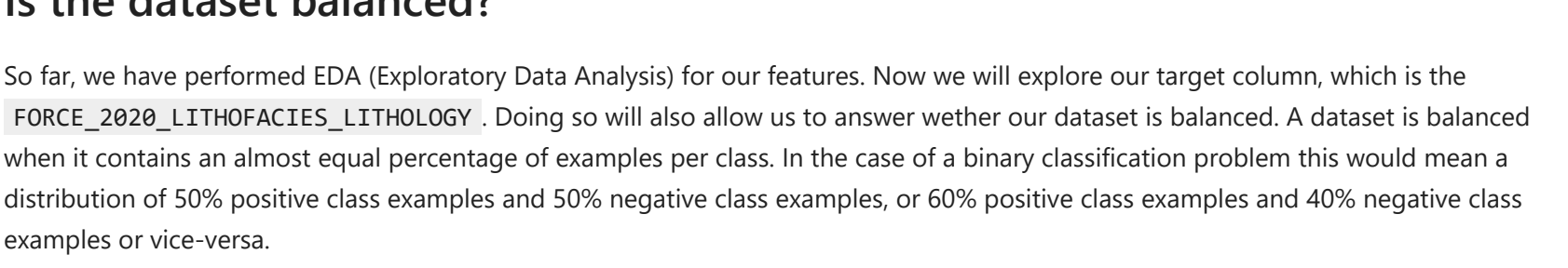
```
In [14]: Index(['WELL', 'DEPTH_MD', 'X_LOC', 'Y_LOC', 'Z_LOC', 'GROUP', 'FORMATION',
'CALI', 'RSHA', 'RMED', 'RDEP', 'RHOB', 'GR', 'NPHI', 'PEF', 'DTC', 'SP', 'BS',
'DRHO', 'FORCE_2020_LITHOFACIES_LITHOLOGY',
'FORCE_2020_LITHOFACIES_CONFIDENCE'],
dtype='object')
```

```
In [15]: fig, ax = plt.subplots()
ax = sns.boxplot(logs_data[logs_data[features]])
ax.set_title('Feature distribution')
plt.xticks(rotation=90)
plt.show()
```



Pairplots

```
In [16]: sns.pairplot(logs_data, vars=logs_ml, diag_kind='kde')
plt.show()
```



Looking at the boxplot and pairplots we can derive the following inferences:

- Value ranges probably are very different, and thus, some features are outscaled, making it difficult to visualize their distribution.
- The features do contain some outliers. An example of this is the **DRHO** data column, with an outlier that might be the result of an incorrect measurement or a tool error.
- Some features have a distribution that resembles a normal distribution (also known as Gaussian distribution), although skewed, other features have multimodal distributions.

To make the visualization of the feature distributions easier we can normalize the data, a data processing step already required in the machine learning workflow.

Is the dataset balanced?

So far, we have performed EDA (Exploratory Data Analysis) for our features. Now we will explore our target column, which is the **FORCE_2020_LITHOFACIES_LITHOLOGY**. Doing so will also allow us to answer whether our dataset is balanced. A dataset is balanced when it contains an almost equal percentage of examples per class. In the case of a binary classification problem this would mean a distribution of 50% positive class examples and 50% negative class examples, or 60% positive class examples and 40% negative class examples or vice-versa.

First, we will calculate the theoretical relative frequency (the percentage of examples for each class) for our dataset to be balanced.

```
In [17]: # For readability transform original lithology codes into labels
# lithology code correspondence was given in the competition starter notebook.
lithology_keys = ['Sandstone', 'Sandstone/Shale',
                  'Shale', 'Marl',
                  'Limestone', 'Dolomite',
                  'Limestone', 'Limestone',
                  'Chalk', 'Chalk',
                  'Anhydrite', 'Anhydrite',
                  'Coal', 'Coal',
                  'Basement']
```

```
In [18]: label = 'FORCE_2020_LITHOFACIES_LITHOLOGY'
unique_lithologies = logs_data[label].unique()
theoretical_pct_lithology = 100 / unique_lithologies.shape[0]
theoretical_pct_lithology
```

```
Out[18]: 8.333333333333333

In [19]: lithologies_labels = logs_data[label].map(lithology_keys)
lithologies_count = lithologies_labels.value_counts(normalize=True) * 100
lithologies_count
```

```
Out[19]: Shale      61.580199
Sandstone    14.432756
Sandstone/Shale 12.851788
Limestone    4.811574
Marl         2.847389
Puref        1.302423
Chalk        0.898155
Halite       0.701659
Coal         0.326353
Dolomite     0.144211
Anhydrite    0.092695
Basement     0.008800
Name: FORCE_2020_LITHOFACIES_LITHOLOGY, dtype: float64
```

```
In [20]: lithologies_count.plot(kind='bar',
                              title='Percentage of data points
```



```
'Chalk': 2,
'Coal': 3,
'Dolomite': 4,
'Halite': 5,
'Limestone': 6,
'Marl': 7,
'Sandstone': 8,
'Sandstone/Shale': 9,
'Shale': 10,
'Tuff': 11]
```

```
In [60]: test_norm["LITHOLOGY"] = test_norm["LITHOLOGY"].map(lithology_code)
```

```
In [61]: test_norm.columns
```

```
Out[61]: Index(['CALI', 'BMED', 'RDEP', 'RBOB', 'GR', 'NPHT', 'PEP', 'DTC', 'SP', 'BS',
              'DRHO', 'Z_LOC', 'LITHOLOGY'],
              dtype='object')
```

```
In [62]: test_data = test_norm.copy()
```

Model performance on test data

Now we will first evaluate baseline performance on the test data with the "zero rule" algorithm, and then use the model we built to predict the lithologies for the test data, and use our score function to evaluate its performace.

```
In [63]: # specify features and target arrays

X_test = test_data.drop('LITHOLOGY', axis=1)
y_test = test_data['LITHOLOGY']

# predict labels
ypr_pred = np.full(y_test.shape[0], 10)
y_pred = rfc.predict(X_test.values)

# calculate accuracy
baseline_accuracy_test = accuracy_score(y_test.values, ypr_pred)
rfc_accuracy_test = accuracy_score(y_test.values, y_pred)

[Parallel(n_jobs=6)]: Using backend ThreadingBackend with 6 concurrent workers.
[Parallel(n_jobs=6)]: Done 20 tasks | elapsed: 0.0s
[Parallel(n_jobs=6)]: Done 100 out of 100 | elapsed: 0.6s finished
```

```
In [64]: models_summary = dict()

ypr_summary = {'Accuracy_training': baseline_accuracy_training,
               'Accuracy_test': baseline_accuracy_test}

rfc_summary = {'Accuracy_training': rfc_accuracy_training,
               'Accuracy_test': rfc_accuracy_test}

models_summary['Zero rule algorithm'] = ypr_summary
models_summary['Random Forest Classifier'] = rfc_summary

models_summary_df = pd.DataFrame.from_dict(models_summary, orient='index')
```

```
Out[64]:
```

	Accuracy_training	Accuracy_test
Zero rule algorithm	0.623315	0.600279
Random Forest Classifier	0.999998	0.717929

Conclusion

We managed to generate a simple model that performed better than the baseline model in both, the training and test data, without taking into account geolocation data (except for the `Z_LOC` column), and with default hyperparameters, although the model shows overfitting.

Future improvements

- Application of techniques to handle imbalanced datasets, both, before model training (oversampling and undersampling), and during model training (for machine learning algorithms that support this functionality).
- Extensive feature engineering and feature selection techniques.
- Try other data imputation techniques.
- Compare performance of other algorithms like Multilayer Perceptron Classifier or KNeighbors Classifier.
- Combine the predictions of different models to see if there is any improvement in accuracy.
- Build a confusion matrix to see if there are any patterns in the misclassification.

Also, in other notebooks, regarding the prediction of well log curves, the columns describing the depth at which the measurement was recorded is not scaled. It would be interesting to scale this column to see if there is any improvement on model performance. As far as I understand, there might be because this column outcales other features.

References

- Burkov, A. (2020).*Machine Learning Engineering*(1st ed.). True Positive Inc.
- Nuwara, Y. (2020,October 27).*A Demo on P-Sonic Log Prediction using Machine Learning in the Volve Field Dataset*
https://github.com/yohanesnuwara/volve-machine-learning/blob/main/notebook/demo_volve_soniclog_prediction.ipynb