

# Task 1: Predicting the Lunar Albedo based on Chemical Composition

author: [Nicklas Meyer](#)

In this task, we are asked to identify relationships between albedo and composition by dividing the data into two equal halves (left and right). We are required to train our models on the left half and predict the right. The task can be viewed as a standard regression one which aims to predict the brightness of each pixel. Finally, we create a framework for combining outputs of constituent models using the stacking framework.

```
In [1]: # Load essential libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.linear_model import LinearRegression, ElasticNet, Lasso
from sklearn.metrics import mean_absolute_error, mean_squared_error
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.model_selection import KFold, cross_val_score
from sklearn.preprocessing import StandardScaler, RobustScaler
from sklearn.base import BaseEstimator, TransformerMixin, RegressorMixin, clone
from sklearn.pipeline import make_pipeline
from sklearn.pipeline import Pipeline

from xgboost.sklearn import XGBRegressor
from lightgbm.sklearn import LGBMRegressor

from keras.models import Sequential
from keras.layers import Dense
from keras.wrappers.scikit_learn import KerasRegressor

plt.style.use('fivethirtyeight') # For better style
%matplotlib inline

Read the datasets provided at the github link
```

```
In [2]: BASE_URL = "https://raw.githubusercontent.com/ML4SCI/ML4Sci_GSoC/main/Messenger/Moon/"

# Load albedo map
albedo_map = pd.read_csv(BASE_URL + "Albedo_Map.csv", header=None)
NUM_ROWS = albedo_map.shape[0]
NUM_COLS = albedo_map.shape[1]
CENTER_HALF = NUM_COLS // 2

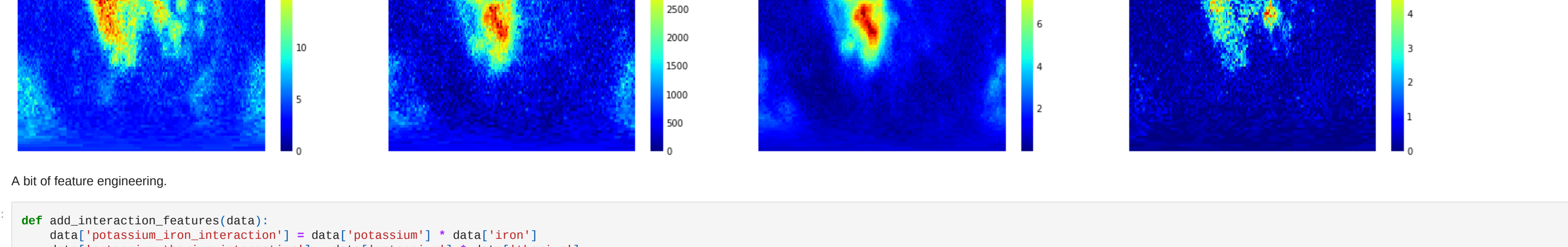
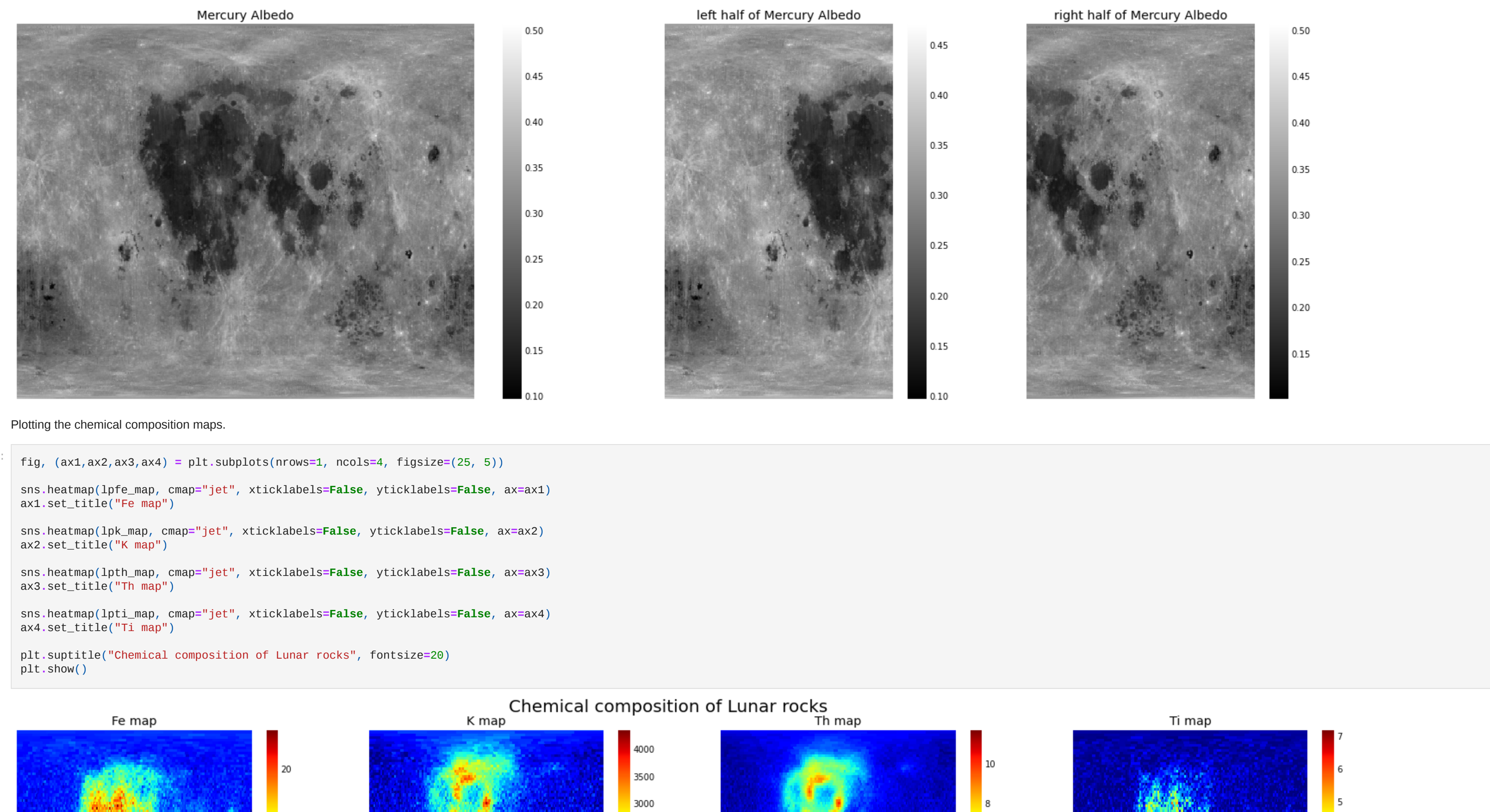
# Load iron map
lprf_map = pd.read_csv(BASE_URL + "LPFe_Map.csv", header=None)

# Load potassium map
lprk_map = pd.read_csv(BASE_URL + "LPK_Map.csv", header=None)

# Load thorium map
lpth_map = pd.read_csv(BASE_URL + "LPTh_Map.csv", header=None)

# Load titanium map
lpti_map = pd.read_csv(BASE_URL + "LPti_Map.csv", header=None)

Plotting the albedo map to confirm if we are reading the file correctly as a sanity check.
```



A bit of feature engineering.

```
In [5]: def add_interaction_features(data):
    data['potassium_iron_interaction'] = data['potassium'] * data['iron']
    data['potassium_thorium_interaction'] = data['potassium'] * data['thorium']
    data['potassium_titanium_interaction'] = data['potassium'] * data['titanium']
    data['iron_thorium_interaction'] = data['iron'] * data['thorium']
    data['iron_titanium_interaction'] = data['iron'] * data['titanium']
    data['thorium_titanium_interaction'] = data['thorium'] * data['titanium']
    return data

We divide the albedo and its corresponding chemical composition maps into two equal portions. The left-half becomes the training data and the right-half becomes the test data.
```

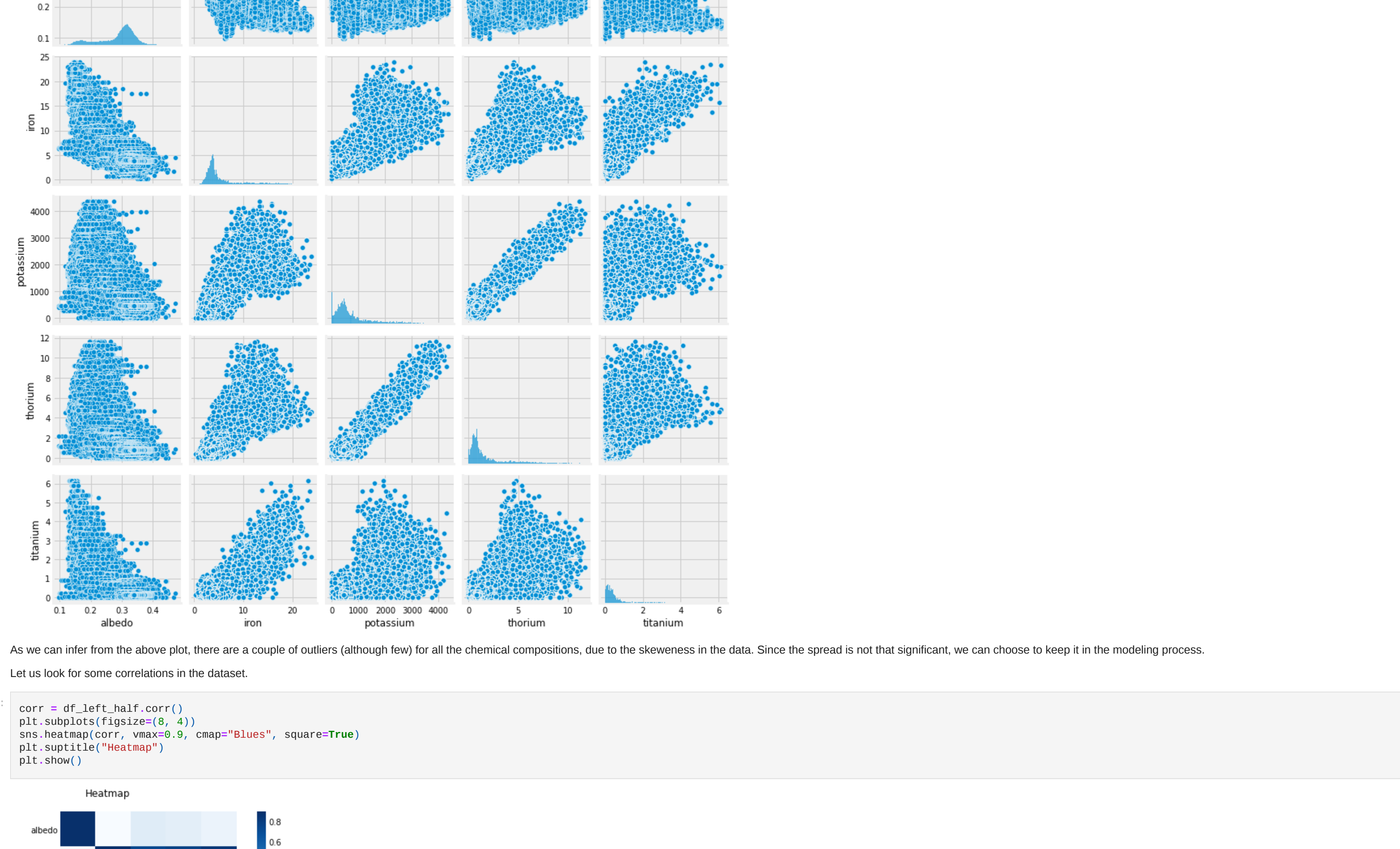
```
In [6]: df_left_half = pd.DataFrame()

df_left_half['albedo'] = albedo_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_left_half['iron'] = lprf_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_left_half['potassium'] = lprk_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_left_half['thorium'] = lpth_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_left_half['titanium'] = lpti_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()

df_right_half = pd.DataFrame()

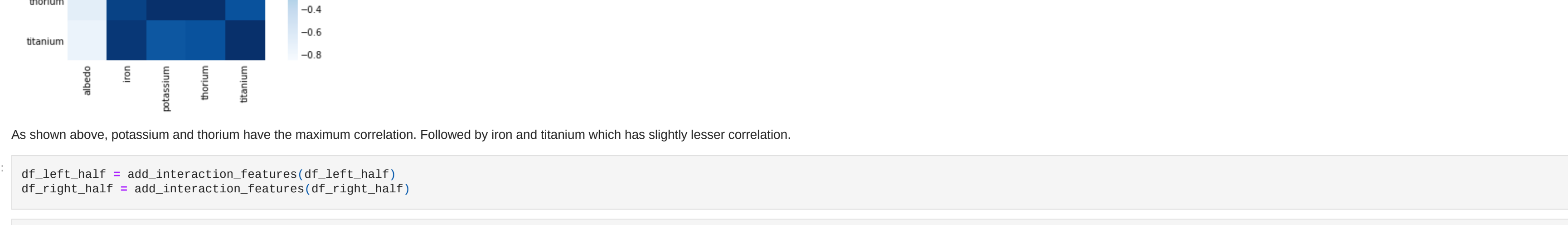
df_right_half['albedo'] = albedo_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_right_half['iron'] = lprf_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_right_half['potassium'] = lprk_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_right_half['thorium'] = lpth_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()
df_right_half['titanium'] = lpti_map.iloc[:, :CENTER_HALF].values.reshape(-1, 1).flatten()

Let us perform a pair-wise scatter plot.
```



As we can infer from the above plot, there are a couple of outliers (although low) for all the chemical compositions, due to the skewness in the data. Since the spread is not that significant, we can choose to keep it in the modeling process.

Let us look for some correlations in the dataset.



As shown above, potassium and thorium have the maximum correlation. Followed by iron and titanium which has slightly lesser correlation.

```
In [9]: df_left_half = add_interaction_features(df_left_half)
df_right_half = add_interaction_features(df_right_half)
```

```
In [10]: features = ['iron', 'potassium', 'thorium', 'titanium', 'potassium_iron_interaction', 'potassium_thorium_interaction',
               'potassium_titanium_interaction', 'iron_thorium_interaction', 'iron_titanium_interaction', 'thorium_titanium_interaction']
```

```
In [11]: X_train = df_left_half[features]
y_train = df_left_half['albedo'].values
X_test = df_right_half[features]
y_test = df_right_half['albedo'].values
```

```
In [12]: #Validation function
n_folds = 2 + 5

def mse_cv(model):
    kf = KFold(n_folds, shuffle=False).get_n_splits(X_train.values)
    mse = -cross_val_score(model, X_train.values, y_train, scoring='neg_mean_squared_error', cv=kf)
    return mse

We experiment with the following set of models:
```

```
In [13]: 1. Linear Regression
2. Lasso Regression
3. Elastic Net Regression
4. Gradient Boosting Regression
5. XGBoost
6. LightGBM
7. Neural Network (MLP)
```

```
In [13]: linear_reg = make_pipeline(StandardScaler(), LinearRegression())
lasso = make_pipeline(RobustScaler(), Lasso(alpha=0.0005, random_state=1))
Elnet = make_pipeline(RobustScaler(), ElasticNet(alpha=0.0005, l1_ratio=0.5, random_state=3))
GBoost = GradientBoostingRegressor(n_estimators=3000, learning_rate=0.05,
                                   max_depth=4, max_features='sqrt',
                                   min_samples_leaf=15, min_samples_split=10,
                                   loss='huber', random_state=1)
```

The procedure, for the training part, may be described as follows:

- Split the total training set into two disjoint sets, namely train and holdout.
- Train several base models on the first part (train)
- Test these base models on the second part (holdout)
- Use the predictions (from 3) (called out-of-folds predictions) as the inputs, and the correct responses (target variable) as the outputs to train a higher level learner called meta-model.

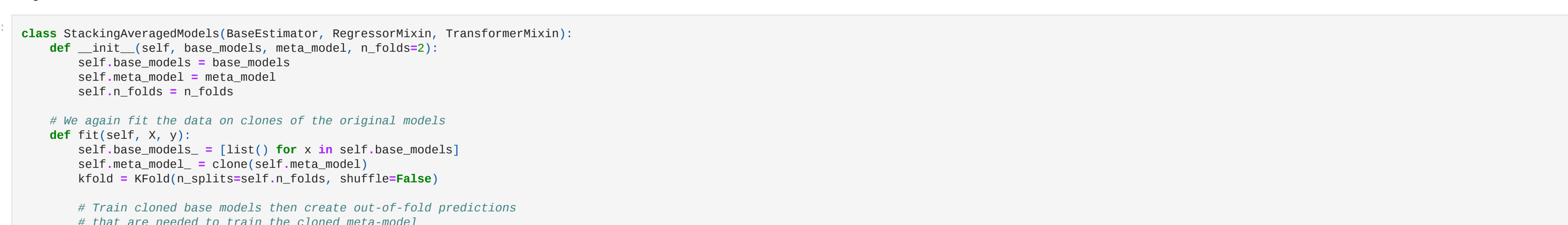


Image taken from [here](#).

```
In [14]: class StackingAveragedModels(BaseEstimator, RegressorMixin, TransformerMixin):
    def __init__(self, base_models, meta_model, n_folds=2):
        self.base_models = base_models
        self.meta_model = meta_model
        self.n_folds = n_folds

    # We again fit the data on clones of the original models
    def fit(self, X, y):
        self.base_models_ = [list() for x in self.base_models]
        self.meta_model_ = clone(self.meta_model)
        kfold = KFold(n_splits=self.n_folds, shuffle=False)

        # Train cloned base models then create out-of-fold predictions
        # that are needed to train the cloned meta-model
        out_of_fold_predictions = np.zeros((X.shape[0], len(self.base_models)))
        for i, model in enumerate(self.base_models):
            for train_index, holdout_index in kfold.split(X, y):
                instance = clone(model)
                self.base_models_[i].append(instance)
                instance.fit(X[train_index], y[train_index])
                y_pred = instance.predict(X[holdout_index])
                out_of_fold_predictions[holdout_index, i] = y_pred

        # Now train the cloned meta-model using the out-of-fold predictions as new feature
        self.meta_model_.fit(out_of_fold_predictions, y)
        return self

    # Use the predictions of all base models on the test data and use the averaged predictions as
    # meta-features for the final prediction which is done by the meta-model
    def predict(self, X):
        meta_features = np.column_stack([
            np.column_stack(model.predict(X) for model in self.base_models_).mean(axis=1)
            for base_models in self.base_models_
        ])
        return self.meta_model_.predict(meta_features)
```

```
In [15]: stacked_averaged_models = StackingAveragedModels(base_models=[linear_reg, Elnet, GBoost],
                                                         meta_model=lasso)

score = mse_cv(stacked_averaged_models)
print("Stacking Averaged models score: {:.4f} ({:.4f}).".format(score.mean(), score.std()))
```

```
In [16]: stacked_averaged_models.fit(X_train.values, y_train)
stacked_train_pred = stacked_averaged_models.predict(X_train.values)
stacked_pred = stacked_averaged_models.predict(X_test.values)
```

```
In [17]: print('MAE:', mean_absolute_error(y_test, stacked_train_pred))
print('MSE:', mean_squared_error(y_test, stacked_train_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, stacked_train_pred)))

MAE: 0.05952907669416993
MSE: 0.0068089997632278
RMSE: 0.0747767723185243
```

```
In [18]: model_xgb = XGBRegressor(colsample_bytree=0.4663, gamma=0.0468,
                             learning_rate=0.05, max_depth=3,
                             min_child_weight=1757, n_estimators=2200,
                             reg_alpha=0.458, reg_lambda=0.8571,
                             subsample=0.5213, silent=True,
                             random_state=1, nthread=-1)
```

```
In [19]: model_xgb.fit(X_train, y_train)
xgb_train_pred = model_xgb.predict(X_train)
xgb_pred = model_xgb.predict(X_test)
```

```
In [20]: print('MAE:', mean_absolute_error(y_test, xgb_train_pred))
print('MSE:', mean_squared_error(y_test, xgb_train_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, xgb_train_pred)))

MAE: 0.05939524868225496
MSE: 0.00623265253734534
RMSE: 0.07494982592532417
```

```
In [21]: model_lgb = LGBMRegressor(objective='regression', num_leaves=5,
                               learning_rate=0.05, n_estimators=728,
                               max_bin=50, bagging_fraction=0.6,
                               bagging_freq=5, feature_fraction=0.2310,
                               feature_fraction_seed=99, bagging_seed=99,
                               min_data_in_leaf=6, min_sum_hessian_in_leaf=1)
```

```
In [22]: model_lgb.fit(X_train, y_train)
lgb_train_pred = model_lgb.predict(X_train)
lgb_pred = model_lgb.predict(X_test.values)
```

```
In [23]: print('MAE:', mean_absolute_error(y_test, lgb_train_pred))
print('MSE:', mean_squared_error(y_test, lgb_train_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, lgb_train_pred)))

MAE: 0.05937702545958335
MSE: 0.00612088944724908
RMSE: 0.07482482877660752
```

```
In [24]: # Define wider model
def wider_model():
    model = Sequential()
    model.add(Dense(52, input_dim=10, kernel_initializer='normal', activation='relu'))
    model.add(Dense(10, kernel_initializer='normal', activation='relu'))
    # compile model
    model.compile(loss='mean_squared_error', optimizer='adam')
    return model
```

```
In [25]: np.random.seed(42)
estimators = []
estimators.append(('standardize', StandardScaler()))
estimators.append(('mlp', MetaRegressor(build_fn=wider_model, epochs=5, batch_size=20, verbose=0)))
m_model = Pipeline(estimators)

score = mse_cv(m_model)
print("NN models score: {:.4f} ({:.4f}).".format(score.mean(), score.std()))

Epoch 1/3 [=====] - 3s 1ms/step - loss: 0.8865
Epoch 2/3 [=====] - 2s 1ms/step - loss: 8.7649e-04
Epoch 3/3 [=====] - 2s 1ms/step - loss: 8.5896e-04
Epoch 4/3 [=====] - 2s 755us/step
Epoch 5/3 [=====] - 3s 1ms/step - loss: 0.0870
Epoch 6/3 [=====] - 2s 1ms/step - loss: 5.0916e-04
Epoch 7/3 [=====] - 2s 1ms/step - loss: 5.0472e-04
Epoch 8/3 [=====] - 2s 782us/step
NN models score: 0.0032 (0.0005)
```

```
In [26]: m_model.fit(X_train.values, y_train)
m_train_pred = m_model.predict(X_train.values)
m_pred = m_model.predict(X_test.values)

Epoch 1/3 [=====] - 5s 1ms/step - loss: 0.0841
Epoch 2/3 [=====] - 4s 1ms/step - loss: 7.3950e-04
Epoch 3/3 [=====] - 5s 1ms/step - loss: 7.2603e-04
Epoch 4/3 [=====] - 5s 760us/step
Epoch 5/3 [=====] - 3s 760us/step
```

```
In [27]: print('MAE:', mean_absolute_error(y_test, m_train_pred))
print('MSE:', mean_squared_error(y_test, m_train_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, m_train_pred)))

MAE: 0.060395441810003
MSE: 0.003115021251535
RMSE: 0.0744424472603612
```

```
In [28]: y_pred = (stacked_pred + xgb_pred + lgb_pred + m_pred) / 4
```

```
In [29]: print('MAE:', mean_absolute_error(y_test, y_pred))
print('MSE:', mean_squared_error(y_test, y_pred))
print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))

MAE: 0.02186035893180384
MSE: 0.00097018137793944
RMSE: 0.03917915164147886
MSE of 0.00091 looks decent.
```

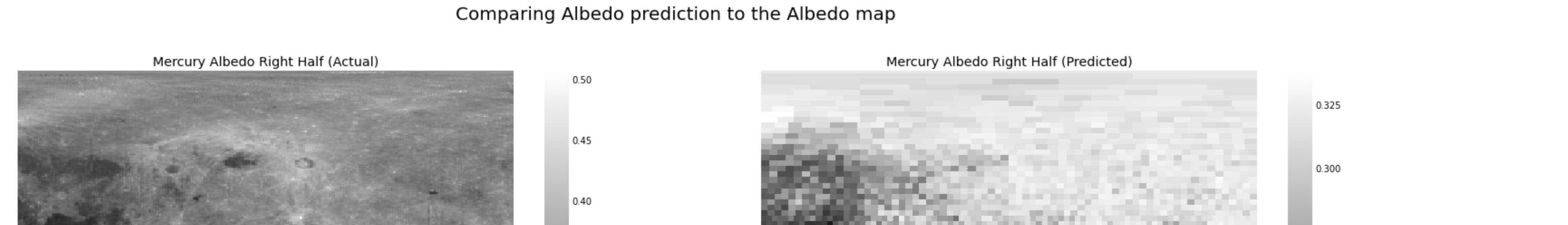
```
In [30]: RIGHT_PREDICTED = y_pred.reshape(NUM_ROWS, CENTER_HALF)
RIGHT_ACTUAL = albedo_map.iloc[:, :CENTER_HALF]
```

We now compare the albedo prediction to the albedo map.



Overall, the (stacked) model performed exceedingly well and was able to re-create a considerable portion of craters for the albedo map corresponding to the right hand side. Although, the pixel quality can be better constructed. We think this can be achieved by integrating geographical features, such as lat/long co-ordinates, surface level information, and so on.

Residuals, in the context of regression models, are the difference between the observed value (y) and the predicted value (ŷ), i.e. the error of the prediction. The residuals plot shows the difference between residuals on the vertical axis and the dependent variable on the horizontal axis, allowing you to detect regions within the target that may be susceptible to more or less error.



As we can infer from the above image, the residuals are roughly in the range ± 0.10.

