# Seeing beyond the visible: Estimating soil parameters from hyperspectral images - a starter pack

November 19, 2021

## 1 Introduction

How to open and understand the dataset

```
[1]: import pandas as pd
     import numpy as np
     import matplotlib.pyplot as plt
```

### 1.1 Basic information

1. Hyperspectral data:
    1. `hsi_path` contains path to hyperspectral masked numpy arrays containing hyperspectral data that underwent masking (i.e., the field area is masked, whereas all irrelevant areas are not masked)
    2. The name of the file (e.g., *'1989.npz'*) refers to the index of the corresponding training sample in the ground-truth table.
2. Ground-truth data:
    1. `gt_path` contains path to ground truth .csv file.
    2. Additionally, `wavelength_path` contains the mapping between a band number and the corresponding wavelength.

```
[2]: hsi_path = 'train_data/1570.npz'
     gt_path = 'train_gt.csv'
     wavelength_path = 'wavelengths.csv'
```

```
[3]: gt_df = pd.read_csv(gt_path)
     wavelength_df = pd.read_csv(wavelength_path)
```

### 1.2 Ground-truth description

`gt_df` contains:

1. `sample_index` - a reference to the numpay array containing the corresponding hyperspectral data.
2. P (for simplicity, we use P while referring to P_2O_5), K, Mg, pH - soil properties levels based on laboratory measurements.

```
[4]: gt_df[gt_df['sample_index']==1570]
```

```
[4]:        sample_index     P       K       Mg     pH
      1570          1570   47.9   165.0   181.0   6.8
```
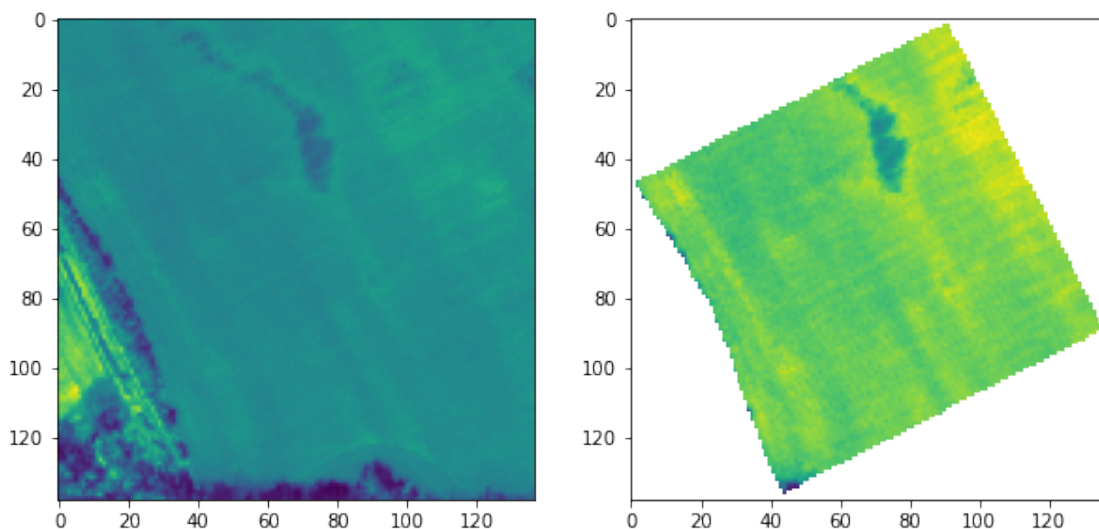
## 1.3 Displaying one hyperspectral band

```
[5]: fig, axs = plt.subplots(1, 2, figsize=(10, 5))
     band_id = 100
     wavelength = wavelength_df.loc[band_id-1]

     with np.load(hsi_path) as npz:
         arr = np.ma.MaskedArray(**npz)

     axs[0].imshow(arr[band_id,:,:].data)
     axs[1].imshow(arr[band_id,:,:])
     plt.suptitle(f'Representation of band {int(wavelength["band_no"])}␣
       ↪({wavelength["wavelength"]} nm)', fontsize=15)
     plt.show()
```

Representation of band 100 (778.54 nm)



## 1.4 Displaying the aggregated spectral curve for a field

```
[6]: fig = plt.figure(figsize=(10, 5))

     masked_scene_mean_spectral_reflectance = [arr[i,:,:].mean() for i in range(arr.
       ↪shape[0])]
```
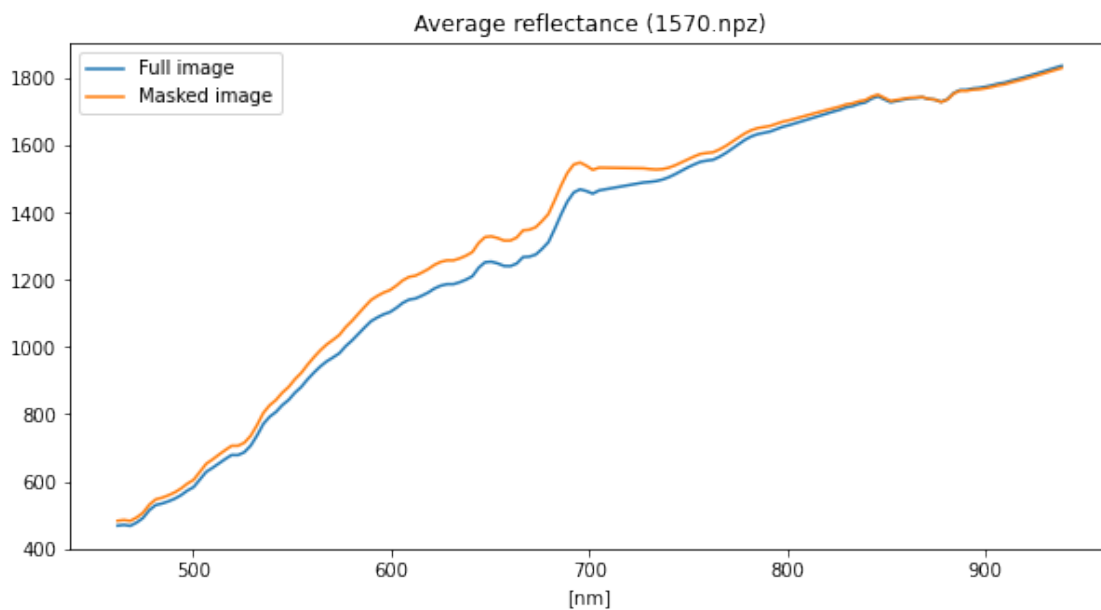
```
full_scene_mean_spectral_reflectance = [arr[i,:,:].data.mean() for i in
 →range(arr.shape[0])]

plt.plot(wavelength_df['wavelength'], full_scene_mean_spectral_reflectance,
 →label='Full image')
plt.plot(wavelength_df['wavelength'], masked_scene_mean_spectral_reflectance,
 →label='Masked image')

plt.xlabel('[nm]')
plt.legend()
plt.title(f'Average reflectance ({hsi_path.split("/")[-1]})')
plt.show()
```



## 2 Generating baseline solution

```
[7]: class BaselineRegressor:
         """
         Baseline regressor, which calculates the mean value of the target from the
     →training
         data and returns it for each testing sample.
         """
         def __init__(self):
             self.mean = 0

         def fit(self, X_train: np.ndarray, y_train: np.ndarray):
```

```python
        self.mean = np.mean(y_train, axis=0)
        self.classes_count = y_train.shape[1]
        return self

    def predict(self, X_test: np.ndarray):
        return np.full((len(X_test), self.classes_count), self.mean)


class SpectralCurveFiltering():
    """
    Create a histogram (a spectral curve) of a 3D cube, using the merge_function
    to aggregate all pixels within one band. The return array will have
    the shape of [CHANNELS_COUNT]
    """

    def __init__(self, merge_function = np.mean):
        self.merge_function = merge_function

    def __call__(self, sample: np.ndarray):
        return self.merge_function(sample, axis=(1, 2))
```

## 2.1   Load the data

```python
[8]:  import os
      from glob import glob

      def load_data(directory: str):
          """Load each cube, reduce its dimensionality and append to array.

          Args:
              directory (str): Directory to either train or test set
          Returns:
              [type]: A list with spectral curve for each sample.
          """
          data = []
          filtering = SpectralCurveFiltering()
          all_files = np.array(
              sorted(
                  glob(os.path.join(directory, "*.npz")),
                  key=lambda x: int(os.path.basename(x).replace(".npz", "")),
              )
          )
          for file_name in all_files:
              with np.load(file_name) as npz:
                  arr = np.ma.MaskedArray(**npz)
              arr = filtering(arr)
              data.append(arr)
```

```python
        return np.array(data)


def load_gt(file_path: str):
    """Load labels for train set from the ground truth file.
    Args:
        file_path (str): Path to the ground truth .csv file.
    Returns:
        [type]: 2D numpy array with soil properties levels
    """
    gt_file = pd.read_csv(file_path)
    labels = gt_file[["P", "K", "Mg", "pH"]].values
    return labels


X_train = load_data("train_data")
y_train = load_gt("train_gt.csv")
X_test = load_data("test_data")

print(f"Train data shape: {X_train.shape}")
print(f"Test data shape: {X_test.shape}")
```

```
Train data shape: (1732, 150)
Test data shape: (1154, 150)
```

## 2.2 Make predictions and generate submission file

```python
[9]: baseline_reg = BaselineRegressor()
     baseline_reg = baseline_reg.fit(X_train, y_train)
     predictions = baseline_reg.predict(X_test)

     submission = pd.DataFrame(data = predictions, columns=["P", "K", "Mg", "pH"])
     submission.to_csv("submission.csv", index_label="sample_index")
```

## 2.3 Calculating the metric

For the purpose of presenting the final metric calculation, we will extract a small *test_set* from the training set.

```python
[10]: X_test = X_train[1500:]
      y_test = y_train[1500:]

      X_train_new = X_train[:1500]
      y_train_new = y_train[:1500]

      # Fit the baseline regressor once again on new training set
      baseline_reg = baseline_reg.fit(X_train_new, y_train_new)
```

```python
baseline_predictions = baseline_reg.predict(X_test)

# Generate baseline values to be used in score computation
baselines = np.mean((y_test - baseline_predictions) ** 2, axis=0)



# Generate random predictions, different from baseline predictions
np.random.seed(0)
predictions = np.zeros_like(y_test)
for column_index in range(predictions.shape[1]):
    class_mean_value = baseline_reg.mean[column_index]
    predictions[:, column_index] = np.random.uniform(low=class_mean_value -␣
 ↪class_mean_value * 0.05,

                                                     high=class_mean_value +␣
 ↪class_mean_value * 0.05,

                                                     size=len(predictions))

# Calculate MSE for each class
mse = np.mean((y_test - predictions) ** 2, axis=0)

# Calculate the score for each class individually
scores = mse / baselines

# Calculate the final score
final_score = np.mean(scores)

for score, class_name in zip(scores, ["P", "K", "Mg", "pH"]):
    print(f"Class {class_name} score: {score}")

print(f"Final score: {final_score}")
```

```
Class P score: 0.9896068600445717
Class K score: 1.004900913045855
Class Mg score: 1.0228518828521695
Class pH score: 1.6431314552511207
Final score: 1.1651227777984292
```