

Sowing Success: How Machine Learning Helps Farmers Select the Best Crops



Measuring essential soil metrics such as nitrogen, phosphorous, potassium levels, and pH value is an important aspect of assessing soil condition. However, it can be an expensive and time-consuming process, which can cause farmers to prioritize which metrics to measure based on their budget constraints.

Farmers have various options when it comes to deciding which crop to plant each season. Their primary objective is to maximize the yield of their crops, taking into account different factors. One crucial factor that affects crop growth is the condition of the soil in the field, which can be assessed by measuring basic elements such as nitrogen and potassium levels. Each crop has an ideal soil condition that ensures optimal growth and maximum yield.

A farmer reached out to you as a machine learning expert for assistance in selecting the best crop for his field. They've provided you with a dataset called `soil_measures.csv`, which contains:

- `"N"` : Nitrogen content ratio in the soil
- `"P"` : Phosphorous content ratio in the soil
- `"K"` : Potassium content ratio in the soil
- `"pH"` value of the soil
- `"crop"` : categorical values that contain various crops (target variable).

Each row in this dataset represents various measures of the soil in a particular field. Based on these measurements, the crop specified in the `"crop"` column is the optimal choice for that field.

In this project, you will build multi-class classification models to predict the type of `"crop"` and identify the single most importance feature for predictive performance.

```
# All required libraries are imported here for you.
import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn import metrics
import matplotlib.pyplot as plt

#1. Read the data into a pandas DataFrame and perform exploratory data analysis

# Load the dataset
crops = pd.read_csv("soil_measures.csv")

# To perform a LogisticRegression the data have to be numerical, with no missing
data

# EDA:
#print(crops.head())
#print(crops.describe())
#print(crops["crop"].value_counts())
#print(crops.isna().sum()) # no missing values
#print(crops.info())

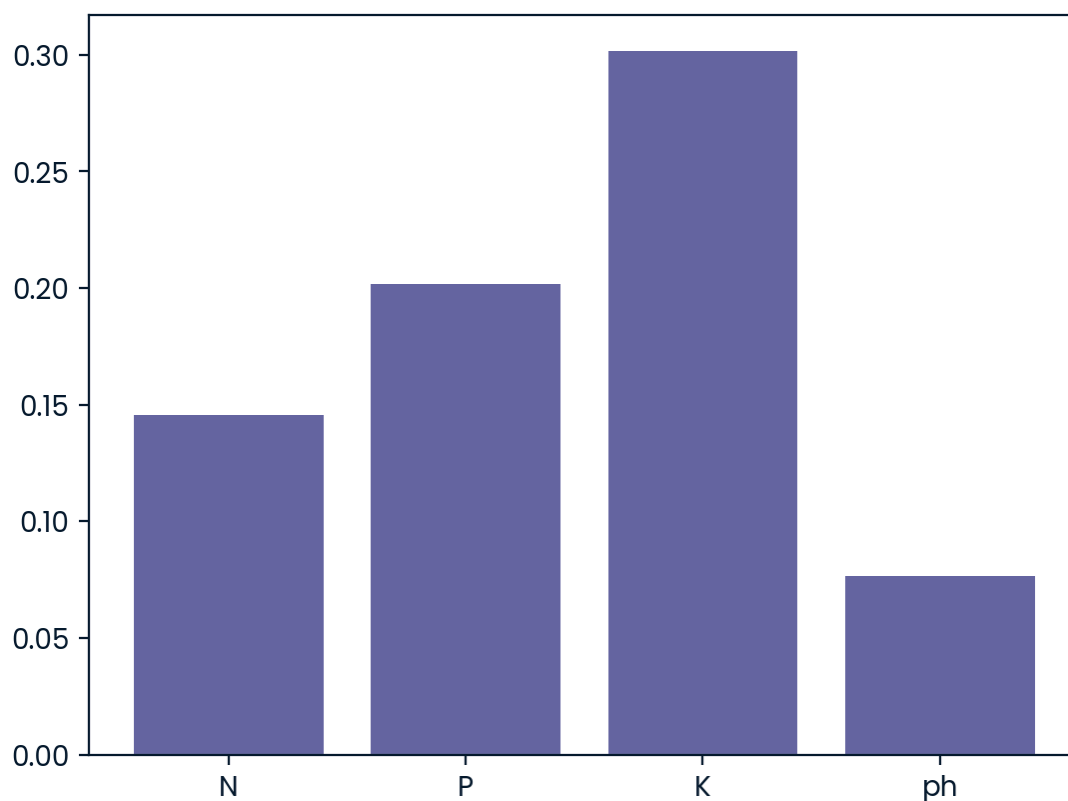
# so scikit-learn handels automatically categorical variable for target but not for
features
# that's why no need to use one hot encoding here
```

#2. Split the data / #3. Evaluate feature performance

```
y = crops["crop"].values
features = crops.drop("crop",axis = 1).columns

results = dict()
for f in features:
    X_f = crops[f].values.reshape(-1,1)
    model = LogisticRegression()
    X_train,X_test,y_train,y_test = train_test_split(X_f,y, random_state = 12)
    model.fit(X_train,y_train)
    y_pred = (model.predict(X_test))
    score = model.score(X_test,y_test)
    results[f] = score

plt.bar(features,results.values())
plt.show()
```



#4. Create the best_predictive_feature variable

```
best_predictive_feature = {"K":results["K"]}
print(best_predictive_feature)
```

```
{'K': 0.3018181818181818}
```