

estimation_of_probability

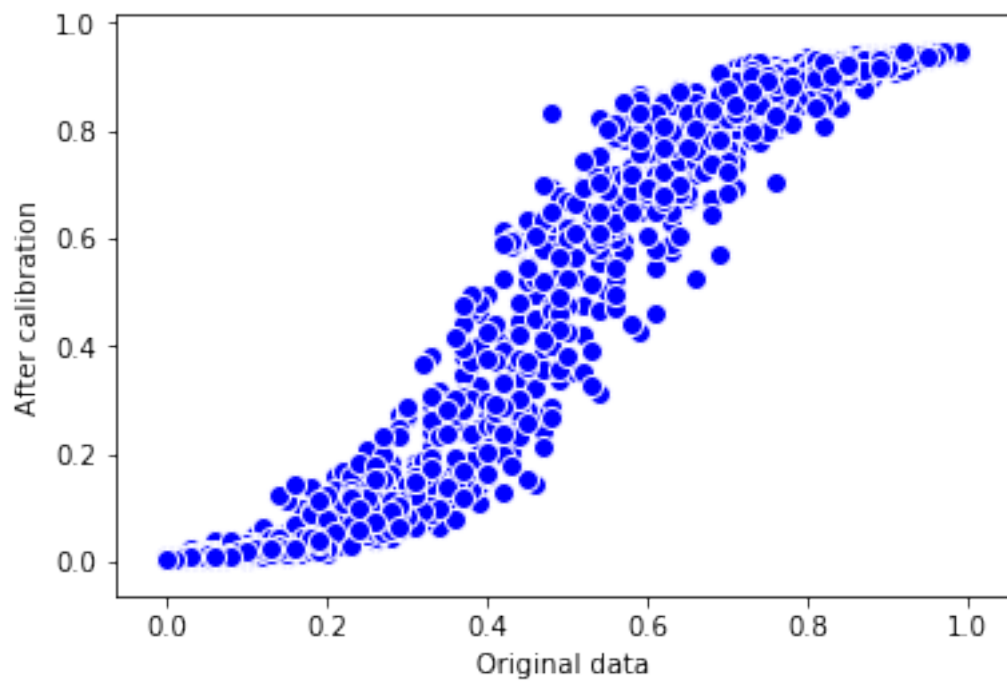
March 10, 2019

```
In [1]: import pandas as pd
import matplotlib as plt
from sklearn.calibration import CalibratedClassifierCV, calibration_curve
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestClassifier
import pickle
import numpy as np
covertypes_dataset = pickle.load(open('covertypes_dataset.pickle','rb'))
hypothesis = RandomForestClassifier(n_estimators=100, random_state=101)
calibration = CalibratedClassifierCV(hypothesis, method='sigmoid', cv=5) # maps the re
covertypes_X = covertypes_dataset.data[:15000,:]
covertypes_Y = covertypes_dataset.target[:15000]
covertypes_test_X = covertypes_dataset.data[15000:25000,:]
covertypes_test_Y = covertypes_dataset.target[15000:25000]

D:\Python\Lib\importlib\_bootstrap.py:219: RuntimeWarning: numpy.ufunc size changed, may indic
return f(*args, **kwargs)
D:\Python\Lib\importlib\_bootstrap.py:219: RuntimeWarning: numpy.ufunc size changed, may indic
return f(*args, **kwargs)
D:\Python\Lib\importlib\_bootstrap.py:219: RuntimeWarning: numpy.ufunc size changed, may indic
return f(*args, **kwargs)

In [2]: hypothesis.fit(covertypes_X, covertypes_Y)
calibration.fit(covertypes_X, covertypes_Y)
prob_raw = hypothesis.predict_proba(covertypes_test_X)
prob_cal = calibration.predict_proba(covertypes_test_X)

In [6]: %matplotlib inline
covertypes = ['Spruce/Fir', 'Lodgepole Pine', 'Ponderosa Pine', 'Cottonwood/Willow', 'Asp
tree_kind = covertypes.index('Ponderosa Pine')
probs = pd.DataFrame(list(zip(prob_raw[:,tree_kind], prob_cal[:,tree_kind])),
                      columns=['Original data', 'After calibration'])
plot = probs.plot(kind='scatter', x=0, y=1, s=64, c='blue', edgecolors='white')
```



In []: