

Solutions 3

Jumping Rivers

During the lecture we fit a logistic regression model to the breast cancer data for classifying tumors in patients. We are going to fit a KNN classifier to the same data.

- Construct the pipeline ready for fitting the model

```
from sklearn.datasets import load_breast_cancer
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.neighbors import KNeighborsClassifier
```

```
cancer = load_breast_cancer()
X_train, y_train = cancer.data, cancer.target
```

```
model = Pipeline([
    ('pre', StandardScaler()),
    ('model', KNeighborsClassifier())
])
```

- We want to find the best value of K for the classifier when optimising for recall, our motivation is that we want to correctly identify as many of the malignant tumours as possible. Start with a grid search over $k = [1, 5, 10, 20, 50, 100]$

```
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import make_scorer, recall_score
```

```
def recall(y_true, y_pred):
    return recall_score(y_true, y_pred, pos_label=0)
```

```
rec = make_scorer(recall)
```

```
clf = GridSearchCV(model, param_grid={
    'model__n_neighbors': [1, 5, 10, 20, 50, 100]
}, cv=10, iid=False, return_train_score=False,
    scoring=rec)
clf.fit(X_train, y_train)
```

```
## GridSearchCV(cv=10, error_score='raise-deprecating',
##              estimator=Pipeline(memory=None,
```

```

##             steps=[('pre',
##                     StandardScaler(copy=True,
##                                     with_mean=True,
##                                     with_std=True)),
##                     ('model',
##                     KNeighborsClassifier(algorithm='auto',
##                                         leaf_size=30,
##                                         metric='minkowski',
##                                         metric_params=None,
##                                         n_jobs=None,
##                                         n_neighbors=5, p=2,
##                                         weights='uniform'))],
##             verbose=False),
##         iid=False, n_jobs=None,
##         param_grid={'model__n_neighbors': [1, 5, 10, 20, 50, 100]},
##         pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
##         scoring=make_scorer(recall), verbose=0)

```

- Create a plot of the K parameter against the average recall score found in the cross validation grid search

```
import pandas as pd
```

```
output = pd.DataFrame(clf.cv_results_)[['param_model__n_neighbors', 'mean_test_score']]
```

```

import seaborn as sns
import matplotlib.pyplot as plt
plt.figure()
sns.lineplot(x='param_model__n_neighbors', y='mean_test_score', data=output)
plt.show()

```

- What region of K looks like it will give the best value?

```
## for me it is between 1 and 20
```

- Re-run your grid search across that region

```

clf = GridSearchCV(model, param_grid={
    'model__n_neighbors': list(range(1, 21))
}, cv=10, iid=False, return_train_score=False,
   scoring=rec)
clf.fit(X_train, y_train)

## GridSearchCV(cv=10, error_score='raise-deprecating',
##             estimator=Pipeline(memory=None,
##                                 steps=[('pre',

```

```

##                                     StandardScaler(copy=True,
##                                     with_mean=True,
##                                     with_std=True)),
##                                ('model',
##                                KNeighborsClassifier(algorithm='auto',
##                                                    leaf_size=30,
##                                                    metric='minkowski',
##                                                    metric_params=None,
##                                                    n_jobs=None,
##                                                    n_neighbors=5, p=2,
##                                                    weights='uniform'))],
##                                verbose=False),
##                                iid=False, n_jobs=None,
##                                param_grid={'model__n_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10,
##                                                                    11, 12, 13, 14, 15, 16, 17, 18,
##                                                                    19, 20]},
##                                pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
##                                scoring=make_scorer(recall), verbose=0)

```

- What is the best parameter choice and the corresponding recall score?

```
clf.best_params_
```

```
## {'model__n_neighbors': 4}
```

```
clf.best_score_
```

```
## 0.9480519480519481
```

- Is this better than the Logistic regression in the notes?

```
## for me it is worse
```

Other techniques

- Try the other classification algorithms that we explored, can you find any that perform better?

Response Optimisation

In this question we aim to use machine learning as a route to optimising a response. In particular, the data we have at hand are mixtures of concrete together with their measured compressive strength. We want to use the available data to propose new formulations of concrete that might be better than those used in the experiment.

The data can be loaded as

```

import jrpyml
concrete = jrpyml.datasets.mixtures.load_data()

• This isn't a traditional experimental design set up, there are some
  mixture formulations that have repeated measures. A sensible
  thing to do in situations like this is to average over the repeated
  responses.

concrete = concrete.groupby([
    item for item in concrete.columns if item != 'CompressiveStrength'
]).agg({
    'CompressiveStrength': 'mean'
}).reset_index()

• Try any number of models aiming to find one that gives good pre-
  dictive performance

from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error, make_scorer, r2_score
from sklearn.model_selection import cross_validate
import pandas as pd

model = RandomForestRegressor(n_estimators=1000)
y_train = concrete['CompressiveStrength']
X_train = concrete.drop(['CompressiveStrength'], axis=1).values

model.fit(X_train, y_train)

## RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
##                          max_features='auto', max_leaf_nodes=None,
##                          min_impurity_decrease=0.0, min_impurity_split=None,
##                          min_samples_leaf=1, min_samples_split=2,
##                          min_weight_fraction_leaf=0.0, n_estimators=1000,
##                          n_jobs=None, oob_score=False, random_state=None,
##                          verbose=0, warm_start=False)

score_fn = make_scorer(mean_squared_error)
scores = cross_validate(model, X_train, y_train,
    scoring = {
        'mse': score_fn,
    },
    cv = 10)
pd.DataFrame(scores).mean()

## fit_time      2.652479
## score_time    0.047236
## test_mse      51.976108
## dtype: float64

```

- Once you have found a model you are happy with, we can aim to use it to optimise mixtures of concrete at different ages. This optimisation problem also has some constraints, namely that the sum of all components must be equal to 1 (as we have proportions of a mixture). Further no input can be less than 0 or greater than 1. One way to impose these constraints is to define an objective function in terms of 6 of the mixture components and infer the final one. Such an objective function might be defined as below. We return the negative model prediction here as our optimisation routine will look to minimise the function.

```
import numpy as np
```

```
def obj(x, age=28):
    # add some constraints, # index 7 is age
    # sum of components excluding age must be 1
    # there are 7 components in total

    if np.any(x < 0) | np.any(x > 1):
        return 1e50
    x = np.append(x, [1-x.sum()], age)
    return -1*model.predict(x.reshape(1,-1))
```

- Given an objective function to minimize, **scipy** has a number of routines to do so. We try to minimize the function from a number of random start points, the aim being to have a better chance at finding a global minimum.

```
from scipy.optimize import minimize
import random
```

```
## 20 random rows of data to start from
start_idx = [random.randint(0,X_train.shape[0]) for _ in range(20)]

outputs = [minimize(obj, X_train[i,:-2], method = 'Nelder-Mead') for i in start_idx]

## Extract the results of the objective function
## once the optimisation routine completes.
np.array([o.fun.flatten()[0] for o in outputs])

## For a given output we might also extract the input array
## (The mixture of concrete in our case)

## array([-67.30419333, -54.034535 , -40.84854 , -60.458135 ,
##        -49.25199167, -55.05255167, -34.87037 , -68.31057333,
##        -40.46383 , -57.993445 , -70.04689833, -68.03809 ,
```

```
##      -71.84044833, -47.13703    , -41.51324    , -73.38313    ,  
##      -27.67683    , -71.16442667, -59.72798    , -72.00106833])  
  
outputs[0].x  
  
## array([0.11716028, 0.11926713, 0.00145396, 0.02872217, 0.00513883,  
##        0.41879427])
```