

lab_wine_partial_finish

May 2, 2019

1 Lab: Hyper-Parameter Optimization with PCA

PCA is often applied as a pre-processing step with classifiers. When using PCA in this manner, one must select the number of PC components to use along with parameters in classifier. In this lab, we will demonstrate how to performing this *hyper-parameter optimization*. In doing the lab, you will learn to:

- Combine PCA with data scaling.
- Compute and visualize PC components
- Select the number of PCs with K-fold cross validation
- Implement the multi-stage classifier pipeline in sklearn
- Perform automatic parameter search using GridSearchCV in combination with a pipeline.

We first download the basic packages.

```
In [19]: import numpy as np
import matplotlib
import matplotlib.pyplot as plt
```

1.1 Downloading the Data

We will use a very simple wine dataset, commonly used in teaching machine learning class. The problem is to classify the type of red wine from features of the wine such as the alcohol and other chemical components. There are three possible wine types.

```
In [20]: from sklearn.datasets import load_wine
from sklearn.model_selection import KFold
data = load_wine()

# TODO print the features names in data.feature_names and data.target_names
print("feature names:", data.feature_names)
print("target names:", data.target_names)
```

```
feature names: ['alcohol', 'malic_acid', 'ash', 'alcalinity_of_ash', 'magnesium', 'total_phenols',
target names: ['class_0' 'class_1' 'class_2']
```

Get the data matrix X from `data.data` and the target values y from `data.target`. Print the number of samples, number of features and number of classes.

```
In [21]: # TODO
X = data.data
y = data.target
print("number of samples = {}".format(X.shape[0]))
print("number of features = {}".format(X.shape[1]))
print("number of classes = {}".format(np.unique(y).shape[0]))

number of samples = 178
number of features = 13
number of classes = 3
```

1.2 Perform PCA for Visualization

Before performing PCA, you should scale the data matrix to remove the mean and normalize the variance of the different components. For this purpose, create a `StandardScaler` object scaling. Then fit the scaling with the entire data `X`. Transform the data and let `Xs` be the scaled data.

```
In [22]: from sklearn.preprocessing import StandardScaler

# TODO
scaling = StandardScaler()
scaling.fit(X)
Xs = scaling.transform(X)
```

Now, fit a PCA on the scaled data matrix `Xs`. You can use the `sklearn` PCA method. In order that we can visualize the results set `n_components=2`. Select `svd_solver='randomized'` and `whiten=True`. Use the `pca.transform` method to find, `Z`, the coefficients of `Xs` in the PCA basis.

```
In [23]: from sklearn.decomposition import PCA

# TODO
ncomp = 2
# Construct the PCA object
pca = PCA(n_components = ncomp,svd_solver='randomized', whiten=True)
pca.fit(Xs)
Z = pca.transform(Xs)
# print(Z)
```

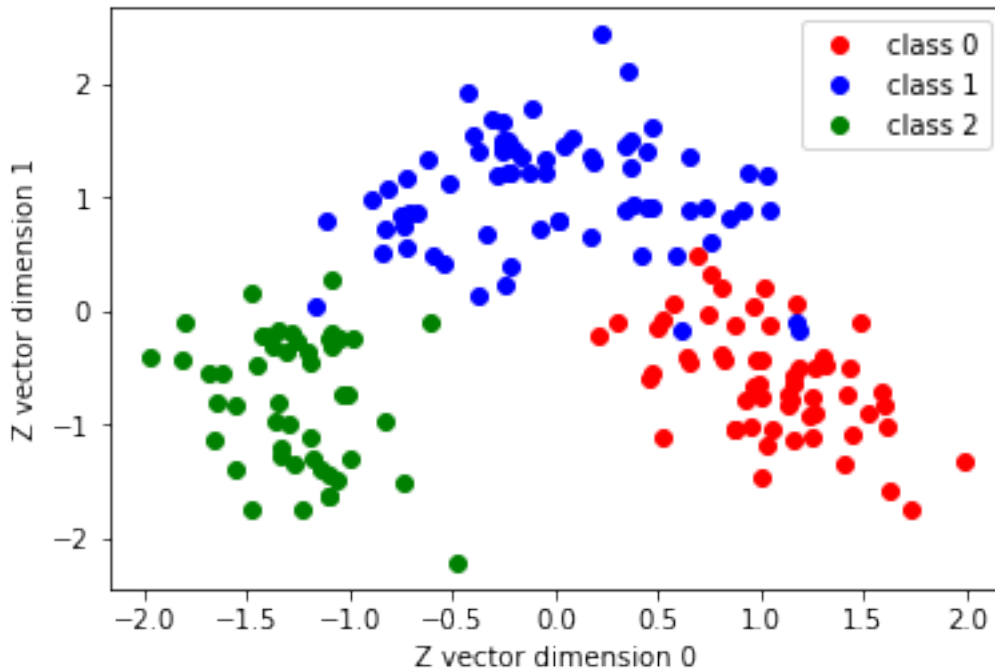
In the transformed basis, each data sample is represented by a two dimensional vector, `Z[i,0]`, `Z[i,1]`. Plot a scatter plot of the transformed data. Use different marker colors for the different classes. If you did everything, you should see that the classes are quite well separated with even two PCs.

```
In [24]: # TODO
class0_index = np.where(y==0)[0]
class1_index = np.where(y==1)[0]
class2_index = np.where(y==2)[0]
plt.plot(Z[class0_index,0],Z[class0_index,1], 'ro')
```

```

plt.plot(Z[class1_index,0],Z[class1_index,1],'bo')
plt.plot(Z[class2_index,0],Z[class2_index,1],'go')
plt.xlabel("Z vector dimension 0")
plt.ylabel("Z vector dimension 1")
plt.legend(["class 0", "class 1", "class 2"])
plt.show()

```



Now, refit the scaled data `Xs` using `n_components=nfeatures` where `nfeatures` is the number of features. This is the maximum number of PCs. Get the singular values from `pca.singular_values_` and plot the portion of variation as a function of the number of PCs. The PoV for using `n` PCs is:

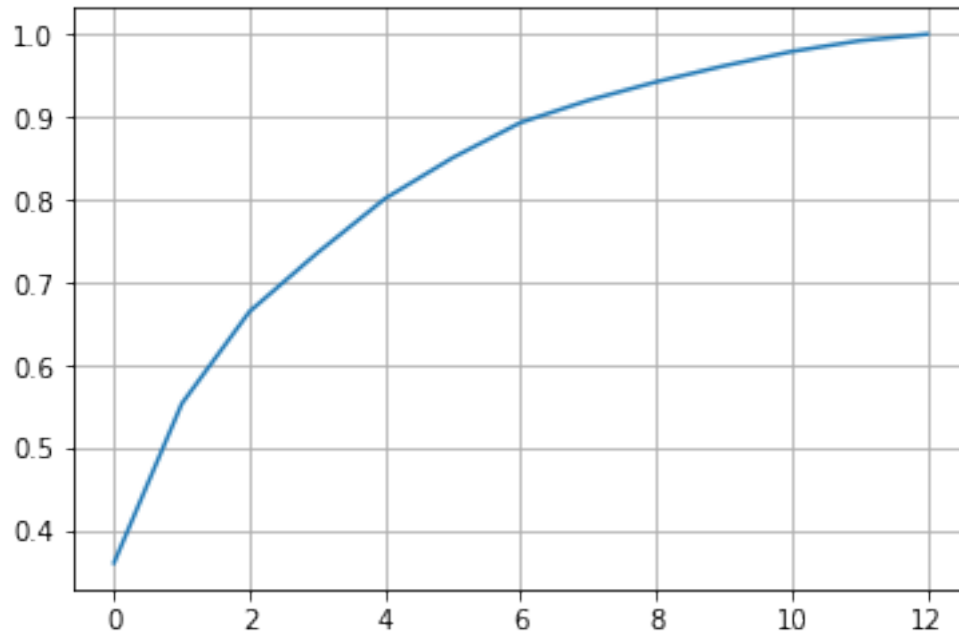
$$\text{PoV}[n] = \frac{\sum_{i=0}^{n-1} s[i]**2}{\sum_{i=0}^{d-1} s[i]**2}$$

where `s[i]` is the `i`-th singular value and `d` is the number of features. You should see that the 4 PCs contains more than 70% of the variance.

```

In [25]: # TODO
nfeatures = Xs.shape[1]
pca = PCA(n_components = nfeatures,svd_solver='randomized', whiten=True)
pca.fit(Xs)
s = pca.singular_values_
lam = s**2
pov = np.cumsum(lam) / np.sum(lam)
plt.plot(pov)
plt.grid()
plt.show()

```



1.3 Using PCA with Classification

We will now use data scaling and PCA as a pre-processing step for logistic classification. The number of PCs to use can be found with cross-validation. Complete the code below which tries different number of PCs components to use and measures the test accuracy for each value.

```
In [26]: from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
nfold = 5

# Create a K-fold object
kf = KFold(n_splits=nfold)
kf.get_n_splits(X)

# Number of PCs to try
ncomp_test = np.arange(2,12)
num_nc = len(ncomp_test)

# Accuracy: acc[icomp,ifold] is test accuracy when using `ncomp = ncomp_test[icomp]
acc = np.zeros((num_nc,nfold))

# Loop over number of components to test
for icomp, ncomp in enumerate(ncomp_test):

    # Look over the folds
    for ifold, I in enumerate(kf.split(X)):
```

```

Itr, Its = I

# TODO: Split data into training
Xtr = X[Itr,:]
Xts = X[Its,:]
ytr = y[Itr]
yts = y[Its]

# TODO: Create a scaling object and fit the scaling on the training data
scaling = StandardScaler()
scaling.fit(Xtr)
Xtrs = scaling.transform(Xtr)
# TODO: Fit the PCA on the scaled training data
pca = PCA(n_components=ncomp ,svd_solver='randomized', whiten=True)
pca.fit(Xtrs)
Ztr = pca.transform(Xtrs)
# TODO: Train a classifier on the transformed training data

# Use a logistic regression classifier
logreg = LogisticRegression(multi_class='auto', solver='lbfgs')
logreg.fit(Ztr, ytr)

# TODO: Transform the test data through data scaler and PCA
Xtss = scaling.transform(Xts)
Zts = pca.transform(Xtss)

# TODO: Predict the labels the test data
yhat = logreg.predict(Zts)
# TODO: Measure the accuracy
acc[icomp, ifold] = np.mean(yhat == yts)

```

Use the `plt.errorbar` function to plot the mean accuracy with error bars corresponding to the standard error of the accuracy as a function of the number of components. Find the optimal number of PCs to use according to the normal rule and one SE rule. If you did it correctly, you should get an accuracy of around 96%.

```

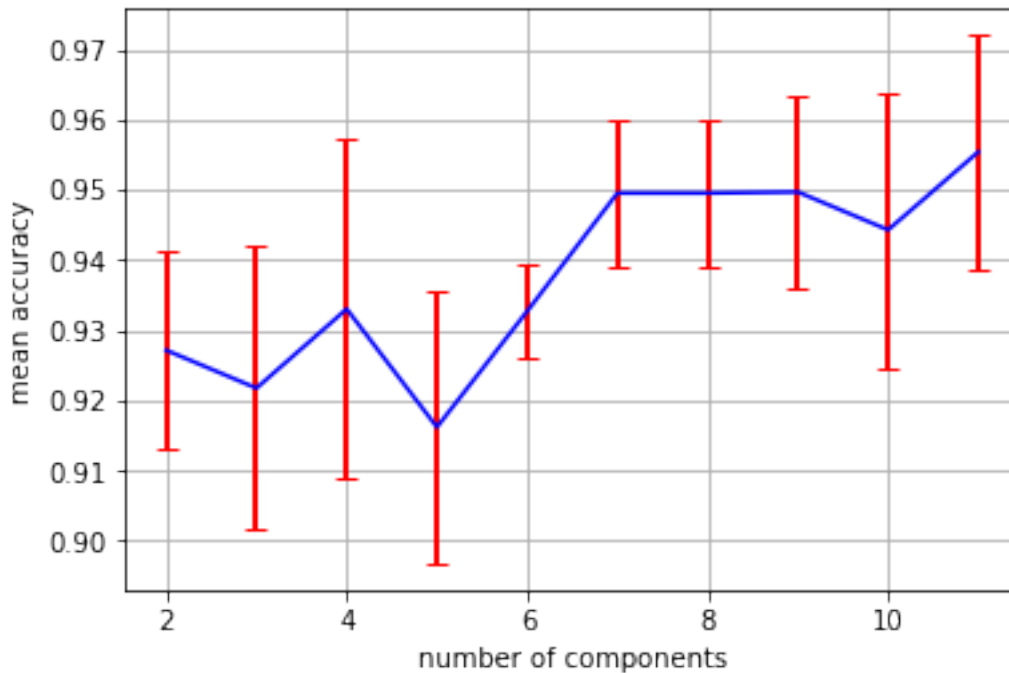
In [27]: # TODO:
acc_mean = np.mean(acc,1)
acc_se = np.std(acc, 1) / np.sqrt(nfold - 1)
plt.errorbar(ncomp_test, acc_mean, acc_se, ecolor='r', color='b', elinewidth=2, capsize=4)
plt.xlabel("number of components")
plt.ylabel("mean accuracy")
plt.grid()
plt.show()
# TODO: Optimal order with the normal rule
i = np.argmax(acc_mean)
opt_order = ncomp_test[i]
print("The optimal order with the normal rule:", opt_order)

```

```

# TODO: Optimal order with one SE rule
acc_tgt = acc_mean[i]-acc_se[i]
# print(acc_tgt)
index = np.where(acc_mean > acc_tgt)[0][0]
opt_order = ncomp_test[index]
print("The optimal order with one SE rule:", opt_order)

```



The optimal order with the normal rule: 11
 The optimal order with one SE rule: 7

1.4 Hyper-Parameter Optimization with GridCV.

We will now try a more complex classifier -- a support vector classifier with a radial basis function. When we use such a classifier, there will be a number of parameters to tune. When the number of parameters to tune becomes large, writing a loop over multiple parameters as we did above becomes cumbersome. The sklearn package has a very nice routine, `GridSearchCV` to perform this sort of parameter search.

Before, we do this we need to create an estimator Pipeline. An estimator pipeline is a sequence of transformations followed by an estimator that will operate on the transformed data. Create the following pipeline:

- Create a `StandardScaler()` object called `scaler` for the first transformation
- Create a `PCA()` object called `pca` for the second transformation
- Create a `SVC()` object called `svc` for the final SVM classifier. Set the parameter `kernel='rbf'`.

Once you have the three steps defined, you can create the pipeline with the command:

```
pipe = Pipeline(steps=[('scaler', scaler), ('pca', pca), ('svc', svc)])
```

```
In [28]: from sklearn.pipeline import Pipeline
         from sklearn.model_selection import GridSearchCV
         from sklearn.svm import SVC

         # TODO
         scaler = StandardScaler()
         pca = PCA()
         svc = SVC(kernel='rbf')
         pipe = Pipeline(steps=[('scaler', scaler), ('pca', pca), ('svc', svc)])
```

We next define all the parameters that we want to search over. Define the following arrays:

- `ncomp_test`: values from 3 to 10 representing number of PCs to test
- `C_test`: values of C in the SVC to test. Use 10^{-2} , 10^{-1} , ... , 10^3
- `gam_test`: values of γ in the SVC to test. Use 10^{-3} , 10^{-2} , ... , 10^1

```
In [29]: # TODO
         ncomp_test = np.arange(3,11)
         c_test = 10**np.arange(-2,4, dtype=float)
         gam_test = 10**np.arange(-3,2, dtype=float)
```

Next, we create a dictionary params of the form:

```
params = {'pca__n_components': ncomp_test, 'svc__C' : c_test, ...}
```

Each key in the dictionary is the of the form `estimator__param` and the value is the values to be tested.

```
In [30]: # TODO
         params = {'pca__n_components': ncomp_test, 'svc__C' : c_test, "svc__gamma": gam_test}
```

Finally, an object `estimator = GridSearchCV(...)` from pipe and params. Set `cv=5`, `train_score=True` and `iid=False`. Fit the estimator from the data X, y . Then the estimator will perform the cross-validation over all the parameters. This may take a minute since we are search over so many parameters.

```
In [31]: # TODO
         estimator = GridSearchCV(pipe, params, cv=5, return_train_score=True, iid=False)
         estimator.fit(X,y)
```

```
Out[31]: GridSearchCV(cv=5, error_score='raise-deprecating',
                      estimator=Pipeline(memory=None,
                      steps=[('scaler', StandardScaler(copy=True, with_mean=True, with_std=True)), ('pca', PCA(svd_solver='auto', tol=0.0, whiten=False)), ('svc', SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
                      decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
                      kernel='rbf', max_iter=-1, probability=False, random_state=None,
```

```

shrinking=True, tol=0.001, verbose=False))],
    fit_params=None, iid=False, n_jobs=None,
    param_grid={'pca__n_components': array([ 3,  4,  5,  6,  7,  8,  9, 10]), 'svc__C': 1.0, 'svc__gamma': 0.1},
    pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
    scoring=None, verbose=0)

```

Print the best test score and best parameters. They are fields in estimator. If you did it correctly, it should be a little higher than the logistic regression (about 0.97 to 0.98 accuracy).

```

In [32]: # TODO
print("best test score: ", estimator.best_score_)
print("best parameter dict:", estimator.best_params_)

```

```

best test score: 0.978078078078078
best parameter dict: {'pca__n_components': 5, 'svc__C': 1.0, 'svc__gamma': 0.1}

```

Finally, you can get the test score for all the parameter choices from

```
test_score = estimator.cv_results_['mean_test_score']
```

Use the `imshow` command to plot the mean test score over `gamma` and `C` for the value `n_components=5`.

```

In [33]: # TODO
test_score = estimator.cv_results_['mean_test_score']
sub_test_score = test_score[3*6*5:4*6*5].reshape(6,5)
plt.imshow(sub_test_score, cmap=plt.cm.jet)
plt.title("mean score over gamma and c when n_components={}".format(5))
plt.xlabel("gamma")
plt.ylabel("C")
plt.colorbar()
plt.show()

```


mean score over gamma and c when n_components=5

