

ML_HW4_Q1_NS

April 9, 2021

1 Question 1

1.1 Kernel Ridge Regression

1.1.1 a)

Apply kernelized ridge regression to the automobile mpg dataset. The training data and test data are provided in auto mpg train.csv and auto mpg test.csv, respectively. The first column is the mpg data while the other 7 columns are the features: 1. mpg: continuous 2. cylinders: multi-valued discrete 3. displacement: continuous 4. horsepower: continuous 5. weight: continuous 6. acceleration: continuous 7. model year: multi-valued discrete 8. origin: multi-valued discrete

We have normalized the feature data to the range [0,1]. Please apply the kernelized ridge regression to this dataset (use mpg as the target, the other 7 columns as features). Report the RMSE (Root Mean Square Error) of the models on the test data. Try (set lamda = 1).

```
[26]: #imports
import pandas as pd
import numpy as np
from sklearn.kernel_ridge import KernelRidge
from sklearn.metrics import accuracy_score, mean_squared_error
```

```
[27]: #we need to normlize the features between 0 and 1, this is the function to do
      ↪ SO

def NormalizeData(data):
    return (data - np.min(data)) / (np.max(data) - np.min(data))
```

```
[28]: #read in train and test data

te = pd.read_csv('auto_mpg_test.csv', sep = ' ')

tr = pd.read_csv('auto_mpg_train.csv', sep = ' ')

tes = pd.DataFrame(te).to_numpy()

tra = pd.DataFrame(tr).to_numpy()

#split features (x) and labels (mpg, y)
```

```

x_tra = tra[:, 1:8]
y_train = tra[:, 0]
x_tes = tes[:, 1:8]
y_test = tes[:, 0]

#normalize features (x_test and x_train)

x_test = NormalizeData(x_tes)
x_train = NormalizeData(x_tra)

print(x_train.shape, y_train.shape, x_test.shape, y_test.shape)

```

(299, 7) (299,) (97, 7) (97,)

Once we can you sklearn (THANK YOU!!!!!!), we can use the built in sklearn.kernel_ridge.KernelRidge function.

Thankfully, we can call the KernelRidge function with the kernel='rbf' instead of it's default linear kernel.

The rbf kernel uses the gaussian distribution asked in the question. See this link for documentation proof: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.pairwise.rbf_kernel.html The default rbf kernel sets sigma = 1, as asked in the question.

```

[29]: #Kernel Ridge Regression from sklearn using rbf (gaussian) kernel
skKRR = KernelRidge(alpha=1.0, kernel='rbf', gamma = 1/2)

#fit kernel to train data
skKRR.fit(x_train, y_train)

#predict test data
skKRR_y_pred = skKRR.predict(x_test)

#report root mean squared error
rmse = np.sqrt(mean_squared_error(y_true=y_test,y_pred=skKRR_y_pred))
print("Root Mean Squared Error for Gaussian Kernel: ",rmse)

```

Root Mean Squared Error for Gaussian Kernel: 3.2102263635542707

Out of curiosity (and because this is a short question due to sklearn), I want to try to see the difference between linear and rbf kernels.

```

[30]: #Kernel Ridge Regression from sklearn using liner kernel
skKRR_lin = KernelRidge(alpha=1.0)

```

```
#fit kernel to train data
skKRR_lin.fit(x_train, y_train)

#predict test data
skKRR_y_pred_lin = skKRR_lin.predict(x_test)

#report root mean squared error
rmse_lin = np.sqrt(mean_squared_error(y_true=y_test,y_pred=skKRR_y_pred_lin))
print("Root Mean Squared Error for Liner Kernel: ",rmse_lin)
```

Root Mean Squared Error for Liner Kernel: 5.207784717373817

Looks like the Gaussian Kernel is superior!