

Assignment 4: Kernels

CS 589 - ML

Shubham Shetty (shubhamshett@umass.edu)
Brinda Murulidhara (bmurulidhara@umass.edu)
Adarsh Kolya (akolya@umass.edu)

April 2021

Preface

Before running our codes, the following packages are imported and test & training data are assigned to variables. Also, some reusable functions are defined -

```
1 # Import Statements
2 import numpy as np
3 from scipy.special import comb
4 from matplotlib import pyplot as plt
5 from sklearn import svm
6 from sklearn.model_selection import KFold
7 from sklearn.metrics import hinge_loss
8 from sklearn.svm import SVC
9
10
11 # Load Data
12 ## 1D Data
13 data = np.load("data_synth.npz")
14 X_trn = data_synth['X_trn']           # 1D Input Training Dataset
15 Y_trn = data_synth['Y_train']        # 1D Output Training Dataset
16 X_val = data_synth['X_val']          # Test Input Dataset
17 Y_val = data_synth['Y_val']          # Test Output Dataset
18
19 ## Big data for SVM testing
20 data_real = np.load("data_real.npz")
21 X_trn_real = data_real['x_trn']       # 686 training inputs of length
22     4
23 Y_trn_real = data_real['y_trn']       # 686 training outputs of
24     length 4
25 X_val_real = data_real['x_tst']       # 686 test inputs
26
27 # Function to perform basis expansion
28 def get_poly_expansion(P):
29     def expand(X):
30         tmp = [np.sqrt(comb(P, p)) * X ** p for p in range(P + 1)]
31         return np.vstack(tmp).T
32     return expand
```

Answer 1

The objective is to find a vector w^* that minimizes :

$$\sum_{n=1}^N (w \cdot x^{(n)} - y^{(n)})^2 + \lambda \|w\|^2$$

Let's rewrite this in matrix form. Let y be a single vector containing all of our outputs:

$$y = (y^{(1)}, \dots, y^{(N)})$$

Consider X to be a matrix with N rows that is formed by $x^{(1)}, x^{(2)} \dots x^{(N)}$.

$$X = \begin{bmatrix} (x^{(1)})^\top \\ \vdots \\ (x^{(N)})^\top \end{bmatrix}$$

The above equation $R(w)$ can be rewritten as :

$$\begin{aligned} R(w) &= \|X \cdot w - y\|^2 + \lambda \|w\|^2 \\ &= (X \cdot w - y)^\top (X \cdot w - y) + \lambda w^\top w \end{aligned}$$

At w^* ,

$$\begin{aligned} \nabla R(w^*) &= 0 \\ 2X^\top (Xw^* - y) + 2\lambda w^* &= 0 \\ X^\top Xw^* + \lambda w^* &= X^\top y \\ (X^\top X + \lambda I) \cdot w^* &= X^\top y \end{aligned}$$

Hence,

$$w^* = (X^\top X + \lambda I)^{-1} X^\top y$$

The above can be rewritten as :

$$w^* = \left(\sum_{n=1}^N (x^{(n)})(x^{(n)})^\top + \lambda I \right)^{-1} \sum_{n=1}^N x^{(n)} y^{(n)}.$$

Answer 2

The objective is to find a vector w^* that minimizes :

$$\sum_{n=1}^N \left(w \cdot h(x^{(n)}) - y^{(n)} \right)^2 + \lambda \|w\|^2$$

Let's rewrite this in matrix form. Let y be a single vector containing all of our outputs:

$$y = (y^{(1)}, \dots, y^{(N)})$$

Consider H to be a matrix with N rows that is formed by $h(x^{(1)}), h(x^{(2)}) \dots h(x^{(N)})$.

$$H = \begin{bmatrix} h(x^{(1)})^\top \\ \vdots \\ h(x^{(N)})^\top \end{bmatrix}$$

The above equation $R(w)$ can be rewritten as :

$$\begin{aligned} R(w) &= \|H \cdot w - y\|^2 + \lambda \|w\|^2 \\ &= (H \cdot w - y)^\top (H \cdot w - y) + \lambda w^\top w \end{aligned}$$

At w^* ,

$$\begin{aligned} \nabla R(w^*) &= 0 \\ 2H^\top (Hw^* - y) + 2\lambda w^* &= 0 \\ H^\top Hw^* + \lambda w^* &= H^\top y \\ (H^\top H + \lambda I) \cdot w^* &= H^\top y \end{aligned}$$

Hence,

$$w^* = (H^\top H + \lambda I)^{-1} H^\top y$$

The above can be rewritten as :

$$w^* = \left(\sum_{n=1}^N (h(x^{(n)})) (h(x^{(n)}))^\top + \lambda I \right)^{-1} \sum_{n=1}^N h(x^{(n)}) y^{(n)}.$$

Answer 3

The objective is to represent $y^{\text{pred}} = w^* \cdot h(x^{\text{pred}})$ in terms of some kernel function k such that $k(x, x') = h(x) \cdot h(x')$ without referencing $h(\cdot)$ or w^* .

From previous solution, we know that

$$w^* = (H^\top H + \lambda I)^{-1} H^\top y$$

where H is a matrix with N rows that is formed by $h(x^{(1)}), h(x^{(2)}) \dots h(x^{(N)})$ and y is a single vector containing all of our outputs.

It is known that for given matrices P & Q ,

$$(PQ + I)^{-1}P = P(QP + I)^{-1}$$

Using this equality, we get

$$w^* = H^\top (HH^\top + \lambda I)^{-1} y$$

An element at position (n, m) in HH^\top is equal to $h(x^{(n)}) \cdot h(x^{(m)})$.

Let us define a matrix K with an every element $K_{nm} = k(x^{(n)}, x^{(m)})$.

Since $h(x^{(n)}) \cdot h(x^{(m)}) = k(x^{(n)}, x^{(m)})$, $HH^\top = K$

Let us define $\alpha = (K + \lambda I)^{-1} y$. Hence,

$$w^* = H^\top \alpha = h(x) \cdot \alpha$$

Naive prediction can be made as,

$$y^{\text{pred}} = w^* \cdot h(x^{\text{pred}})$$

Substituting value of w^* we get,

$$\begin{aligned} y^{\text{pred}} &= h(x) \cdot \alpha \cdot h(x^{\text{pred}}) \\ \implies y^{\text{pred}} &= \alpha \cdot h(x) \cdot h(x^{\text{pred}}) \\ \implies y^{\text{pred}} &= \alpha \cdot k(x, x^{\text{pred}}) \end{aligned}$$

Hence $y^{\text{pred}} = \alpha \cdot k(x, x^{\text{pred}})$, where α and the kernel function are both independent of w^* and the basis expansion function $h(\cdot)$.

Answer 4

Given that

$$h(x) = [c_0, c_1x, c_2x^2, \dots, c_Px^P], c_p = \sqrt{\binom{P}{p}}$$

Claim : $k(x, x') = (1 + x \cdot x')^P$

Proof :

$$\begin{aligned} h(x) \cdot h(x') &= [c_0, c_1x, c_2x^2, \dots, c_Px^P] \cdot [c_0, c_1x', c_2x'^2, \dots, c_Px'^P] \\ &= \sum_{p=0}^P (c_p^2 x^p x'^p) \\ &= \sum_{p=0}^P \binom{P}{p} (x \cdot x')^p \end{aligned} \tag{1}$$

Equation (1) can be simplified using binomial theorem to -

$$\begin{aligned} h(x) \cdot h(x') &= (1 + x \cdot x')^P \\ &= k(x, x') \end{aligned}$$

$$\boxed{k(x, x') = (1 + x \cdot x')^P}$$

Answer 5

```
1 def eval_basis_expanded_ridge(x, w, h):  
2     y = np.dot(w, h(x)[0])  
3     return y
```

Answer 6

```
1 def train_basis_expanded_ridge(X, Y, l, h):
2     H = h(X)
3
4     C = np.zeros((H[0].shape[0], H[0].shape[0]))
5     for item in H:
6         C = np.add(C, np.matmul(np.array([item]).T, np.array([item])))
7
8     I = np.identity(len(C))
9
10    w = np.linalg.solve(np.add(C, l * I),
11                        np.matmul(np.transpose(H), Y))
12
13    return w
```

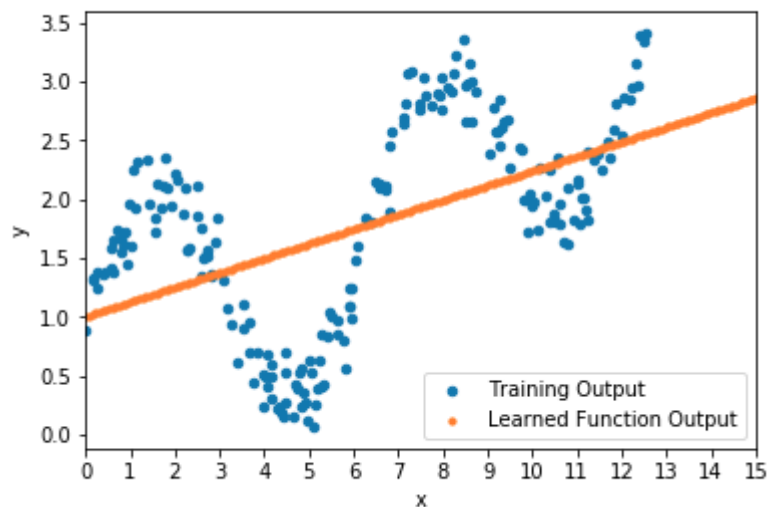

Answer 7

```
1 def basis_expanded_ridge_reg():
2     for P in [1, 2, 3, 5, 10]:
3         # Define polynomial basis expansion function
4         h = get_poly_expansion(P)
5
6         # Train basis expanded ridge regression function
7         W = train_basis_expanded_ridge(X_trn, Y_trn, 0.1, h)
8
9         # Print optimal weights and generate plot
10        print(f"P: {P}\nW: {W}\n")
11        Y = []
12        x_new = np.linspace(0,15,200)
13        H = h(x_new)
14        for x in x_new:
15            Y.append(eval_basis_expanded_ridge(x,W,h))
16        Y = np.array(Y)
17        plt.figure(1, figsize=(6, 4))
18        plt.scatter(X_trn, Y_trn, s=20, label="Training Output")
19        plt.scatter(x_new, Y, s=10, label="Learned Function Output")
20        plt.xlim([0, 15])
21        plt.xlabel("x")
22        plt.ylabel("y")
23        plt.xticks(np.arange(0, 16, 1))
24        plt.legend()
25        plt.show()
```

Below are the values of W seen for each value of P and corresponding plot of the final learned function as a function of x superimposed on the training data-

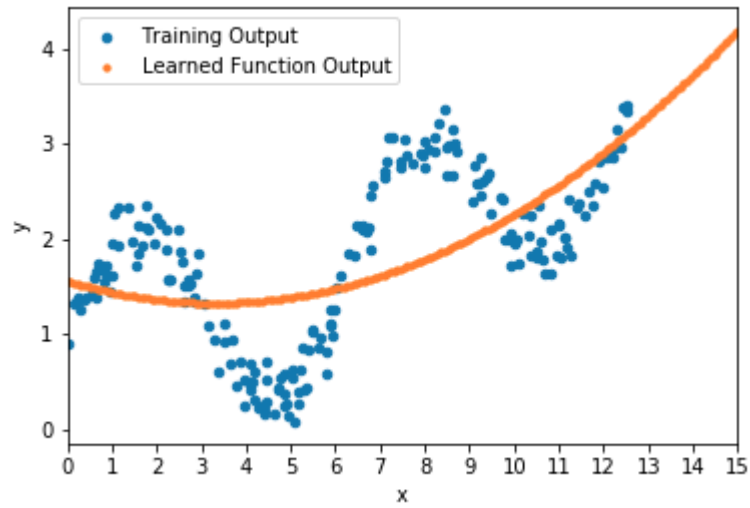
1. $P = 1$

$W = [1.00565302, 0.12351259]$



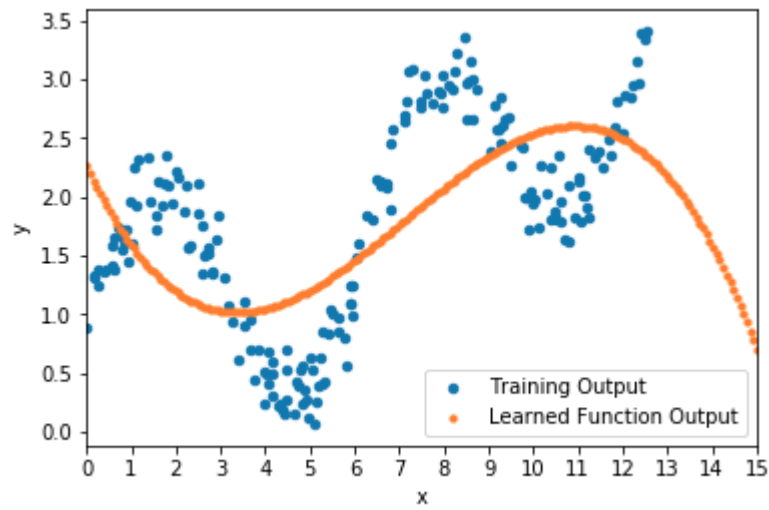
2. $P = 2$

$$W = [1.55636445, -0.09905134, 0.02105954]$$



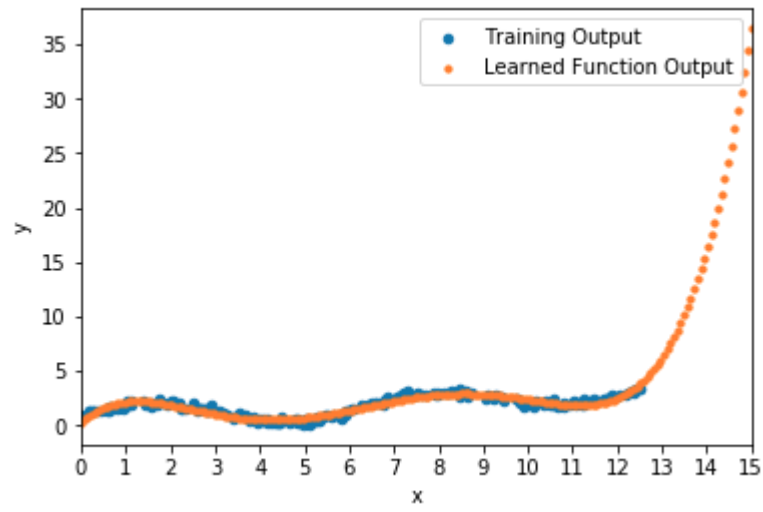
3. $P = 3$

$$W = [2.2585207, -0.4731082, 0.09166343, -0.0074039]$$



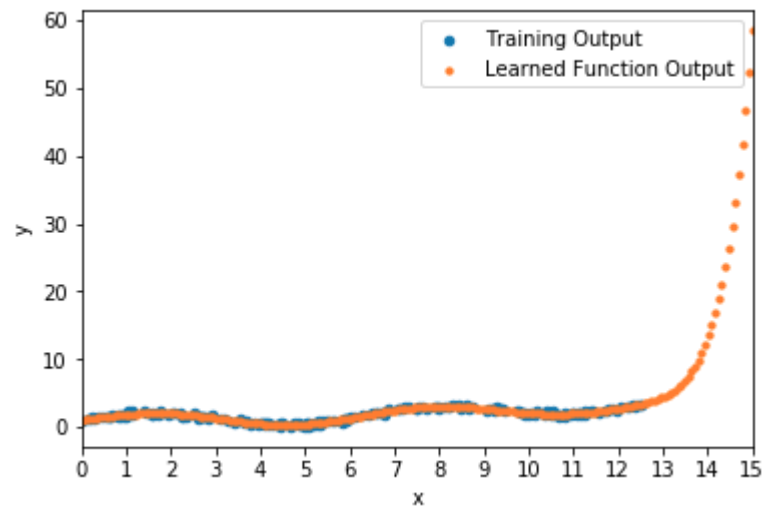
4. $P = 5$

$$W = [2.32031732e-01, 1.70733532e+00, -7.50673264e-01, 1.66018570e-01, -2.11820885e-02, 1.49990203e-03]$$



5. $P = 10$

$W = [1.00199719e+00, 3.46174032e-01, -6.79997048e-02, 4.04746152e-02, -2.41324409e-02, 7.30654951e-03, -1.30797499e-03, 1.49023358e-04, -1.05582815e-05, 3.71804211e-07, 2.73112041e-09]$



Answer 8

```
1 def get_poly_kernel(P):  
2     def k(x, xp):  
3         kernel_value = (1+np.inner(x, xp))**P  
4         return kernel_value  
5     return k
```

Answer 9

```
1 x = 0.5
2 xp = 0.7
3 k = get_poly_kernel(5)
4 h = get_poly_expansion(5)
5 out1 = k(x,xp)
6 out2 = np.inner(h(x),h(xp))
7 print("output 1", out1)
8 print("output 2", out2)
```

The outputs observed are -

- output 1 **4.484033437500002**
- output 2 **[[4.48403344]]**

Answer 10

```
1 from preface import *
2
3 def train_kernel_ridge(X,Y,l,k) :
4     N = X.shape[0]
5
6     # Compute K
7     K = np.empty((N,N))
8     for n in range(N) :
9         for m in range(N) :
10             K[n][m] = k(X[n], X[m])
11
12     a = np.linalg.solve(K + l*np.identity(N), Y)
13     return a
```

Answer 11

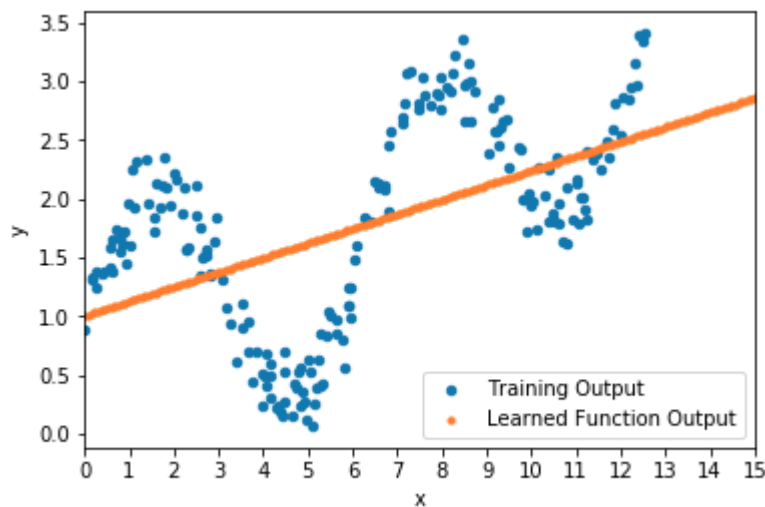
```
1 def eval_kernel_ridge(x_trn, x, a, k):  
2     y = np.dot(a, np.array(list(map(lambda x1: k(x1, x), x_trn))))  
3     return y
```

Answer 12

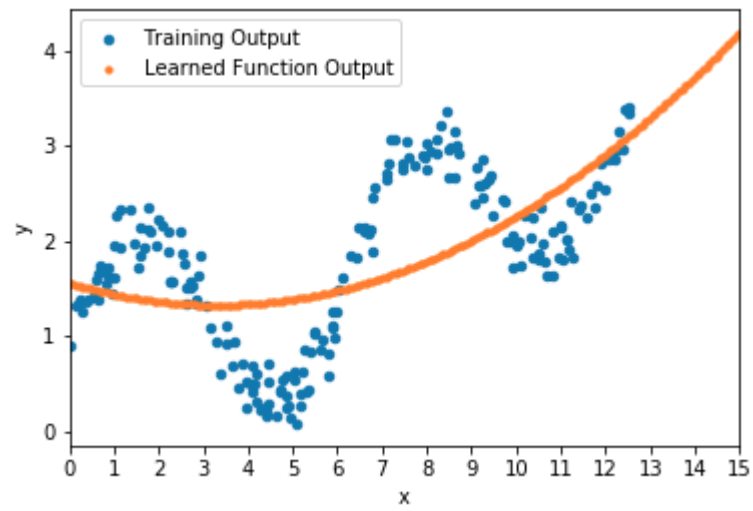
```
1 def kernel_ridge_reg():
2     for P in [1, 2, 3, 5, 10]:
3         # Define polynomial kernel function
4         k = get_poly_kernel(P)
5
6         # Train Kernel Ridge Regression Model
7         a = train_kernel_ridge(X_trn, Y_trn, 0.1, k)
8         print(f"P: {P}")
9         Y = []
10        X_new = np.linspace(0,15,200)
11
12        # Evaluate for new data
13        for X in X_new:
14            Y.append(eval_kernel_ridge(X_trn, X, a, k))
15        Y = np.array(Y)
16
17        # Plot results
18        plt.figure(1, figsize=(6, 4))
19        plt.scatter(X_trn, Y_trn, s=20, label="Training Output")
20        plt.scatter(X_new, Y, s=10, label="Learned Function Output")
21        plt.xlim([0, 15])
22        plt.xlabel("x")
23        plt.ylabel("y")
24        plt.xticks(np.arange(0, 16, 1))
25        plt.legend()
26        plt.show()
```

Plot of the final learned function as a function of x superimposed on the training data for each value of P is given below -

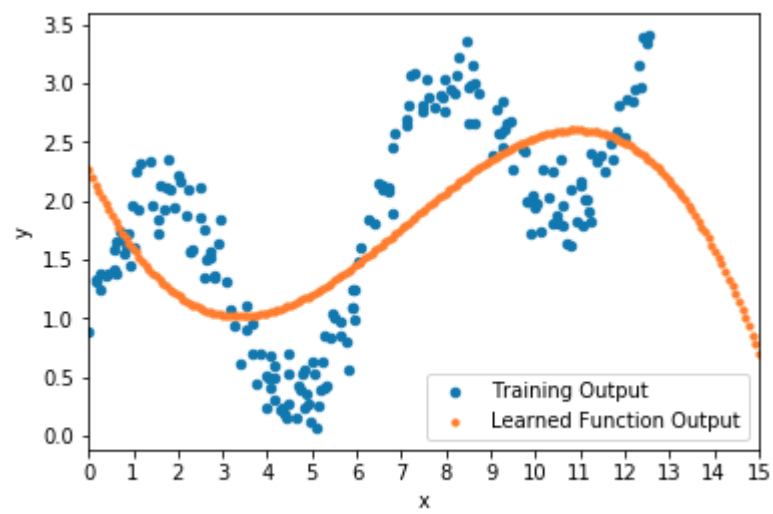
1. $P = 1$



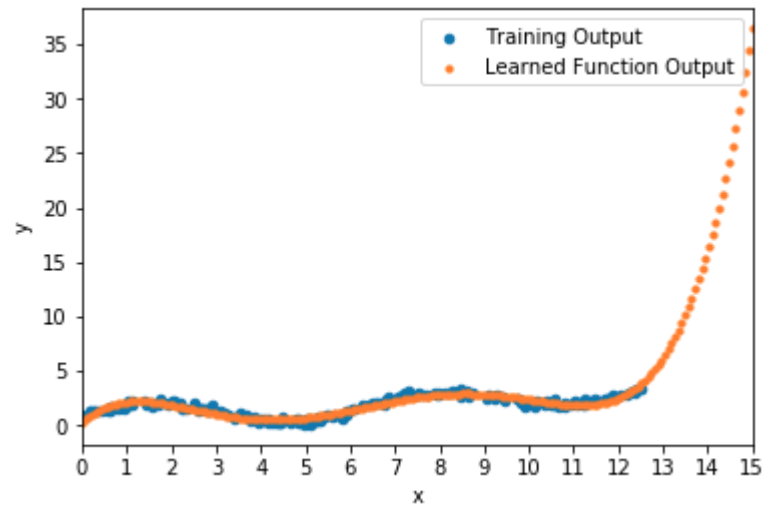
2. $P = 2$



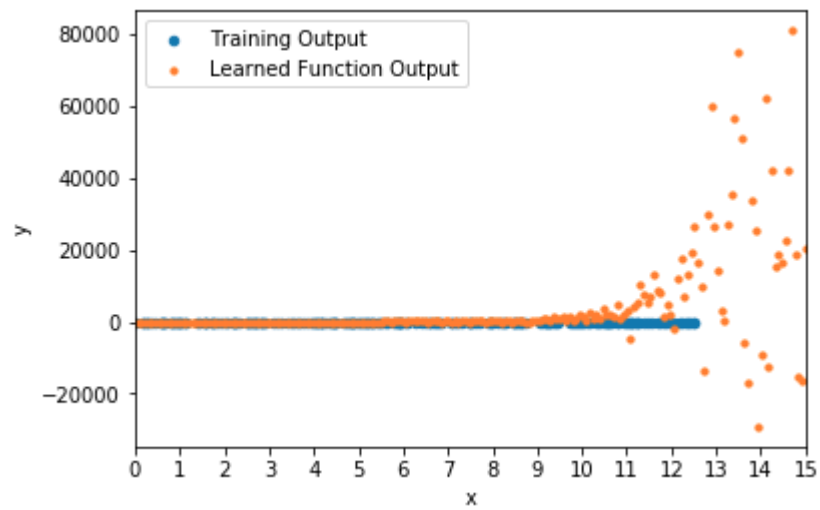
3. $P = 3$



4. $P = 5$



5. $P = 10$



Answer 13

From the previous results, we can see clearly from the plots that the results using kernel ridge regression are similar to results obtained using basis-expanded ridge regression for all P except for the case when $P = 10$.

If we consider the values for α and K generated for each P , we can check whether K is positive definite with respect to α . Following code evaluates $\alpha^\top K \alpha$ for each kernel

-

```
1 def is_kernel_valid(X_trn, Y_trn):
2     print("P\tPositive Definite Term")
3     for P in [1, 2, 3, 5, 10]:
4         k = get_poly_kernel(P)
5         a = train_kernel_ridge(X_trn, Y_trn, 0.1, k)
6         N = X_trn.shape[0]
7         K = np.empty((N,N))
8         for n in range(N) :
9             for m in range(N) :
10                 K[n][m] = k(X_trn[n], X_trn[m])
11         mer = round(np.matmul(np.matmul(a.transpose(),K),a),10)
12         print(f"{P}\t{mer}\n")
```

Results observed from this code are -

P	$\alpha^\top K \alpha$
1	1.0265933456
2	2.4325249636
3	5.3332049447
5	3.5564436669
10	-28.3624896136

According to Mercer's Theorem, a kernel function is valid if and only if K is positive definite. From above table, we can see that for $P = 10$, the value $\alpha^\top K \alpha < 0$. Hence the kernel is not positive definite for $P = 10$, and thus is not safe to use for prediction.

Answer 14

```
1 from sklearn import svm
2 from sklearn.metrics import hinge_loss
3 from sklearn.model_selection import KFold
4
5 def train_svm(X_trn, y_trn):
6     l_vals = [2, 20, 200]
7     splits = 5
8     kf = KFold(n_splits=splits, shuffle=True)
9     print("Lambda\t\tHinge Loss\n")
10    for l in l_vals:
11        clf = svm.SVC(kernel='linear', C=1 / (2 * l))
12        sum_hinge_loss = 0
13        for train_index, test_index in kf.split(X_trn):
14            # Split train-test
15            X_train, X_test = X_trn[train_index], X_trn[test_index]
16            y_train, y_test = y_trn[train_index], y_trn[test_index]
17
18            # Train the model
19            clf.fit(X_train, y_train)
20            predictions = clf.decision_function(X_test)
21
22            sum_hinge_loss += hinge_loss(y_test, predictions)
23
24        print(f"{l}\t\t{sum_hinge_loss / splits}")
25 train_svm(X_trn_real, Y_test_real)
```

Mean validation-set hinge loss using linear kernel :

λ	Hinge loss
2	0.0413724451
20	0.0539554495
200	0.1069876216

Answer 15

Validation-set hinge loss for polynomial basis function of degree 3 :

$\gamma \backslash \lambda$	2	20	200
0.001	0.0558986632	0.0464007172	0.0728558537
0.01	0.0416938698	0.05478073	0.0551566034
1	0.0443868818	0.0255245765	0.0603758981

Answer 16

Validation-set hinge loss for polynomial basis function of degree 5 :

$\gamma \backslash \lambda$	2	20	200
0.001	0.2375459622	0.5529637811	0.2828776961
0.01	0.3191144169	0.744456554	0.2543535828
1	0.0571467599	0.1683186749	0.1648740333

Answer 17

Validation-set hinge loss for radial basis function :

$\gamma \backslash \lambda$	2	20	200
1	0.3843586144	0.8345589747	0.8839056438
0.01	0.0776329575	0.3602880686	0.8367416977
0.001	0.2512348443	0.772343879	0.8779857099

Answer 18

- We chose the polynomial kernel with degree 3, $\gamma = 0.01$ and $\lambda = 2$.
- The estimated 0-1 error is: 0.0043796
- The observed generalization error on the leaderboard is: 0.00584