

$$\begin{aligned} a) \hat{w}_{krr} &= \underset{w}{\operatorname{argmin}} \| \Phi w - y \|_2^2 + \lambda \| w \|_2^2 \\ &= \underset{w}{\operatorname{argmin}} (\Phi w - y)^T (\Phi w - y) + \lambda w^T w \\ &= \underset{w}{\operatorname{argmin}} w^T \Phi^T \Phi w - 2w^T \Phi^T y + y^T y + 2\lambda w \end{aligned}$$

$$\text{let } f(w) = w^T \Phi^T \Phi w - 2w^T \Phi^T y + y^T y + 2\lambda w$$

$$\nabla w f = 2\Phi^T \Phi w - 2\Phi^T y + 2\lambda w = 0$$

$$\Phi^T \Phi w + \lambda w = \Phi^T y$$

$$(\Phi^T \Phi + \lambda I) w = \Phi^T y$$

$$\hat{w}_{krr} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

$$\text{let } \Phi = U \Sigma V^T$$

$$\hat{w}_{krr}^T = (V \Sigma^T U^T U \Sigma V^T + \lambda I)^{-1} V \Sigma^T U^T y$$

$$= (V \Sigma^T \Sigma V^T + \lambda I)^{-1} V \Sigma^T U^T y$$

$$= (V \Sigma^T \Sigma V^T + \lambda V V^T)^{-1} V \Sigma^T U^T y$$

$$= (V (\Sigma^T \Sigma + \lambda I) V^T)^{-1} V \Sigma^T U^T y$$

$$= V (\Sigma^T \Sigma + \lambda I)^{-1} V^T V \Sigma^T U^T y$$

$$= V (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T U^T y$$

$$\text{From class, we know that } (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T = \Sigma^T (\Sigma \Sigma^T + \lambda I)^{-1}$$

$$\begin{aligned} \Rightarrow \hat{w}_{krr} &= V \Sigma^T (\Sigma \Sigma^T + \lambda I)^{-1} U^T y \\ &= V \Sigma^T U^T U (\Sigma \Sigma^T + \lambda I)^{-1} U^T y \\ &= \Phi^T U (\Sigma \Sigma^T + \lambda I)^{-1} U^T y \\ &= \Phi^T (U (\Sigma \Sigma^T + \lambda I) U^T)^{-1} y \\ &= \Phi^T (U \Sigma \Sigma^T U^T + \lambda U U^T)^{-1} y \\ &= \Phi^T (U \Sigma V^T V \Sigma^T U^T + \lambda I)^{-1} y \\ &= \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y \end{aligned}$$

$$\text{Thus, } \hat{w}_{krr} = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y$$

For this solution to be valid,  $(\Sigma \Sigma^T + \lambda I)^{-1}$  must exist. While  $\Sigma$  may have some singular values  $= 0$ , causing  $\Sigma \Sigma^T$  to have 0s on the diagonal, adding a  $\lambda > 0$  to all diagonal entries of  $\Sigma \Sigma^T$  means all values on diagonal of  $\Sigma \Sigma^T + \lambda I$  are  $> 0$ , since all  $\sigma_i \geq 0$  so  $\sigma_i + \lambda > 0$ , meaning  $(\Sigma \Sigma^T + \lambda I)^{-1}$  exists regardless of what the rank of  $\Phi$  is.

b)

$$\begin{aligned} i) \hat{w}_{krr} &= \Phi^T \alpha = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y \\ \Phi^T \alpha &= \Phi^T (K + \lambda I)^{-1} y \text{ since we define } K = \Phi \Phi^T \\ (\Phi^T)^+ \Phi^T \alpha &= (\Phi^T)^+ \Phi^T (K + \lambda I)^{-1} y \end{aligned}$$

we know  $\Phi \in \mathbb{R}^{d \times n}$  and  $(\Phi^T)^+ \in \mathbb{R}^{n \times d}$ , so  $(\Phi^T)^+ \Phi^T \in \mathbb{R}^{n \times n}$  since  $\Phi$  has  $n$  linearly independent rows,  $\Phi^T$  has  $n$  linearly independent columns, so  $(\Phi^T)^+ \Phi^T = I$ .

$$\text{Thus, } \alpha = (K + \lambda I)^{-1} y$$

$$ii) \hat{y} = \hat{w}_{krr}^T \Phi(z) = (\Phi^T \alpha)^T \Phi(z) = \alpha^T \Phi \Phi(z)$$

$$\text{we know } \Phi \Phi(z) = \begin{pmatrix} \Phi(x_1)^T \Phi(z) \\ \vdots \\ \Phi(x_n)^T \Phi(z) \end{pmatrix} = \begin{pmatrix} K(x_1, z) \\ \vdots \\ K(x_n, z) \end{pmatrix} = K(x, z)$$

$$\Rightarrow \hat{y} = \alpha^T \Phi \Phi(z) = \alpha^T K(x, z) = ((K + \lambda I)^{-1} y)^T K(x, z)$$

$$\Rightarrow \hat{y} = (K + \lambda I)^{-1} y)^T K(x, z)$$

c) We defined  $K = \Phi \Phi^T$  where  $K_{ij} = K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$ , or represents the similarity between the vectors. This is useful because we can find the weight vector  $\alpha$  using  $K$  which requires a simple evaluation between vector pairs using the kernel function. This allows to avoid computing and storing  $\Phi(x_i)$  which grows on the order of power such as requiring  $O(p^d)$  to store a  $d$ -degree polynomial from our  $p$ -dimensional input.

d) code attached

2)  
a)

i)  $g(x) = g_1(x) + g_2(x)$

$g_1(x)$  convex:  $g_1(w) \geq g_1(v) + \nabla g_1(v)^T (w-v) \quad \forall w, v \in \mathbb{R}^p$

$g_2(x)$  convex:  $g_2(w) \geq g_2(v) + \nabla g_2(v)^T (w-v) \quad \forall w, v \in \mathbb{R}^p$

$g_1(w) + g_2(w) \geq g_1(v) + g_2(v) + \nabla g_1(v)^T (w-v) + \nabla g_2(v)^T (w-v) \quad \forall w, v \in \mathbb{R}^p$

$g_1(w) + g_2(w) \geq g_1(v) + g_2(v) + \nabla (g_1(v) + g_2(v))^T (w-v) \quad \forall w, v \in \mathbb{R}^p$

$\Rightarrow g(w) \geq g(v) + \nabla g(v)^T (w-v) \quad \forall w, v \in \mathbb{R}^p$

Thus,  $g(x)$  is a convex function

ii)  $g(w) = w^T K w$

$g(w) = g(v) + \nabla g^T (w-v) + \frac{1}{2} (w-v)^T \nabla^2 g (w-v)$

$= v^T K v + 2(Kv)^T (w-v) + (w-v)^T K (w-v)$

1st-order Taylor expansion includes the 1st 2 terms:

$g(v) + \nabla g^T (w-v) = v^T K v + 2(Kv)^T (w-v)$

Need to show  $\forall w, v \in \mathbb{R}^p \quad g(w) \geq g(v) + \nabla g^T (w-v)$

Since  $K$  is positive semi-definite,  $(w-v)^T K (w-v) \geq 0$

$\Rightarrow g(w) = v^T K v + 2(Kv)^T (w-v) + (w-v)^T K (w-v)$

$\geq v^T K v + 2(Kv)^T (w-v)$

$= g(v) + \nabla g^T (w-v)$

$\Rightarrow g(w) \geq g(v) + \nabla g^T (w-v) \quad \forall w, v \in \mathbb{R}^p$  so  $g(w)$  is a convex function

b) At  $t$  steps,  $f(\alpha^{(t)}) = \sum_{i=1}^n (1 - y_i k_i^T \alpha^{(t)})_+ + \lambda \alpha^{(t)T} K \alpha^{(t)}$

Since  $y_i k_i^T \alpha^{(t)} > 1, 1 - y_i k_i^T \alpha^{(t)} < 0 \Rightarrow (1 - y_i k_i^T \alpha^{(t)}) = 0$

Thus,  $f(\alpha^{(t)}) = \sum_{i=1}^n (1 - y_i k_i^T \alpha^{(t)})_+ + \lambda \alpha^{(t)T} K \alpha^{(t)} \Rightarrow f(\alpha^{(t)}) = \lambda \alpha^{(t)T} K \alpha^{(t)}$

$\alpha^{(t+1)} = \alpha^{(t)} - \nabla f|_{\alpha^{(t)}}$

$= \alpha^{(t)} - \nabla (2\lambda K \alpha^{(t)})$

$\Rightarrow \alpha^{(t+1)} = \alpha^{(t)} - 2\lambda K \alpha^{(t)}$

c)  $\nabla_{\alpha} f = \sum_{i=1}^n \mathbb{I} \{1 - y_i k_i^T \alpha > 0\} (-y_i k_i) + 2\lambda K \alpha$

At optimal  $\nabla f|_{\alpha^*} = 0 \Rightarrow \sum_{i=1}^n \mathbb{I} \{1 - y_i k_i^T \alpha^* > 0\} (-y_i k_i) + 2\lambda K \alpha^* = 0$

Suppose  $\lambda K \alpha^* = 0$ . Then either  $\lambda = 0$  which is not the case, or  $K \alpha^* = 0$ , meaning  $k_i^T \alpha^* = 0 \quad \forall i$ , so all  $n$  samples are classified the same way, also unlikely to be the case. Therefore,  $2\lambda K \alpha^* \neq 0$ , so  $\sum_{i=1}^n \mathbb{I} \{1 - y_i k_i^T \alpha^* > 0\} (-y_i k_i) \neq 0$

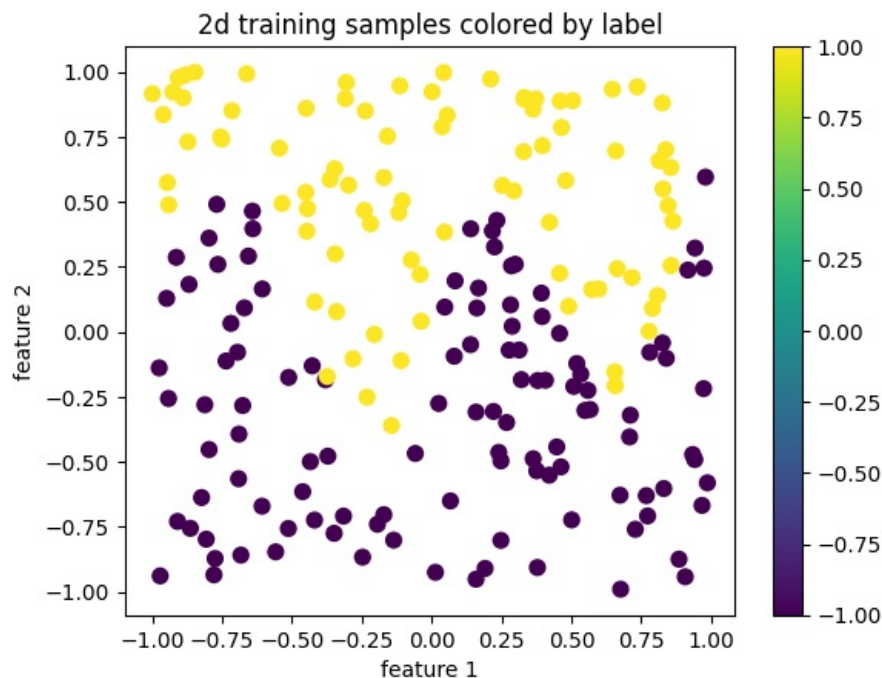
so some samples have  $1 - y_i k_i^T \alpha^* > 0 \Rightarrow y_i k_i^T \alpha^* < 1$ , so they are classified correctly (as  $\alpha^{(t)}$  classified all correctly) but are close to the decision boundary, so some loss is added, so not all terms = 0 in loss as we can obtain even perfect classification with loss  $> 0$ .

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import scipy.io as sio
import sys
import numpy.linalg as la
from tabulate import tabulate
import random
from PIL import Image
import math
```

1d

```
In [2]: n = 200
p = 2
X = 2 * (np.random.rand(n, p) - 0.5)
y = np.sign(X[:, 1] - (X[:, 0]**2 / 2 + np.sin(X[:, 0] * 7) / 2))

plt.figure(1)
plt.scatter(X[:, 0], X[:, 1], 50, c=y)
plt.colorbar()
plt.xlabel("feature 1")
plt.ylabel("feature 2")
plt.title("2d training samples colored by label")
plt.show()
```



```
In [3]: def krr(sigma, lam):
    """ YOUR CODE STARTS HERE """
    K = np.zeros((n, n))
    for i in range(n):
        for j in range(n):
            # note: ||x_i - x_j||^2 = (x_i - x_j)^T (x_i - x_j) = (x_i - x_j) dot (x_i - x_j)
            K[i, j] = math.e ** (-(np.dot(X[i] - X[j], X[i] - X[j])) / (2 * sigma**2))
    alpha = la.inv(K + lam * np.eye(n)) @ y
    yhat = K @ alpha
    """ YOUR CODE ENDS HERE """

    y2 = np.array(np.sign(yhat))
    plt.figure(2)
    plt.scatter(X[:, 0], X[:, 1], 50, c=y2)
    plt.colorbar()
    plt.xlabel("feature 1")
    plt.ylabel("feature 2")
    plt.title("2d training samples colored by PREDICTED label")
    plt.show()

    ntest = 2000
    Xtest = 2 * (np.random.rand(ntest, p) - 0.5)

    """ YOUR CODE STARTS HERE """
    Ktest = np.zeros((ntest, n))
    for i in range(ntest):
        for j in range(n):
            # note: ||x_i - x_j||^2 = (x_i - x_j)^T (x_i - x_j) = (x_i - x_j) dot (x_i - x_j)
```

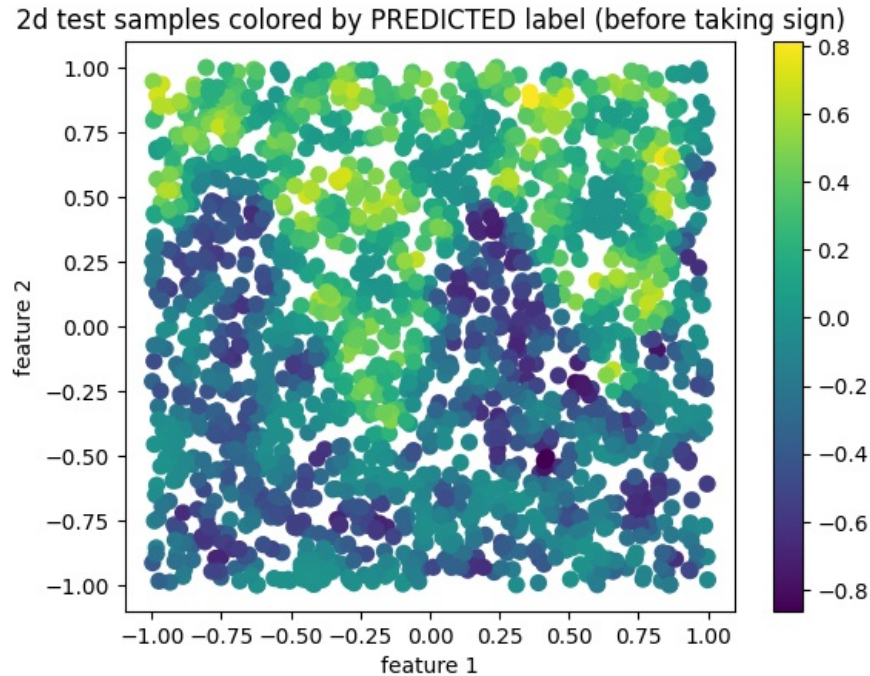
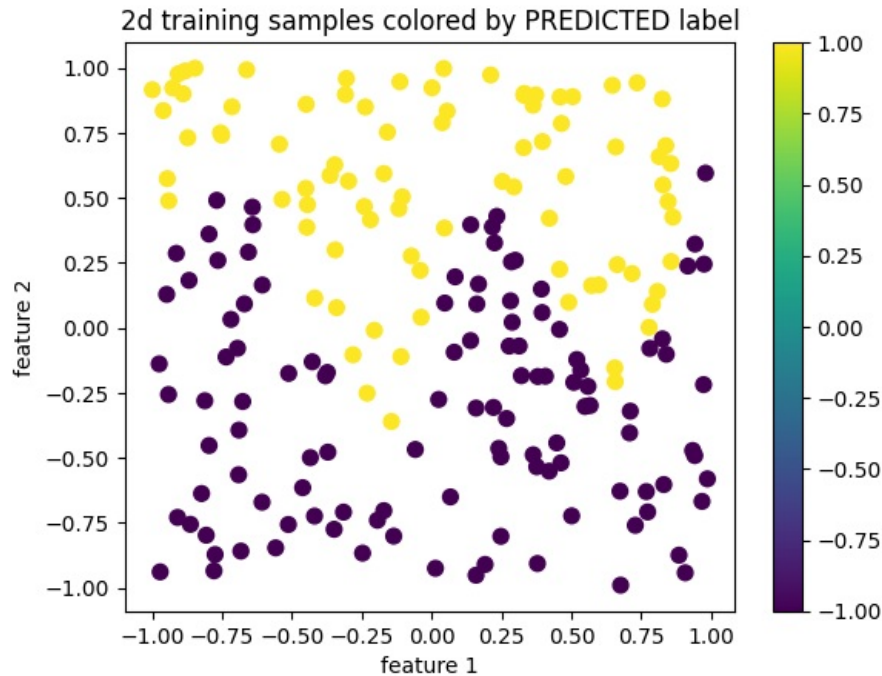
```

Ktest[i, j] = math.e ** (-(np.dot(Xtest[i] - X[j], Xtest[i] - X[j])) / (2 * sigma**2))
ytest = Ktest @ alpha
### YOUR CODE ENDS HERE ###

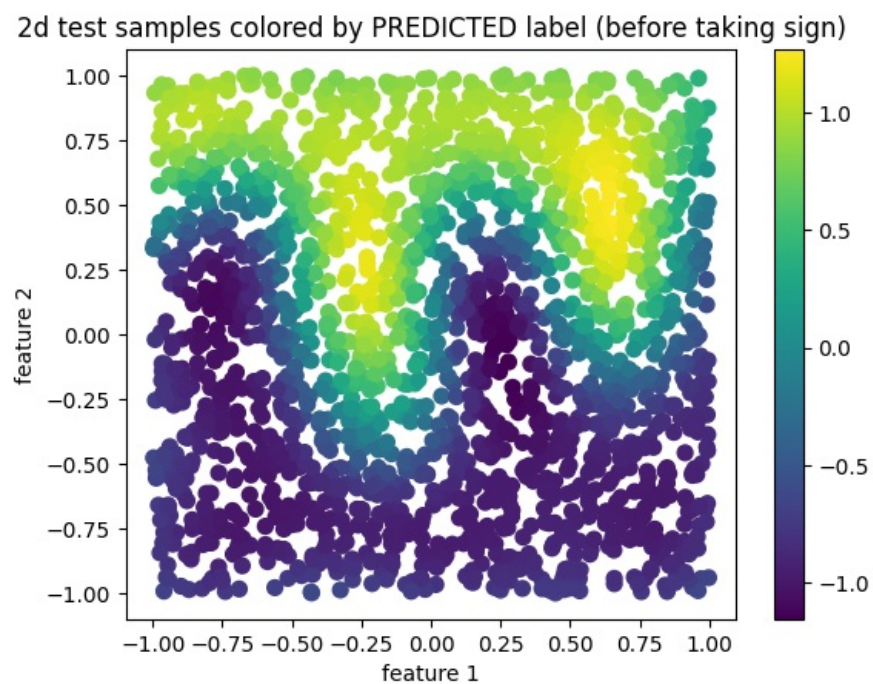
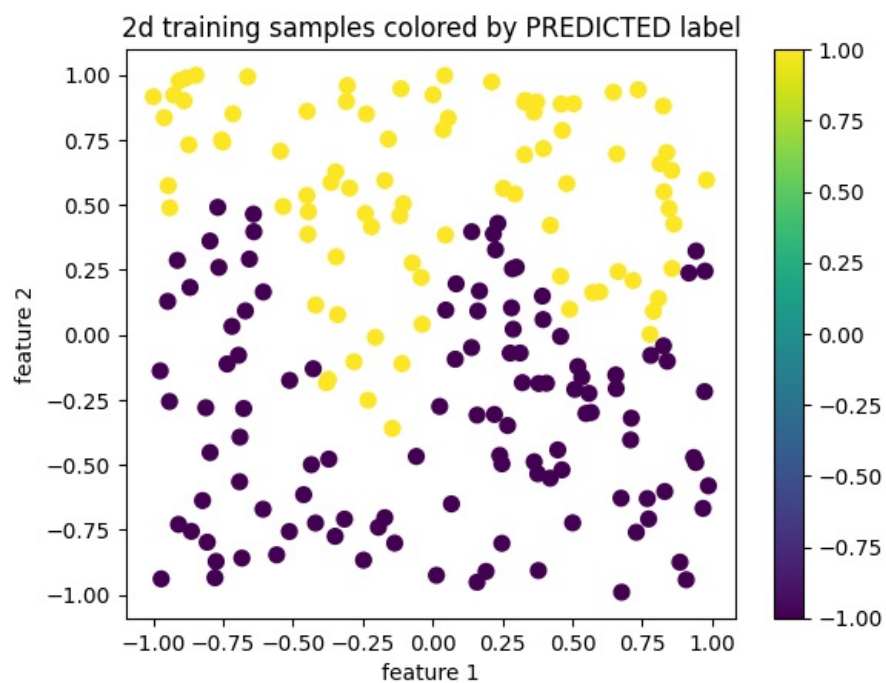
plt.figure(3)
plt.scatter(Xtest[:, 0], Xtest[:, 1], 50, c=np.array(ytest))
plt.colorbar()
plt.xlabel("feature 1")
plt.ylabel("feature 2")
plt.title("2d test samples colored by PREDICTED label (before taking sign)")
plt.show()

```

In [4]: `krr(sigma=0.05, lam=1)`

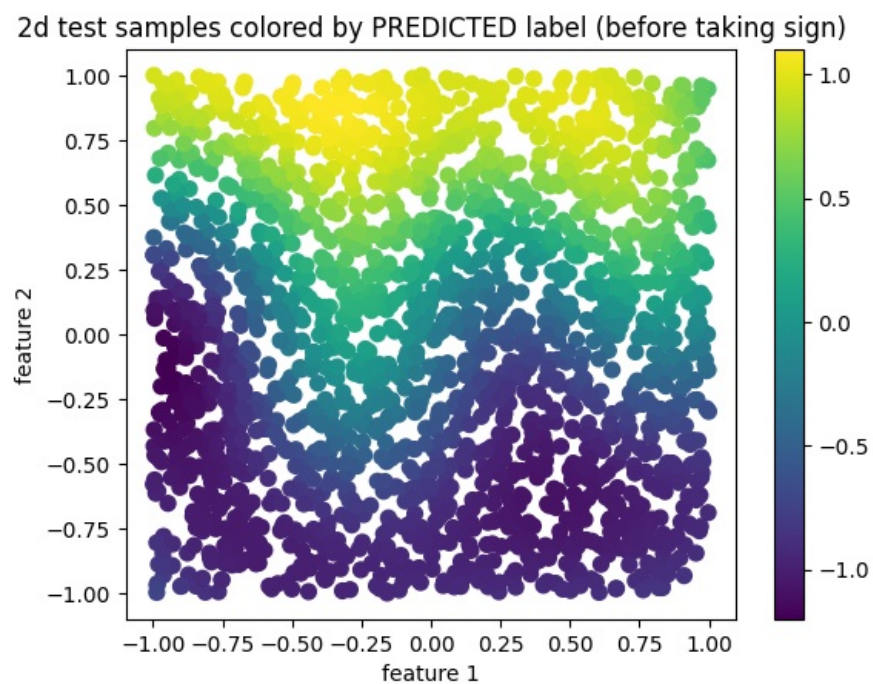
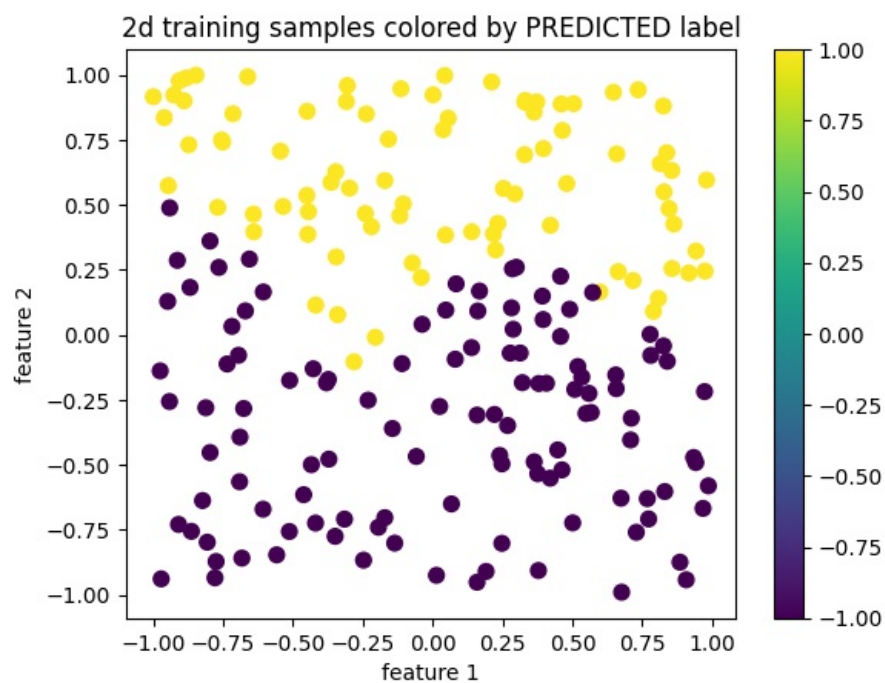


In [5]: `krr(sigma=0.25, lam=1)`

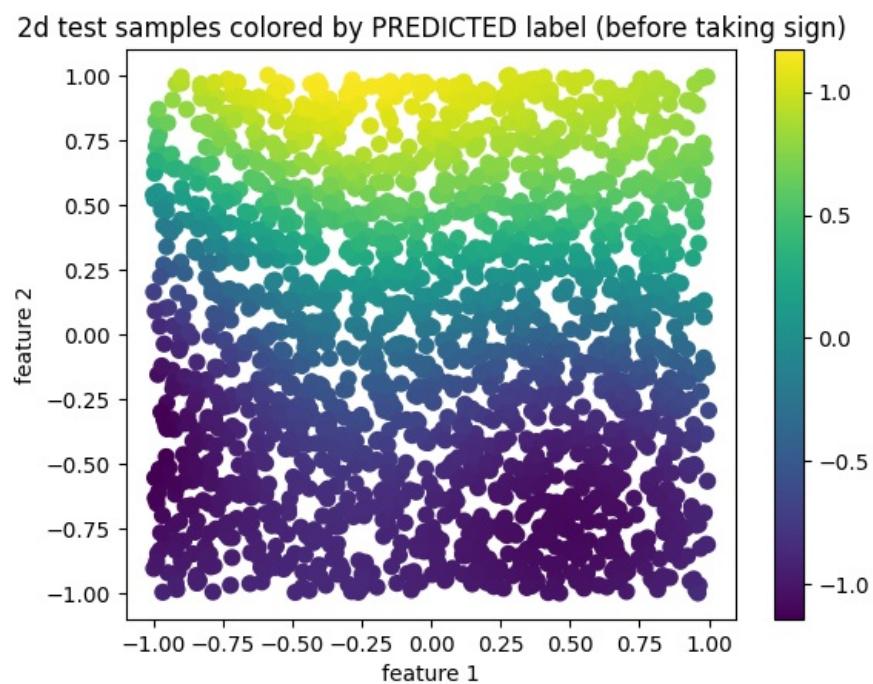
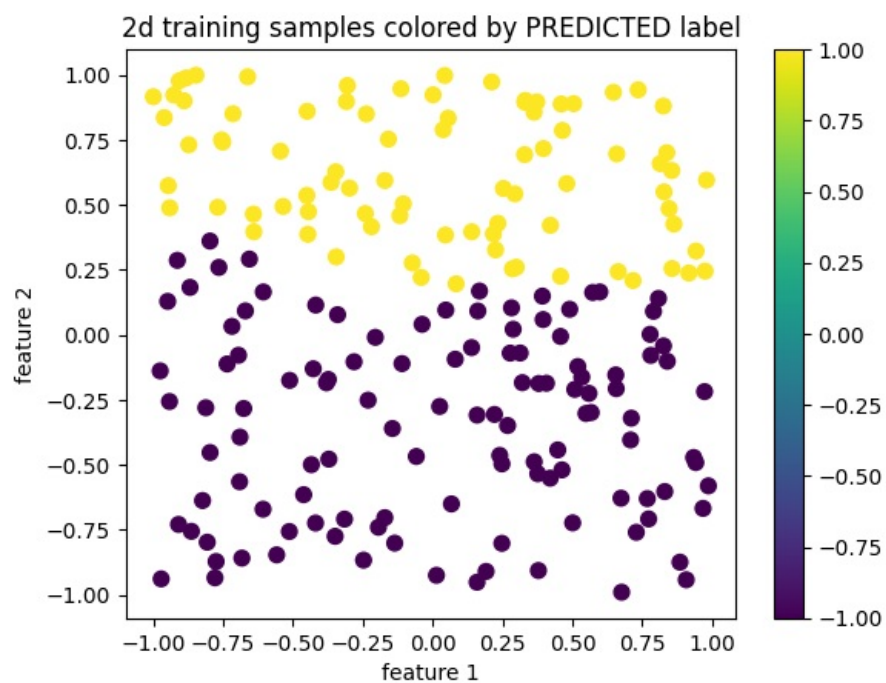


```
In [6]: krr(sigma=0.5, lam=1)
```

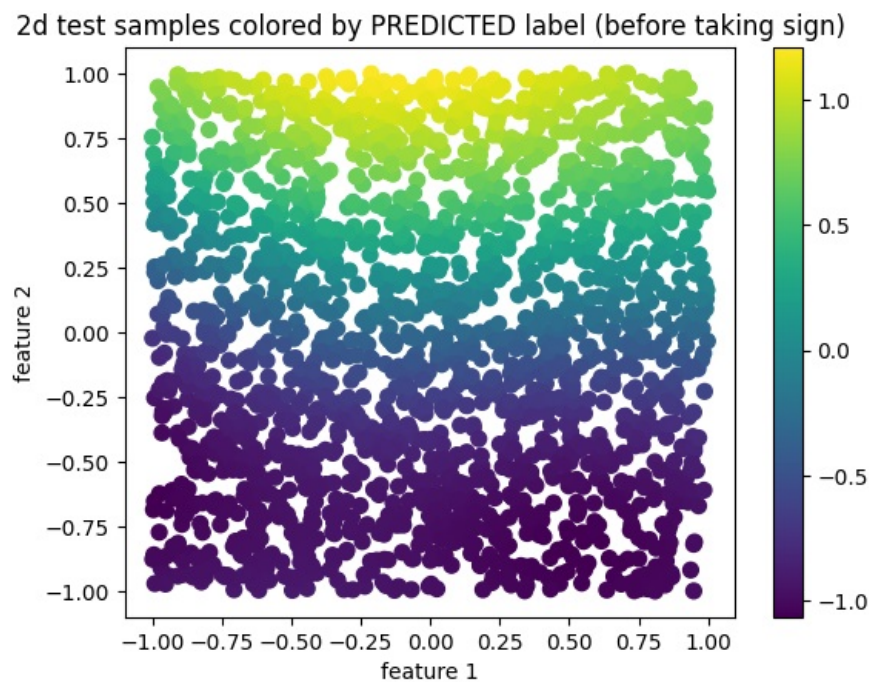
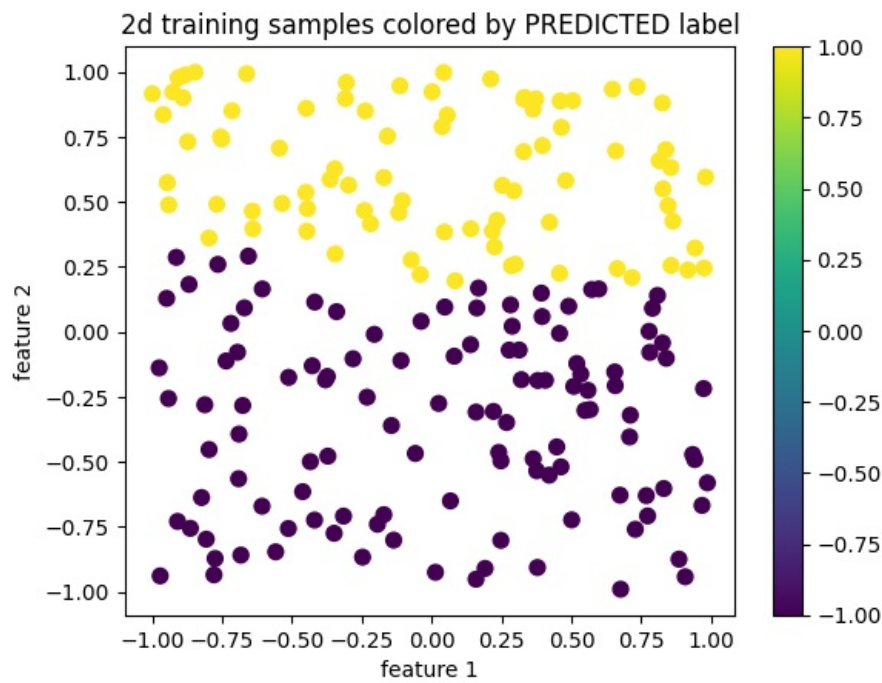




```
In [7]: krr(sigma=0.75, lam=1)
```



```
In [8]: krr(sigma=1, lam=1)
```



We see that the model becomes more linear as  $\sigma$  increases. When  $\sigma = 0.05$ , we see that the decision boundary in the test samples is not very distinct. On the other hand, when  $\sigma = 0.8$  or  $\sigma = 1$ , for example, the decision boundary is significantly more linear compared to  $\sigma = 0.25$ , having a highly nonlinear decision boundary.

If we think of  $K(x_i, x_j)$  as providing a similarity score between  $x_i$  and  $x_j$ , we see that as  $\|x_i - x_j\|$  decreases,  $K(x_i, x_j)$  increases, which makes sense because if  $x_i$  and  $x_j$  are closer together, they should have a higher similarity score. Given this, we see that as  $\sigma$  increases,  $K(x_i, x_j)$  increases, meaning a higher  $\sigma$  treats the same samples as more similar to each other. This means the decision boundary can incorporate information of both  $x_i$  and  $x_j$  without becoming too nonlinear since the  $K(x_i, x_j)$  treats them as similar already. Therefore, we see that with higher  $\sigma$ , the decision boundary is more linear.