

COMP0083: Kernel assignment: advanced topics in machine learning

January 24, 2025

Contents

Question 1	2
1-1	2
1-2	2
Question 2	3
2-1	3
2-2	8

Question 1

Question 1.1

From the given figure, the dataset consists of a red class enclosed by a ring of blue points. A simple feature space that allows error-free linear classification for this dataset is defined as:

$$\phi_1(x_1, x_2) = x_1^2 + x_2^2. \quad (1)$$

In this feature space, one can separate two classes by choosing a threshold value r_0 such that

- For $\phi_1(x_1, x_2) < r_0$, classify as red,
- For $\phi_1(x_1, x_2) \geq r_0$, classify as blue.

Question 1.2

We can apply eigendecomposition of inner product matrix K . Let us assume $\lambda_1, \dots, \lambda_m$ are eigenvalues of K and v_1, \dots, v_m are corresponding eigenvectors. Then, we can decompose K as

$$K = Q\Lambda Q^{-1} = Q\Lambda Q^T, \quad (2)$$

where Λ is a diagonal matrix of eigenvalues λ_i and Q is a matrix of eigenvectors $[v_1, \dots, v_m]$. As K is positive semidefinite, we can split the original diagonal matrix apart and get

$$K = Q\Lambda^{\frac{1}{2}} \left(\Lambda^{\frac{1}{2}}\right)^T Q^T = \Phi\Phi^T, \quad (3)$$

where we have $\Phi = Q\Lambda^{\frac{1}{2}}$.

In the question, the matrix K can be express as inner product of feature space $\phi(x_i)$. Thus, we can deduce

$$\begin{aligned} K_{ij} &= \langle \phi(x_i), \phi(x_j) \rangle_{\mathcal{H}} = (\Phi\Phi^T)_{ij}, \\ \Rightarrow \phi(x_i) &= \Phi_i, \\ \Rightarrow \phi(x_i) &= [\sqrt{\lambda_1}v_{i1}, \dots, \sqrt{\lambda_m}v_{im}]. \end{aligned} \quad (4)$$

where we find the feature space representation in terms of eigenvalue and eigenvector of matrix K .

Question 2

Question 2.1 Incomplete Cholesky for efficient COCO

Let us calculate the exact cost of COCO first. The solution of COCO is given as

$$\begin{bmatrix} 0 & \frac{1}{n}\tilde{K}\tilde{L} \\ \frac{1}{n}\tilde{L}\tilde{K} & 0 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \gamma \begin{bmatrix} \tilde{K} & 0 \\ 0 & \tilde{L} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}. \quad (5)$$

Thus, the computation cost is from solving the eigenvalue problem in above solution. It involves matrix multiplications. To compute general $n \times n$ matrix multiplications, we have

```

1 input A and B, both n by n matrices
2 initialize C to be an n by n matrix of all zeros
3 for i from 1 to n:
4     for j from 1 to n:
5         for k from 1 to n:
6             C[i][j] = C[i][j] + A[i][k]*B[k][j]
7 output C (as A*B)

```

This needs n^3 multiplications and $n^3 - n^2$ additions. So the overall cost is in the order of $O(n^3)$.

To approximate our solution via incomplete Cholesky, we can have two different methods. Since both methods work well, we are going to introduce both.

Method 1: QR decomposition on matrix

In this method, we directly go to the final step of the problem. To solve Equation 5, we first simplify the equation as

$$Av = \gamma Bv, \quad (6)$$

where

$$A = \begin{bmatrix} 0 & \frac{1}{n}\tilde{K}\tilde{L} \\ \frac{1}{n}\tilde{L}\tilde{K} & 0 \end{bmatrix}, \quad B = \begin{bmatrix} \tilde{K} & 0 \\ 0 & \tilde{L} \end{bmatrix}. \quad (7)$$

Now, we apply QR decomposition on matrix B

$$B = D^\top C^\top C D = D^\top D. \quad (8)$$

Substitute this back to our problem, we can further derive

$$\begin{aligned} Av &= \gamma D^\top D v, \\ \Rightarrow AD^{-1}V &= \gamma D^\top V, \\ \Rightarrow (D^\top)^{-1} AD^{-1}V &= \gamma V. \end{aligned} \quad (9)$$

Here, we rescale the vector from v to $V = Dv$, transforming the original problem into an eigenvalue problem of lower dimensionality. The computational cost can be analysed as follows. Applying the Incomplete Cholesky decomposition to the matrix B yields a $m \times n$ matrix D , where m is smaller than n . Subsequently, the eigenvalue problem is solved for a reduced $m \times m$ matrix. This process involves two primary computational costs: (1) the Incomplete Cholesky decomposition, which has a complexity of $O(nm^2)$, and (2) solving the eigenvalue problem, with a cost of $O(m^3)$. Together, these steps significantly reduce the computational burden compared to solving the original problem directly.

Method 2: Feature map projection

Instead of decompose matrix B, we apply Incomplete Cholesky on both K and L .

$$\begin{aligned} K &= X^\top X = (QR)^\top QR = R^\top R, \\ L &= Y^\top Y = (PS)^\top PS = S^\top S, \end{aligned} \quad (10)$$

where we let $X = QR$ and $Y = PS$. As we make projections for both X and Y , we can modify the original problem into a new basis. In that case, we will have a new equation for this problem. First of all, we define $\theta = RH\alpha$ and $\eta = SH\beta$. These changes will largely simplify our equations. Then, we can start to derive our eigenvalue equations. We write Equation 5 as

$$\begin{aligned} -\frac{1}{n}\tilde{K}\tilde{L}\beta + \lambda\tilde{K}\alpha &= 0, \\ -\frac{1}{n}\tilde{L}\tilde{K}\alpha + \gamma\tilde{L}\beta &= 0. \end{aligned} \quad (11)$$

which is easier for our substitutions. We have the following derivation.

$$\begin{aligned} -\frac{1}{n} (HR^\top RH) (HS^\top SH) \beta + \lambda (HR^\top RH) \alpha &= 0, \\ -\frac{1}{n} (HS^\top SH) (HR^\top RH) \alpha + \gamma (HS^\top SH) \beta &= 0. \end{aligned} \quad (12)$$

Multiply the first equation by α^\top , and the second by β^\top . The first one is then written as

$$\begin{aligned} -\frac{1}{n} (\alpha^\top HR^\top) (RHS^\top) (SH\beta) + \lambda (\alpha^\top HR^\top) (RH\alpha) &= 0, \\ \Rightarrow -\frac{1}{n} \theta^\top RHS^\top \eta + \lambda \theta^\top \theta &= 0. \end{aligned} \quad (13)$$

Similarly, we get the second equation as

$$\begin{aligned} -\frac{1}{n} (\beta^\top HS^\top) (SHR^\top) (RH\alpha) + \gamma (\beta^\top HS^\top) (SH\beta) &= 0, \\ \Rightarrow -\frac{1}{n} \eta^\top SHR^\top \theta + \gamma \eta^\top \eta &= 0. \end{aligned} \quad (14)$$

The solution can be expressed as

$$\begin{bmatrix} 0 & \frac{1}{n} RHS^\top \\ \frac{1}{n} SHR^\top & 0 \end{bmatrix} \begin{bmatrix} \theta \\ \eta \end{bmatrix} = \gamma \begin{bmatrix} \theta \\ \eta \end{bmatrix}. \quad (15)$$

The witness functions f and g are

$$\begin{aligned} f(x) &= \langle f, X \rangle = f^\top X = \theta^\top R, \\ g &= \langle g, Y \rangle = g^\top Y = \eta^\top S. \end{aligned} \quad (16)$$

In this case, we reduce the dimension of the problem again. We get the same order of computational cost as method 1. In summary, we reduce the computational cost by using Incomplete Cholesky from $O(n^3)$ to $O(nm^2)$.

Implement the incomplete Cholesky-based COCO

Here is the code for Gaussian kernel.

```
1 def gaussian_kernel(x, sigma):
2     n = len(x)
3     K = np.zeros((n, n))
4     for i in range(n):
5         for j in range(n):
6             K[i, j] = np.exp(-((x[i] - x[j]) ** 2).sum() / sigma**2)
7     return K
```

Gaussian Kernel

We follow the code provided in extra readings.

```
1 def incomplete_cholesky(K, eta):
2     dim = K.shape[0]
3     R = np.zeros((dim, dim))
4     d = np.diag(K).copy()
5     I = []
6     j = 0
7     while True:
8         a = np.max(d)
9         if a <= eta:
10             break
11         pivot_index = np.argmax(d)
12         I.append(pivot_index)
13         nu = np.sqrt(a)
14         for i in range(dim):
15             R[j, i] = (K[I[j], i] - R[:, j, i] @ R[:, I[j]]) / nu
16             d[i] -= R[j, i] ** 2
17         j += 1
18     R = R[:, j, :]
19     return R
```

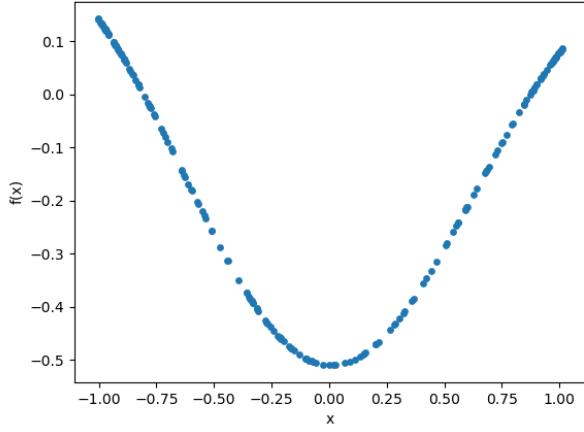
Incomplete cholesky

Here is code for first method. We apply incomplete cholesky on matrix B.

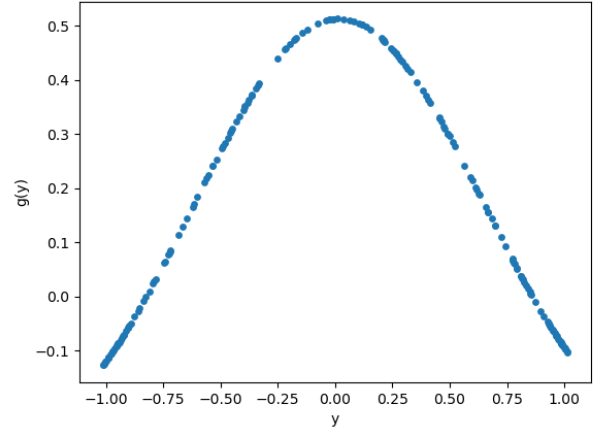
```
1 def coco_method1():
2     # generate the data, with N samples
3     N = 200
4     t = np.random.uniform(0, 2 * np.pi, N)
5     n1 = np.random.normal(0, 0.01, N)
6     n2 = np.random.normal(0, 0.01, N)
7     x = np.sin(t) + n1
8     y = np.cos(t) + n2
9
10    # define the parameters for gaussian kernel and incomplete cholesky
11    sigma = 1
12    eta = 1e-6
13
14    # define K and L
15    K = gaussian_kernel(x, sigma)
16    L = gaussian_kernel(y, sigma)
17    H = np.eye(N) - np.ones((N, N)) / N
18    K_tilde = H @ K @ H
19    L_tilde = H @ L @ H
20    # define zero blocks
21    z = np.zeros((N, N))
22    # define matrix for eigen-problem
23    # A v = gamma * B v
24    A = np.block([[z, K_tilde @ L_tilde / N],
25                  [L_tilde @ K_tilde / N, z]])
26    B = np.block([[K_tilde, z],
27                  [z, L_tilde]])
28
29    # incomplete Cholesky solution
30    # apply QR decomposition on B -> R^T R
31    R = incomplete_cholesky(B, eta)
32    # A v = gamma * R.T * R v, let u = Rv -> v = inv(R)u
33    # A inv(R)u = gamma * R.T * u
34    # inv(R.T) A inv(R) u = gamma u
35    C = np.linalg.pinv(R.T) @ A @ np.linalg.pinv(R)
36    eigvals, eigvecs = np.linalg.eig(C)
37    coco = np.abs(np.real(eigvals[0]))
38    eigvecs_v = np.linalg.pinv(R) @ eigvecs[:, 0]
39    alpha = eigvecs_v[:N]
40    beta = eigvecs_v[N:]
41    # witness functions
42    fx = K @ H @ alpha
43    gx = L @ H @ beta
44    # Corrections
45    correlation = np.corrcoef(fx, gx)
```

Method 1: QR decomposition on matrix

Here are the plots via method 1.



(a) $f(x)$ vs x



(b) $g(y)$ vs y

Figure 1: COCO: Plot of f and g under method 1

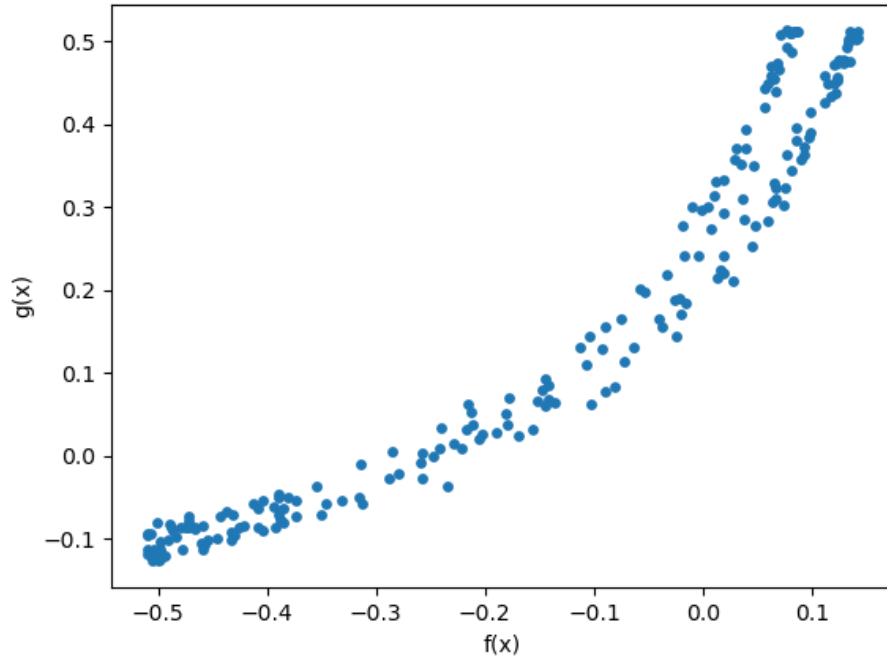


Figure 2: COCO: Plot of f vs g under method 1

Here is code for method 2. To avoid repeated code, we only include the changed part.

```

1 def coco_method2():
2     # same code as method 1
3     # incomplete cholesky on both K and L
4     R = incomplete_cholesky(K, eta)
5     S = incomplete_cholesky(L, eta)
6     # define zero blocks
7     m_R = R.shape[0]
8     m_S = S.shape[0]
9     z_R = np.zeros((m_R, m_R))
10    z_S = np.zeros((m_S, m_S))
11    # define matrix for eigen-problem
12    A = np.block([[z_R, R @ H @ S.T / N],
13                  [S @ H @ R.T / N, z_S]])
14
15    eigvals, eigvecs = np.linalg.eig(A)

```

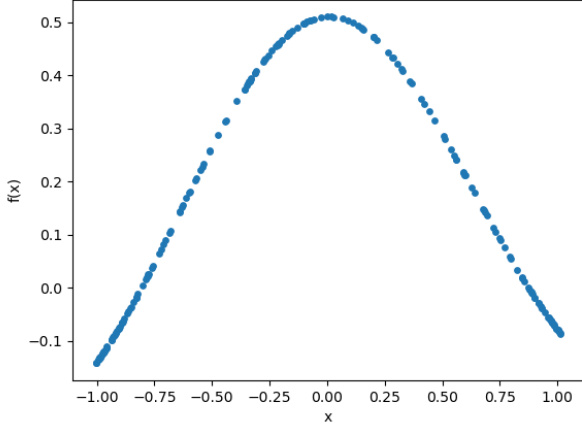
```

16 N_max = np.argmax(eigvals)
17 coco = np.abs(np.real(eigvals[N_max]))
18
19 eigvecs_v = eigvecs[:, N_max]
20 theta = eigvecs_v[:m_R]
21 eta = eigvecs_v[m_R:]
22
23 fx = theta.T @ R
24 gx = eta.T @ S
25 correlation = np.corrcoef(fx, gx)

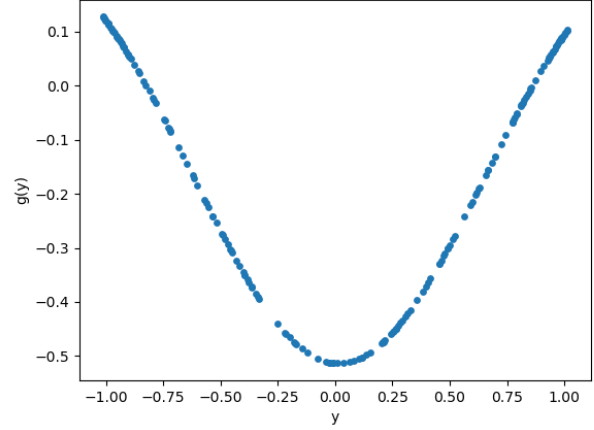
```

Method 2: Feature map projection

Here are plots from method 2.



(a) $f(x)$ vs x



(b) $g(y)$ vs y

Figure 3: COCO: Plot of f and g under method 2

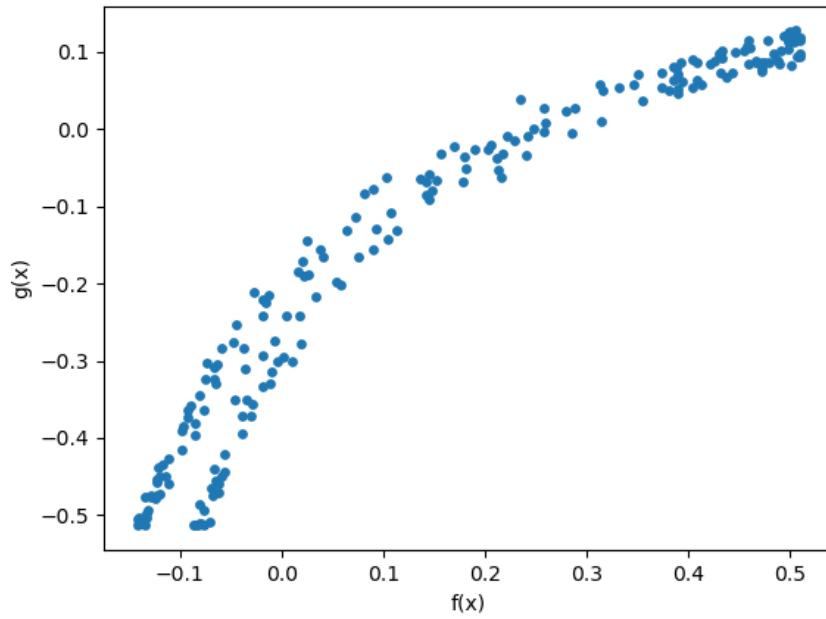


Figure 4: COCO: Plot of f vs g under method 2

In summary, we get same results from both method. The COCO value is 0.0949 and the correlation is 0.9503.

Question 2.2

We use the same method in the lecture notes. To maximize the canonical correlation with two constraints, we can write the corresponding Lagrangian:

$$\mathcal{L}(f, g, \lambda, \gamma) = -f^\top \widehat{C}_{XY} g + \frac{\lambda}{2} (f^\top \widehat{C}_{XX} f - 1) + \frac{\gamma}{2} (g^\top \widehat{C}_{YY} g - 1). \quad (17)$$

Let us write this in terms of the Gram matrices \widetilde{K} and \widetilde{L} :

$$\begin{aligned} f^\top \widehat{C}_{XY} g &= \frac{1}{n} \alpha^\top H X^\top (X H Y^\top) Y H \beta = \frac{1}{n} \alpha^\top \widetilde{K} \widetilde{L} \beta, \\ f^\top \widehat{C}_{XX} f &= \frac{1}{n} \alpha^\top H X^\top (X H X^\top) X H \alpha = \frac{1}{n} \alpha^\top \widetilde{K}^2 \alpha, \\ g^\top \widehat{C}_{YY} g &= \frac{1}{n} \beta^\top H Y^\top (Y H Y^\top) Y H \beta = \frac{1}{n} \beta^\top \widetilde{L}^2 \beta. \end{aligned} \quad (18)$$

Now, the Lagrangian is simplified as

$$\mathcal{L}(f, g, \lambda, \gamma) = -\frac{1}{n} \alpha^\top \widetilde{K} \widetilde{L} \beta + \frac{\lambda}{2} \left(\frac{1}{n} \alpha^\top \widetilde{K}^2 \alpha - 1 \right) + \frac{\gamma}{2} \left(\frac{1}{n} \beta^\top \widetilde{L}^2 \beta - 1 \right). \quad (19)$$

We must minimise the above Lagrangian with respect to α, β . Differentiating with respect to α and β and setting the resulting expressions to zero, we get

$$\begin{aligned} -\frac{1}{n} \widetilde{K} \widetilde{L} \beta + \frac{\lambda}{n} \widetilde{K}^2 \alpha &= 0, \\ -\frac{1}{n} \widetilde{L} \widetilde{K} \alpha + \frac{\gamma}{n} \widetilde{L}^2 \beta &= 0. \end{aligned} \quad (20)$$

Multiply the first equation by α^\top , and the second by β^\top ,

$$\begin{aligned} \alpha^\top \widetilde{K} \widetilde{L} \beta &= \lambda \alpha^\top \widetilde{K}^2 \alpha \\ \beta^\top \widetilde{L} \widetilde{K} \alpha &= \gamma \beta^\top \widetilde{L}^2 \beta \end{aligned} \quad (21)$$

This shows that we need to set $\lambda = \gamma$. Then, we need to solve following expression:

$$\begin{bmatrix} 0 & \widetilde{K} \widetilde{L} \\ \widetilde{L} \widetilde{K} & 0 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \gamma \begin{bmatrix} \widetilde{K}^2 & 0 \\ 0 & \widetilde{L}^2 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}. \quad (22)$$

When γ is at maximum, we obtain the solution for this eigenvalue problem.

Assume a Gaussian kernel, and that the points are also non-pathologically distributed so that K and L have full rank. This means that both \widetilde{K} and \widetilde{L} are invertible. Using $\lambda = \gamma$ in Equations 20, we multiply the first equation by \widetilde{K}^\top and the second by \widetilde{L}^\top

$$\begin{aligned} \widetilde{L} \beta - \gamma \widetilde{K} \alpha &= 0, \\ \widetilde{K} \alpha - \gamma \widetilde{L} \beta &= 0. \end{aligned} \quad (23)$$

Then, we combine above two equations. we can multiply the second equation in Equation 23 by γ

$$\gamma \widetilde{K} \alpha - \gamma^2 \widetilde{L} \beta = 0, \quad (24)$$

as $\gamma \widetilde{K} \alpha = \widetilde{L} \beta$, we can further deduce

$$\begin{aligned} \widetilde{L} \beta - \gamma^2 \widetilde{L} \beta &= 0, \\ \Rightarrow (1 - \gamma^2) \widetilde{L} \beta &= 0, \\ \Rightarrow 1 - \gamma^2 &= 0, \\ \Rightarrow \gamma &= \pm 1. \end{aligned} \quad (25)$$

Thus, we have a trivial solution for γ for any choice of α, β .

To solve this problem, one can add a regularization term into the constrains. The new constrains can be shown as:

$$\begin{aligned} \langle f, \widehat{C}_{XX} f \rangle_{\mathcal{F}} + \kappa \|f\|_{\mathcal{F}} &= 1, \\ \langle g, \widehat{C}_{YY} g \rangle_{\mathcal{G}} + \kappa \|g\|_{\mathcal{G}} &= 1. \end{aligned} \quad (26)$$

Combining the information from the question, we can write the new Lagrangian as

$$\mathcal{L}(f, g, \lambda, \gamma) = -\frac{1}{n} \alpha^\top \tilde{K} \tilde{L} \beta + \frac{\lambda}{2} \left(\alpha^\top \tilde{K}^2 \alpha + \kappa \alpha^\top \tilde{K} \alpha - 1 \right) + \frac{\gamma}{2} \left(\beta^\top \tilde{L}^2 \beta + \kappa \beta^\top \tilde{L} \beta - 1 \right), \quad (27)$$

where we remove $\frac{1}{n}$ factor in last two terms via rescaling. The following step is the same as above. Differentiating with respect to α and β and setting the resulting expressions to zero, we get

$$\begin{aligned} -\frac{1}{n} \tilde{K} \tilde{L} \beta + \lambda \left(\tilde{K}^2 \alpha + \kappa \tilde{K} \alpha \right) &= 0, \\ -\frac{1}{n} \tilde{L} \tilde{K} \alpha + \gamma \left(\tilde{L}^2 \beta + \kappa \tilde{L} \beta \right) &= 0. \end{aligned} \quad (28)$$

We multiply the first equation by α^\top , and the second by β^\top ,

$$\begin{aligned} \alpha^\top \frac{1}{n} \tilde{K} \tilde{L} \beta &= \lambda \alpha^\top \left(\tilde{K}^2 + \kappa \tilde{K} \right) \alpha, \\ \beta^\top \frac{1}{n} \tilde{L} \tilde{K} \alpha &= \gamma \beta^\top \left(\tilde{L}^2 + \kappa \tilde{L} \right) \beta. \end{aligned} \quad (29)$$

Then, we derive the form asked in the question

$$\begin{bmatrix} 0 & \frac{1}{n} \tilde{K} \tilde{L} \\ \frac{1}{n} \tilde{L} \tilde{K} & 0 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \lambda \begin{bmatrix} \tilde{K}^2 + \kappa \tilde{K} & 0 \\ 0 & \tilde{L}^2 + \kappa \tilde{L} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}. \quad (30)$$

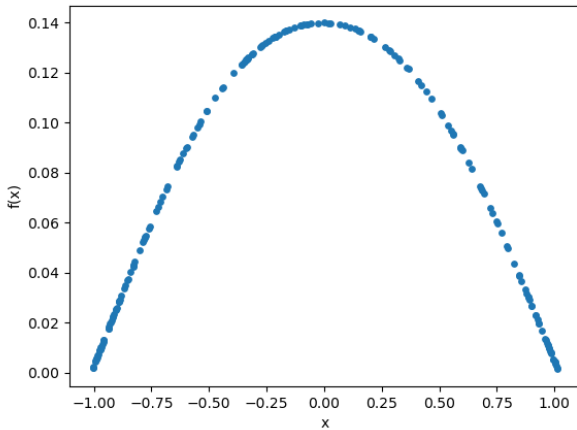
Implement the kernel CCA

In this section, we modify the previous code to incorporate the extra regularization term. Since the majority of the code remains unchanged, we will focus only on the modifications and explain the adjusted parts. For consistency with COCO, we will use two methods to implement the solution.

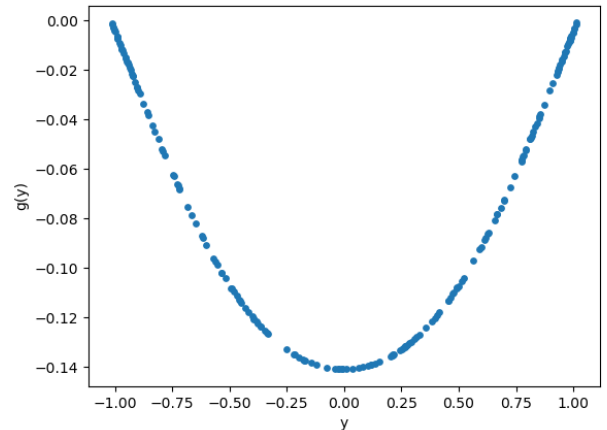
Let us begin with Method 1. This approach is straightforward to modify, as the changes are limited to the elements of the matrices. The rest of the code remains unchanged.

```
1 def cca():
2     # only difference part are shown
3     A = np.block([[z, K_tilde @ L_tilde / N],
4                   [L_tilde @ K_tilde / N, z]])
5     B = np.block([[K_tilde @ K_tilde + kapa * K_tilde, z],
6                   [z, L_tilde @ L_tilde + kapa * L_tilde]])
```

CCA Method 1



(a) $f(x)$ vs x



(b) $g(y)$ vs y

Figure 5: CCA: Plot of f and g under method 1

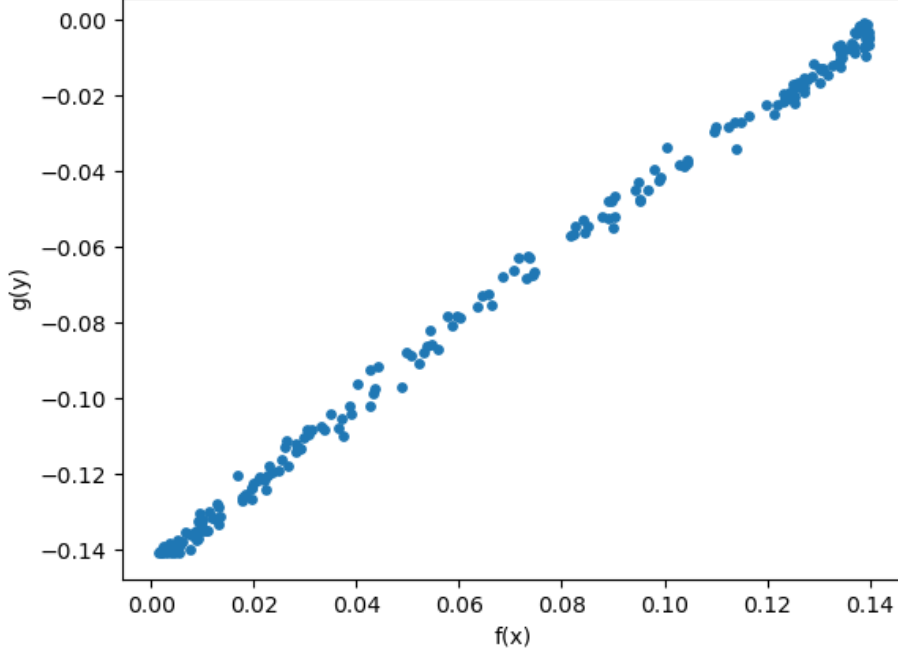


Figure 6: CCA: Plot of f vs g under method 1

For method 2, we need to rewrite our solution in new basis. We can start from Equation 29 and write it in terms of θ and η :

$$\begin{aligned} \frac{1}{n} \theta^\top R H S^\top \eta &= \lambda \theta^\top (R H R^\top + \kappa) \theta, \\ \frac{1}{n} \eta^\top S H R^\top \theta &= \gamma \eta^\top (S H S^\top + \kappa) \eta. \end{aligned} \quad (31)$$

Then, the solution is

$$\begin{bmatrix} 0 & \frac{1}{n} R H S^\top \\ \frac{1}{n} S H R^\top & 0 \end{bmatrix} \begin{bmatrix} \theta \\ \eta \end{bmatrix} = \gamma \begin{bmatrix} R H R^\top + \kappa & 0 \\ 0 & S H S^\top + \kappa \end{bmatrix} \begin{bmatrix} \theta \\ \eta \end{bmatrix}. \quad (32)$$

The corresponding code is

```

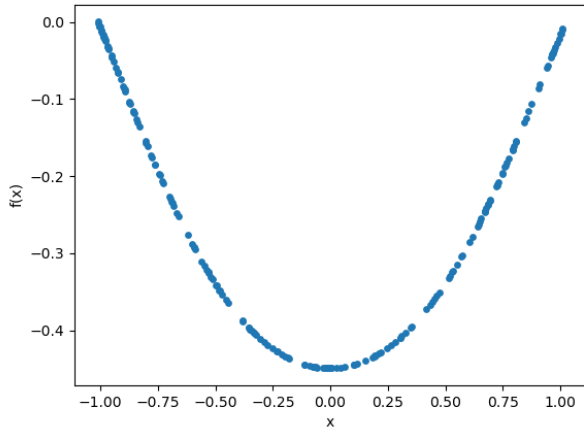
1 def cca_method2():
2     # only different part shown here
3     A = np.block([[z_R, R @ H @ S.T / N],
4                   [S @ H @ R.T / N, z_S]])
5
6     B = np.block([[R @ H @ R.T + kapa*np.eye(m_R), z_RS],
7                   [z_SR, S @ H @ S.T + kapa*np.eye(m_S)]])
8
9     C = np.linalg.inv(B) @ A
10    eigvals, eigvecs = np.linalg.eig(C)

```

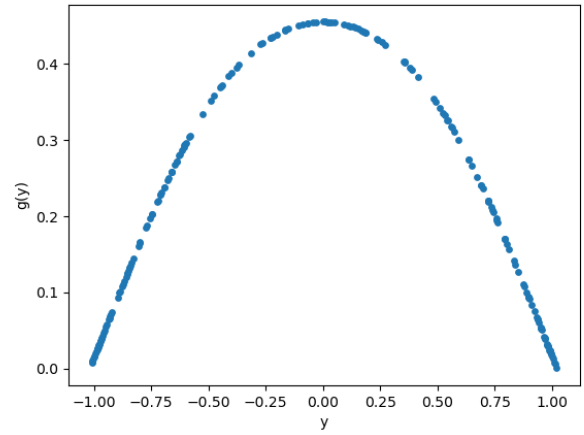
CCA Method 2

Here are the plots.

The COCO value is 0.9876 and the correlation is 0.9980. Comparing with previous plots and value of correlations, we find that witness functions $f(x)$ and $g(y)$ derived from KCCA are smoother. Also, we have a higher correlation than COCO. This improvement comes from the regularization terms which makes our code more stable.



(a) $f(x)$ vs x



(b) $g(y)$ vs y

Figure 7: CCA: Plot of f and g under method 2

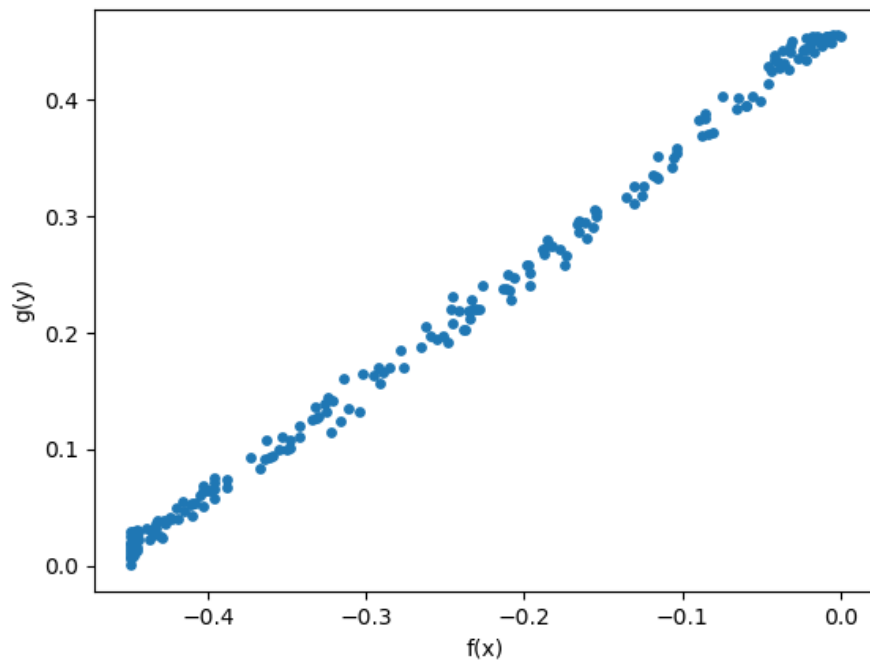


Figure 8: CCA: Plot of f vs g under method 2