

Part 1

● 2. Print the hyperparameters you found :

○ Params:

- ◆ $c = 3$, $\gamma = 0.002$ ($K = 5$)

○ Code :

◆ Reference

- ✧ In the [guide](#), it says that since doing a complete grid-search may still be time-consuming, it's recommended to [use a coarse grid first](#). After identifying a “better” region on the grid, [a finer grid search](#) on that region can be conducted.

Therefore, I conducted [2 grid search](#) in my code.

◆ Grid search function

```
1 def grid_search(kfold_data, cs, gammas):
2     cov = np.zeros((len(cs), len(gammas)))
3     max_acc = 0.0
4     best = {'c': -1, 'g': -1}
5     for i, c in enumerate(cs):
6         for j, g in enumerate(gammas):
7             # Taking average on the validation accuracy
8             total_err = 0.0
9             for kf in kfold_data:
10                 train_idx, val_idx = kf[0], kf[1]
11                 clf = SVC(C=c, kernel='rbf', gamma=g)
12                 clf.fit(x_train[train_idx], y_train[train_idx])
13                 y_pred = clf.predict(x_train[val_idx])
14                 total_err += sum(np.logical_xor(y_train[val_idx], y_pred))
15             avg_acc = 1.0 - (total_err / len(y_train))
16             if avg_acc > max_acc:
17                 max_acc = avg_acc
18                 best = {'c': c, 'g': g}
19             # save in cov-matrix
20             cov[i, j] = avg_acc
21     print(f"c : {best['c']}, gamma : {best['g']}, max_acc : {max_acc}")
22     return cov, best
```

◆ Do the grid search twice on {c1, gamma1} and {c2, gamma2} respectively.

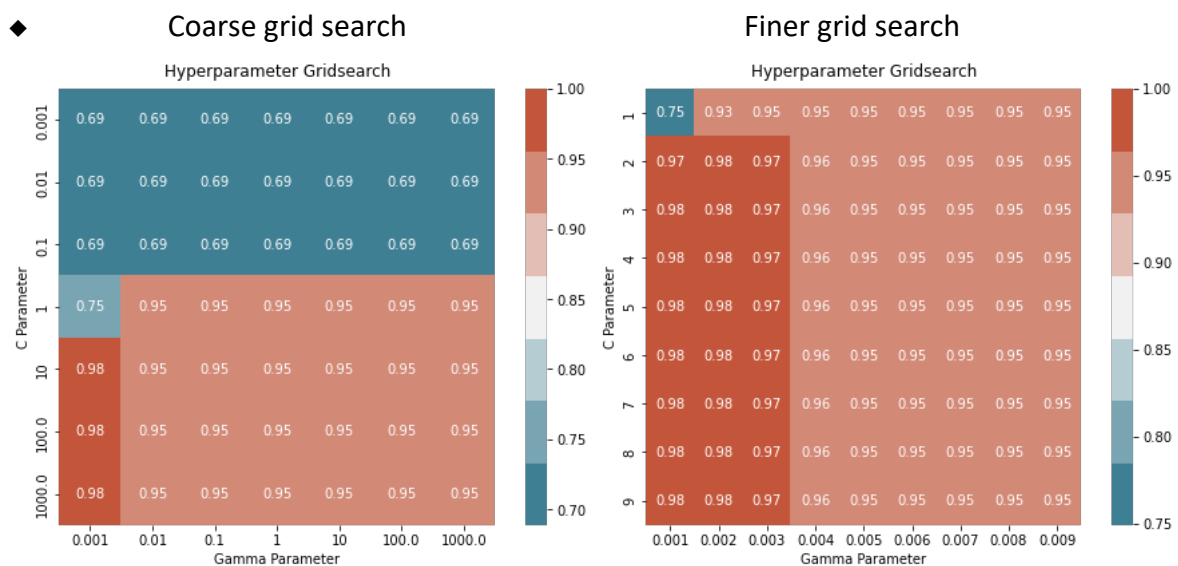
- ✧ np.round is used to fix the [floating point problem](#).
- ✧ Avg validation accuracy = 0.981

```
1 kfold_data = cross_validation(x_train, y_train, k=5)
2 # Coarse search
3 c1 = [1e-3, 1e-2, 1e-1, 1, 10, 1e2, 1e3]
4 gamma1 = [1e-3, 1e-2, 1e-1, 1, 10, 1e2, 1e3]
5 cov1, best1 = grid_search(kfold_data, c1, gamma1)
6
7 # Finer search
8 c2 = np.arange(1, 10, 1)
9 gamma2 = np.round(np.arange(0.001, 0.01, 0.001), 3)
10 cov2, best2 = grid_search(kfold_data, c2, gamma2)
11 pep8(_ih)

c : 10, gamma : 0.001, max_acc : 0.98
c : 3, gamma : 0.002, max_acc : 0.9818181818181818
```

- 3. Plot the grid search results of your SVM.

- Plot :



- Code :

- Heatmap function

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4
5 def heatmap_sns(c, gamma, cov):
6
7     plt.figure(figsize=(8, 6))
8     ax = sns.heatmap(
9         cov,
10         vmin=np.min(cov), vmax=1.0, annot=True,
11         xticklabels=gamma,
12         yticklabels=c,
13         cmap=sns.diverging_palette(220, 20, n=7),
14         square=True
15     )
16     ax.set_title('Hyperparameter Gridsearch', fontdict={'fontsize': 12}, pad=9)
17     ax.set_xlabel("Gamma Parameter")
18     ax.set_ylabel("C Parameter")
```

- Plot 2 maps

```
1 # Coarse grid search          1 # Fine grid search
2 heatmap_sns(c1, gamma1, cov1) 2 heatmap_sns(c2, gamma2, cov2)
```

- 4. Train your model by {C: 3, gamma: 0.002} on the whole training data and evaluate the performance on the test set.

- Result

- Test acc = 0.916, Val acc = 0.981 (see Q2)

```
1 # Train the model by {C: 3, gamma: 0.002} on the whole training data
2 best_model = SVC(C=best2['c'], kernel='rbf', gamma=best2['g'])
3 best_model.fit(x_train, y_train)
4
5 # Testing accuracy
6 y_pred = best_model.predict(x_test)
7 print("Accuracy score: ", accuracy_score(y_pred, y_test))
8 pep8(_ih)
```

Accuracy score: 0.9166666666666666

Part 2

Let's take a look at some properties of kernel first.

Given valid kernels $k_j(x, x')$ for $j \in \mathbb{N}$, the following kernels will also be valid.

① $K(x, x') = \alpha k_1(x, x')$ for $\alpha \geq 0$

Pf: For gram matrix $K = \alpha k_1$

$$\forall u \in \mathbb{R}^n, u^T K u = \alpha \cdot u^T k_1 u \geq 0$$

② $K(x, x') = k_1(x, x') + \alpha$ for $\alpha \geq 0$

Pf: Let ϕ_1 denote the feature map of k_1 .

Then, using the feature map $\Phi: x \mapsto [\phi_1(x), \sqrt{\alpha}]^T$,

$$\text{We have: } \langle \Phi(x), \Phi(x') \rangle = \langle \phi_1(x), \phi_1(x') \rangle + \alpha = k_1(x, x') + \alpha = K(x, x')$$

③ $K(x, x') = \sum_{j=1}^m \alpha_j k_j(x, x')$, $\alpha_j \geq 0$

Pf: For gram matrix $K = \sum_{j=1}^m \alpha_j k_j$

$$\forall u \in \mathbb{R}^n, u^T K u = \sum_{j=1}^m \alpha_j u^T k_j u \geq 0$$

due to the positivity of α_j and the validity of the kernels k_j

④ $K(x, x') = k_1(x, x') \odot k_2(x, x')$

Pf: By construction, the gram matrix is given by

$K = K_1 \odot K_2$, where \odot denotes the Hadamard (entrywise) product.

Given that K_1 and K_2 are symmetric positive semi-definite matrices, their eigendecompositions $K_1 = \sum_{i=1}^n \lambda_i u_i u_i^T$ and $K_2 = \sum_{j=1}^n \mu_j v_j v_j^T$ have positive eigenvalues $\lambda_i \geq 0$ and $\mu_j \geq 0$. This leads to:

$$\begin{aligned} K &= \sum_{i=1}^n \sum_{j=1}^n \lambda_i \mu_j (u_i u_i^T) \odot (v_j v_j^T) \\ &= \sum_{i=1}^n \sum_{j=1}^n \lambda_i \mu_j (u_i \odot v_j) (u_i \odot v_j)^T \\ &= \sum_{k=1}^{n^2} \gamma_k \omega_k \omega_k^T, \text{ where } \gamma_k = \lambda_{\lfloor \frac{k}{n} \rfloor} \mu_{k \bmod n} \geq 0 \end{aligned}$$

$$\text{and } \omega_k = u_{\lfloor \frac{k}{n} \rfloor} \odot v_{k \bmod n}$$

$$\text{Thus, } \forall u \in \mathbb{R}^n, u^T K u = \sum_{k=1}^{n^2} \gamma_k u^T \omega_k \omega_k^T u = \sum_{k=1}^{n^2} \gamma_k (\omega_k^T u)^2 \geq 0$$

⑤ $K(x, x') = q(k_1(x, x'))$, $q(\cdot)$ is a polynomial with nonnegative coefficients.

Pf: Suppose $q(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$, we now obtain

$$K(x, x') = a_n (k_1(x, x'))^n + a_{n-1} (k_1(x, x'))^{n-1} + \dots + a_1 k_1(x, x') + a_0$$

By repeatedly using rule ①, ③, ④, we can verify $K(x, x')$ a valid kernel.

⑥ $K(x, x') = \exp(k_1(x, x'))$

Pf: According to Taylor expansion, $K(x, x') = \sum_{n=0}^{\infty} \frac{1}{n!} (k_1(x, x'))^n$

By rule ⑤, we can prove $K(x, x')$ a valid kernel.

⑦ $K(x, x') = f(x) k_1(x, x') f(x')$

Pf: Since $k_1(x, x')$ is a valid kernel, it can be written as

$$k_1(x, x') = \phi_1(x)^T \phi_1(x')$$

We can obtain $K(x, x') = f(x) k_1(x, x') f(x')$

$$= [f(x) \phi_1(x)]^T [f(x') \phi_1(x')]$$

1.

$$(a) k(x, x') = (k_1(x, x')^2) + (k_1(x, x') + 1)^2$$

(i) $(k_1(x, x'))^2 = k_1(x, x')k_1(x, x')$, is a valid kernel by rule ④ ^{k·k}

(ii) $(k_1(x, x') + 1)$ is a valid kernel by rule ③ ^{k+α}
 $\Rightarrow (k_1(x, x') + 1)^2$ is a valid kernel by rule ④ ^{k·k}

(iii) $(k_1(x, x'))^2 + (k_1(x, x') + 1)^2$ is a valid kernel by rule ③ ^{Σαk} ■

$$(b) k(x, x') = (k_1(x, x'))^2 + \exp(\|x\|^2) * \exp(\|x'\|^2)$$

(i) $(k_1(x, x'))^2$ is a valid kernel by rule ④ ^{k·k}

(ii) $\exp(\|x\|^2) + \exp(\|x'\|^2)$ is a valid kernel

Pf: Using a explicit feature map $\phi: x \mapsto f(x) = e^{\|x\|^2}$, we have

$$e^{\|x\|^2} * e^{\|x'\|^2} = f(x)f(x') = \phi(x)^T \phi(x'), \text{ which is a valid kernel. } \Sigma \alpha k$$

(iii) $(k_1(x, x'))^2 + [\exp(\|x\|^2) \exp(\|x'\|^2)]$ is a valid kernel by rule ③ ■

2.

① If K is a positive semi-definite matrix, we have

$K = V\Lambda V^T$ where V is an orthonormal matrix with eigenvector v_t , and Λ is a diagonal matrix contains all eigenvalues $\lambda_t \geq 0$

Consider the feature map $\phi: x_n \mapsto (\sqrt{\lambda_t} v_{tn})_{n=1}^N \in \mathbb{R}^N$

we find that:

$$\phi(x_n)^T \phi(x_m) = \sum_{t=1}^d \lambda_t v_{tn} v_{tm} = (V\Lambda V^T)_{nm} = K_{nm} = k(x_n, x_m)$$

② By definition, we define the gram matrix $K = \Phi^T \Phi$, which is an N by N symmetric matrix with elements $K_{nm} = \phi(x_n)^T \phi(x_m) = k(x_n, x_m)$

Let $q(x) = x^T K x = x^T \Phi^T \Phi x = (\Phi x)^T (\Phi x) = \|\Phi x\|^2 \geq 0$ for $\forall x \in \mathbb{R}^N$.

Thus, K is positive semi-definite.

Therefore, K is positive semi-definite $\Leftrightarrow k(x_n, x_m)$ is a valid kernel. ■

3.

From setting $\frac{\delta}{\delta \omega} J(\omega)$ to zero, we have

$$\omega = \frac{-1}{\lambda} \sum_{n=1}^N \{ \omega^T \phi(x_n) - t_n \} \phi(x_n) = \sum_{n=1}^N a_n \phi(x_n) = \phi^T a, \quad \text{--- ①}$$

$$\text{where } \phi = \begin{bmatrix} \phi(x_n) \\ \vdots \end{bmatrix}^T, a = (a_1, \dots, a_N)^T \text{ with } a_n = \frac{-1}{\lambda} \{ \omega^T \phi(x_n) - t_n \}$$

$$a_n = \frac{-1}{\lambda} \{ \omega^T \phi(x_n) - t_n \}$$

$$= \frac{-1}{\lambda} \{ \omega_1 \phi_1(x_n) + \omega_2 \phi_2(x_n) + \dots + \omega_M \phi_M(x_n) - t_n \}$$

$$= -\frac{\omega_1}{\lambda} \phi_1(x_n) - \frac{\omega_2}{\lambda} \phi_2(x_n) - \dots - \frac{\omega_M}{\lambda} \phi_M(x_n) + \frac{t_n}{\lambda}$$

$$= (c_n - \frac{\omega_1}{\lambda}) \phi_1(x_n) + (c_n - \frac{\omega_2}{\lambda}) \phi_2(x_n) + \dots + (c_n - \frac{\omega_M}{\lambda}) \phi_M(x_n) \quad \text{--- ②}$$

$$\text{where } c_n = \frac{t_n}{\lambda} \cdot \frac{1}{\phi_1(x_n) + \phi_2(x_n) + \dots + \phi_M(x_n)}$$

Now, we know that a is a linear combination of $\phi(x_n)$

Therefore, for dual representation $J(a)$, we have:

$$J(a) = \frac{1}{2} a^T K K a - a^T K t + \frac{1}{2} t^T t + \frac{\lambda}{2} a^T K a$$

$$= \frac{1}{2} a^T \phi \phi^T \phi \phi^T a - a^T \phi \phi^T t + \frac{1}{2} t^T t + \frac{\lambda}{2} a^T \phi \phi^T a, \quad K = \phi \phi^T$$

$$= \frac{1}{2} \omega^T \phi \phi^T \omega - \omega \phi^T t + \frac{1}{2} t^T t + \frac{\lambda}{2} \omega^T \omega \quad \text{by}$$

$$= J(\omega)$$

4.

$$\textcircled{1} \|x - x'\|^2 = x^T x + x'^T x' - 2x^T x'$$

$$\textcircled{2} k(x, x') = \exp(-\|x - x'\|^2 / 2\sigma^2)$$

$$= \exp\left(\frac{-1}{2\sigma^2} x^T x\right) \exp(x^T x') \exp\left(\frac{-1}{2\sigma^2} x'^T x'\right)$$

$$= \exp\left(\frac{-x^T x}{2\sigma^2}\right) \sum_{k=0}^{\infty} \frac{(x^T x')^k}{k!} \exp\left(\frac{-x'^T x'}{2\sigma^2}\right)$$

$$= \sum_{k=0}^{\infty} \exp\left(\frac{-x^2}{2\sigma^2}\right) \frac{\sqrt{\frac{2^k}{k!}}}{\sqrt{\frac{2^k}{k!}}} (x^k)^T (x'^k) \exp\left(\frac{-x'^2}{2\sigma^2}\right), \quad x \in \mathbb{R}^1$$

$$= \phi(x)^T \phi(x'), \quad \text{where } \phi(x) = \exp\left(\frac{-x^2}{2\sigma^2}\right) \begin{bmatrix} \frac{1}{\sqrt{\pi}} x \\ \frac{1}{\sqrt{\pi}} x^2 \\ \vdots \end{bmatrix}$$

Alternative way to show that $k(x, x')$ is valid:

$$k(x, x') = \exp\left(\frac{-x^2}{2\sigma^2}\right) \exp(x^T x') \exp\left(\frac{-x'^2}{2\sigma^2}\right)$$

(i) $x^T x'$ is a linear kernel $\Rightarrow x^T x'$ is a valid kernel by rule $\textcircled{1}$

(ii) $\exp(x^T x')$ is a valid kernel by rule $\textcircled{6}$

(iii) $f(x) \exp(x^T x') f(x')$, $f(x) = \exp\left(\frac{-x^2}{2\sigma^2}\right)$ is a valid kernel by rule $\textcircled{7}$

5.

$$L(x, \lambda) = (x-2)^2 + \lambda[(x+3)(x-1)-2]$$

$$= (1+\lambda)x^2 + (2\lambda-4)x + (4-5\lambda)$$

$$\text{Set } \nabla L(x, \lambda) = 2x(1+\lambda) + (2\lambda-4) = 0$$

$$\Rightarrow x = \frac{2-\lambda}{(1+\lambda)}$$

$$L(x(\lambda), \lambda) = (\lambda-2)^2/(1+\lambda) + (2\lambda-4)(2-\lambda)/(1+\lambda) + (4-5\lambda)$$

$$= -(\lambda-2)^2/(1+\lambda) + (4-5\lambda)$$

$$= \begin{cases} \frac{-(\lambda-2)^2}{(1+\lambda)} + (4-5\lambda), & \lambda > -1 \\ -\infty, & \lambda \leq -1 \end{cases}$$

Thus, the dual problem is :

$$\text{maximize } -(\lambda-2)^2/(1+\lambda) + (4-5\lambda)$$

$$\text{subject to } \lambda \geq 0$$

Check :

① The optimal value p^* to the original problem is

$$\text{constraint: } x^2 + 2x - 5 \leq 0 \Rightarrow -1\sqrt{6} \leq x \leq 1+\sqrt{6}$$

$$\text{Thus, } p^* = (-1+\sqrt{6}-2)^2 = (-3+\sqrt{6})^2 = 15-6\sqrt{6}$$

② The optimal value p^* to the dual problem is :

$$(i) L = -(\lambda-2)^2/(1+\lambda) + (4-5\lambda)$$

$$\frac{\partial L}{\partial \lambda} = 0 = \frac{-2(\lambda-2)(1+\lambda) + (\lambda-2)^2}{(1+\lambda)^2} - 5$$

$$\Rightarrow (\lambda-2)(-\lambda-4) - 5\lambda = 10\lambda - 5 = 0, \lambda \neq -1$$

$$\Rightarrow -6\lambda^2 - 12\lambda + 3 = 0, \lambda \neq -1$$

$$\Rightarrow 2\lambda^2 + 4\lambda - 1 = 0, \lambda \neq -1$$

$$\Rightarrow \lambda = \frac{-4 \pm 2\sqrt{6}}{4} = -1 \pm \frac{\sqrt{6}}{2} \text{ since } \lambda \geq 0$$

(ii)

$$\text{For } \lambda = -1 + \frac{\sqrt{6}}{2}, L = p^*$$

$$L = -(-3 + \frac{\sqrt{6}}{2})^2 / \frac{\sqrt{6}}{2} + (4 - \frac{5}{2}\sqrt{6})$$

$$\frac{2}{\sqrt{6}} = (-\frac{3}{2} + 3\sqrt{6}) / \frac{\sqrt{6}}{2} + 4 - \frac{5}{2}\sqrt{6}$$

$$= \frac{-9}{2\sqrt{6}} + 6 + 9 - \frac{5}{2}\sqrt{6}$$

$$= 15 - 6\sqrt{6}$$