

NYCU Machine Learning HW3

312605001

機器人碩士 歐庭維

Part 1 : Linear SVM

1. Initialization: C = 1

```
alpha = [1.      0.0667 1.      0.      1.      0.      1.      0.      0.      0.
0.      0.      0.      1.      0.      0.      0.0667 0.      0.0667 0.
1.      0.      1.      0.      0.      0.      1.      0.      0.      0.
0.      1.      0.      0.      0.      1.      0.2     0.      1.      0.
0.      0.      0.      0.      1.      0.      1.      0.      1.      0. ]
alpha_sum = 14.4
b = [10.54, 10.54, 10.54, 10.54]
Mode: linear , Classification Rate (CR): 94.00%
```

2. Initialization: C = 10

```
alpha = [ 0.  0.  9.  0.  0.  0. 10.  0.  0.  0.  0.  0.  0.  0.  0.  0.
0.  0. 10.  0.  9.  0.  0.  0.  8.  0.  0.  0.  0. 10.  0.  0.
0.  0.  0.  0.  0.  0.  0. 10.  0.  0.  0. 10.  0.]
alpha_sum = 76.0
b = [15.14, 15.14, 15.14]
Mode: linear , Classification Rate (CR): 96.00%
```

3. Initialization: C = 100

```
alpha = [0.00000e+00 3.00000e-04 4.44438e+01 0.00000e+00 4.00000e-04 0.00000e+00
1.00000e+02 0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00 1.00000e-04
0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00 3.00000e-04 0.00000e+00
3.00000e-04 0.00000e+00 1.00000e+02 0.00000e+00 4.44438e+01 0.00000e+00
0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00 1.00000e-04 0.00000e+00
0.00000e+00 1.00000e+02 1.00000e-04 1.00000e-04 0.00000e+00 0.00000e+00
0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00 2.00000e-04
0.00000e+00 0.00000e+00 1.00000e+02 0.00000e+00 0.00000e+00 0.00000e+00
8.88884e+01 0.00000e+00]
alpha_sum = 577.7778
b = [11.004, 11.0044, 11.0041, 11.0038, 11.004, 11.004, 11.0044, 11.005, 11.0056, 11.0052, 11.0049, 11.0044]
Mode: linear , Classification Rate (CR): 92.00%
```

Part 2 : RBF kernel-based SVM

1. Initialization: C = 10, sigma = 5

```
alpha = [10.      8.9618 10.      0.      10.      0.      10.      0.      0.
0.      0.      0.      0.      10.      0.      0.      8.9618 0.
8.9618 0.      10.      0.      10.      0.      0.      0.      10.
0.      6.595  0.      0.      10.      0.      0.      0.      10.
10.      0.      10.      0.2905 0.      10.      0.      0.      10.
0.      10.      0.      10.      0.      ]
alpha_sum = 193.771
b = [0.0678, 0.0678, 0.0678, 0.0678, 0.0678]
Mode: rbf , Classification Rate (CR): 90.00%
```

```
alpha = [ 0.      0.      3.8295  0.      0.      0.      9.4708  0.      0.
  0.      0.      0.      0.      0.      0.      0.      0.
  0.      0.      10.     0.      3.8295  0.      0.      0.
  0.      0.      0.      0.      10.     0.      0.      0.
  0.      0.      0.      0.      0.      0.      0.      10.
  0.      0.      0.      7.1297  0.      ]
alpha_sum = 54.2594
b = [13.7329, 13.7329, 13.7329, 13.7328]
Mode: polynomial , Classification Rate (CR): 94.00%
```

3. Initialization: C = 10, p = 3

```
alpha = [ 0.      0.      3.8735  0.      0.      0.      8.5547  0.      0.
 0.      0.      0.      0.      0.      0.      0.      0.      0.
 0.      0.     10.      0.      3.8735  0.2324  0.      0.      0.
 0.      0.      0.      0.     10.      0.      0.      0.      0.
 0.      0.      0.      0.      0.      0.      0.      0.     10.
 0.      0.      0.      6.5342  0.      ]
alpha_sum = 53.0684
b = [6.4345, 6.4346, 6.4345, 6.4346, 6.4344]
Mode: polynomial , Classification Rate (CR): 96.00%
```

4. Initialization: C = 10, p = 4

```
alpha = [ 0.      0.      3.2052  0.      0.      0.      6.111  0.      0.
 0.      0.      0.      0.      0.      0.      0.      0.      0.
 0.      0.     10.      0.      3.2052  2.0531  0.      0.      0.
 0.      0.      0.      0.     10.      0.      0.      0.      0.
 0.      0.      0.      0.      0.      0.      0.      0.     10.
 0.      0.      0.      4.5745  0.      ]
alpha_sum = 49.1489
b = [4.6677, 4.6678, 4.6677, 4.668, 4.6676]
Mode: polynomial , Classification Rate (CR): 96.00%
```

5. Initialization: C = 10, p = 5

```
alpha = [ 0.      0.      6.7459  0.      0.      0.      0.      0.      0.
 0.      0.      0.      0.      0.      0.      0.      0.      0.
 0.      0.     10.      0.      6.7459  0.4632  0.      0.      0.
 0.      0.      0.      0.     10.      0.      0.      0.      0.
 0.      0.      0.      0.      0.      1.7114  0.      0.     10.
 0.      0.      0.      2.2436  0.      ]
alpha_sum = 47.91
b = [-2.8302, -2.8302, -1.2454, -7.9984, -3.5436]
Mode: polynomial , Classification Rate (CR): 34.00%
```

Part 4 : Discussion and results presenting

1. Linear SVM 與 kernel-based SVM 所訓練的 hyperplane 有何差異?

- Linear kernel-based SVM** 方法是指在特徵空間中找到最佳的超平面去將數據區分開，**Linear SVM** 的決策邊界在二維中 ($P = 1$) 時是一條直線，在三維中是平面 ($P = 2$)，在高維空間中是高維平面 ($P > 2$)。
- Kernel-based SVM** 方法是使用 **kernel function** 將原始的特徵空間 **mapping** 到高維度的空間，並在這個高維度空間中以一個線性超平面區分數據，與 **Linear SVM** 的超平面不同的是，**kernel-based SVM** 的超平面在原始的特徵空間中是非線性的，但是經過 **mapping** 到高維空間的時候會變成線性的，這樣使其能夠處理在原始特徵空間中非線性可分離的數據。

2. 隨著 kernel parameter 的改變，RBF kernel 與 polynomial kernel 所訓練的 hyperplane 可能有什麼變化? 其與分類率的變化有何關聯? 請嘗試解釋之。

- RBF kernel** 的參數 **gamma**(紅字手寫部分)，主要功能是控制 **RBF** 函數的形狀，當 **gamma** 變大時，代表所訓練的超平面可能會變得更複雜，有可能更好的擬合數據的分布狀況。相同的，當 **sigma** 越大，代表 **rbf kernel** 更平緩 (泛化能力強)，而在實驗中，當 **sigma = 1** 時的 **CR** 最佳，有 96%，但隨著參數的調整，也有可能由於考慮太多 **noise** 而產生 **overfitting** 的情況，使 **CR** 降低，可由 **part2** 結果得知。

② 輻射基底函數(Radial Basis Function kernel , RBF)

80% 作用

$$\begin{aligned} \rightarrow K(\mathbf{x}_i, \mathbf{x}_j) &= \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_i) \\ &= e^{\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}} = e^{-r \|\vec{x}_i - \vec{x}_j\|^2} \end{aligned}$$

其中 σ : parameter of the RBF kernel , $\sigma \in [-\infty, \infty]$ \Rightarrow tuning different σ values \Leftrightarrow choosing different ϕ

Note

- b. Polynomial kernel 的參數代表多項式的次數，同樣的也控制多項式的複雜程度，當 degree 變大時，代表 kernel 以及超平面會變得更複雜，能夠更好的擬合數據，但同樣也有可能會有 overfitting 的問題。由 part 3 的結果可以得知，當 $p = 1 \sim p = 4$ 時，CR 都大致能維持在 96%，但當 $p = 5$ 時 CR 掉到了 34%，我推測就是由於 overfitting 的發生導致 CR 的下降。

$$\rightarrow K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^P \text{ or } (\mathbf{x}_i^T \mathbf{x}_j)^P \equiv \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

3. 設定 kernel parameter 時，是否有方法避免 hyperplane 過度擬合(overfitting)的現象發生?若有請詳細討論。

通常在設定 kernel parameters 時，會透過大量調整不同的參數去找到最佳的 CR，但這通常代表要嘗試非常多次，當取的參數不夠多時，也無法得知當前看似最好的參數是不是 CR 表現的"global maximum"或是"local maximum"，或是 overfitting 發生的時機。為了找到最好的 CR，可以透過 Grid-search 以及交叉驗證來對幾乎所有的參數組合作迭代，這會花很多時間，但能確保能找到一組參數組合有最好的表現，當然也可以透過簡化模型或是增加數據量來避免 overfitting 的問題，但是這樣做的工作量會增加，以及能否真的能避免也會是個需要實驗的問題，這就看使用者的需求決定。