



Mohammad Javad Ranjbar

810101173

HW4

Machine Learning, Fall 2022

# سوال ۱:

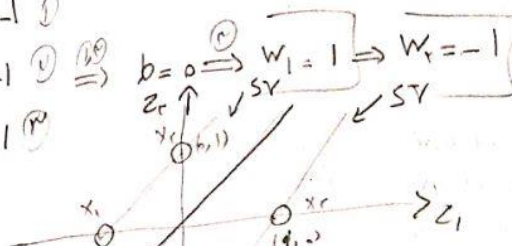
مسئله دواد (دوچهره) ۸/۱۰/۱۳  
تمرین سری چهارم یادگیری ماشین

می‌دانیم که معادله نقاط ردی (SV) به صورت زیر خواهد بود:  $x_1, x_2, x_3$  و  $x_4$  را ۱- در نظر بگیریم.

$$w^T x_1 + b = -1 \Rightarrow -w_1 + b = -1 \quad (1)$$

$$w^T x_2 + b = -1 \Rightarrow w_2 + b = -1 \quad (2) \quad \Rightarrow \quad b = 0 \Rightarrow w_1 = 1 \Rightarrow w_2 = -1$$

$$w^T x_3 + b = 1 \Rightarrow w_1 + b = 1 \quad (3)$$



مستقیم است که معادله خط جدا کننده به صورت  $Z_1 - Z_2 = 0$  خواهد بود در نتیجه داریم

$$W = [1, -1], \quad b = 0$$

با توجه به این مقادیر مقدار مارجین برابر با  $\frac{2}{\sqrt{2}}$  خواهد بود  
حال معادله لاگرانژین به صورت زیر خواهد بود

$$L(w, b, \alpha) = \frac{1}{2} w^T w + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b))$$

$$\frac{\partial L}{\partial w} = w + \sum_{i=1}^n \alpha_i (-y_i) x_i = 0 \Rightarrow w = - \sum_{i=1}^n \alpha_i y_i x_i \quad (1)$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^n \alpha_i y_i = 0 \quad (2)$$

Dual Problem به صورت زیر خواهد بود

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Subject to

$$\alpha_i \geq 0$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

حل برای مسائل مشابه

$$\begin{aligned} \alpha_1 + \alpha_2 + \alpha_3 &= 0 \rightarrow \alpha_1 - \alpha_3 = -1 \\ -\alpha_1 + 0 + \alpha_3 &= 1 \rightarrow \alpha_1 + \alpha_3 = -1 \\ 0 + \alpha_2 + 0 &= -1 \rightarrow \alpha_2 = -1 \end{aligned} \Rightarrow \alpha_1 = -1, \alpha_3 = 0 \Rightarrow \alpha = \begin{bmatrix} -1 \\ -1 \\ 0 \end{bmatrix}$$

سوال ۲:

(۲) این با توجه به خاصیت ضرب داخلی است که اگر  $\phi(x_i), \phi(x_j)$  برای  $x_i, x_j$  باشند مشخص است که ضرب داخلی آن خاصیت  $K(x_i, x_j) = \phi(x_i)^T \phi(x_j) = \phi(x_j)^T \phi(x_i)$  را دارا خواهد بود.

$$\begin{aligned} \|\phi(x_i) - \phi(x_j)\|^2 &= (\phi(x_i) - \phi(x_j))^T (\phi(x_i) - \phi(x_j)) = \phi(x_i)^T \phi(x_i) + \phi(x_j)^T \phi(x_j) - \phi(x_i)^T \phi(x_j) - \phi(x_j)^T \phi(x_i) \\ &= 2 - 2 e^{-\frac{1}{2}(x_i - x_j)^2} < 2 \end{aligned}$$

که با توجه به همیشه مثبت و نزولی بودن تابع  $e^{-x}$  کمین عبارت همیشه کوچکتر از ۲ خواهد بود.  
(ب) از آنجا که این نقطه از نقاط آموزش فعلی زیادی دارد معیار فاصله آن یعنی  $\|x_{far} - x_i\|^2$  مقدار بسیار زیادی خواهد داشت و همینطور پس  $e^{-\frac{1}{2}(x_{far} - x_i)^2}$  به سمت ۰ میل خواهد کرد پس داریم

$$\|x_{far} - x_i\|^2 \gg 0 \Rightarrow K(x_{far}, x_i) = e^{-\frac{1}{2}(x_{far} - x_i)^2} \simeq 0 \Rightarrow \sum y_i \alpha_i K(x_{far}, x_i) \simeq 0$$

پس داریم

$$\langle \hat{W}, \phi(x_i) \rangle + b = \sum y_i \alpha_i K(x_{far}, x_i) + b = f(x; \alpha, b) = b$$

سوال ۳:

۳

در SVM Soft Margin با  $C=1$  : برای شکل ۴ می باشد

این SVM با Soft Margin می باشد پس یا شکل ۲ می تواند باشد یا شکل ۴. حال با توجه به اینکه مقدار  $C$  نسبت به  $\lambda$  عدد توکلتری است یعنی جابجایی خطای بیشتر دایم پس باید شکلی را انتخاب کنیم که دو لای را آسانتر از جدا می کند که شکل ۴ است

Soft margin با  $C=1$  : برای شکل ۲ می باشد

این SVM دایم است پس یا شکل ۲ می تواند باشد یا شکل ۴. حال با توجه به مقدار بزرگ  $C$  باید اجازه دهی خطای کمی داشته باشیم در نتیجه شکلی را انتخاب می کنیم که دو لای را سختگیرانه تر جدا می کند و بیشترین شکل ۲ انتخاب می شود

Hard margin با کرنل  $K(x_i, x_j) = x_i \cdot x_j + (x_i \cdot x_j)^2$  : برای شکلی می باشد

از آنجاکه تابع کرنل می تواند دومی می باشد باید به شکل سهی دور باشد که تنها شکلی این را فراهم می کند شکل ۵ است

Hard margin با کرنل  $K(x_i, x_j) = \exp(-\lambda |x_i - x_j|)$  : برای شکل ۶ است

با توجه به غیر خطی بودن این دایمی سبز و زیاد بودن مقدار  $\lambda$  جابجایی که فاصله کمی از هم دارند مقدار کرنل بیشتر از حالت  $\lambda = \frac{1}{2}$  خواهد بود پس برای طبقه بندی نیاز به مقدار زیادی SVM نداریم پس سن شکل ۶ را باید شکل ۶ را انتخاب کنیم

Hard margin با کرنل  $K(x_i, x_j) = \exp(-\lambda |x_i - x_j|)$  : برای شکل ۱ است

این دایمی سبز غیر خطی بوده و باید یا برای شکل ۱ باشد یا ۶. حال با توجه به کم بودن مقدار  $\lambda = \frac{1}{2}$  یعنی برای فواصل نزدیک هم مقدار کرنل به شدت کم خواهد شد در نتیجه به SVM زیادی برای طبقه بندی نیاز داریم که بین شکل ۶، ۱ باید شکل ۱ را انتخاب کنیم.

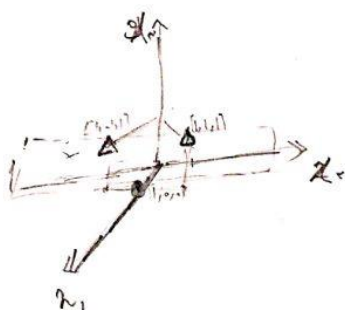
سوال ۴:

(۴)

الف) نقاط  $x_1$  با ابعاد گزین به شکل زیر در فضا آمده:

$$P(x_1) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad P(x_2) = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}, \quad P(x_3) = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

حل اگر دستگاه درجه سه کنیم



مشخص است که در شکل یک نقطه را می توان با یک معادله از یکدیگر جدا کرد پس خطی جدایی پذیر هستند  
ب) معادله لانه را می نویسیم

$$L(w, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (y_i (w^T \phi(x_i) + b) - 1)$$

حل برای یافتن شرط نسبت به آن مشتق گرفته و برابر با صفر می گذاریم:

$$\frac{\partial L(w, \alpha)}{\partial w} = w + \sum_{i=1}^n \alpha_i y_i \phi(x_i) = 0 \rightarrow$$

$$\frac{\partial L(w, \alpha)}{\partial b} = \sum_{i=1}^n \alpha_i y_i = 0$$

حل جابجایی می کنیم:

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + \alpha_1 x_1 \times \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \alpha_2 x_2 \times \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} + \alpha_3 x_3 \times \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 0$$



$$W_1 = \alpha_1 + \alpha_2 + \alpha_3 = 0 \quad (1)$$

$$W_2 = 0 = \alpha_2 + \alpha_3 = 0 \quad (2)$$

$$W_3 = 0 + \alpha_2 + \alpha_3 = 0 \quad (3)$$

$$-\alpha_1 + \alpha_2 + \alpha_3 = 0 \quad (4)$$

بر اساس ارقام داریم  $W_1 = 0$  که اکثر معادله قید را می‌دهد، می‌توانیم معادلات زیر را به دست می‌آوریم

$$-1 \times [W_1, W_2, W_3] \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + b = 1 \Rightarrow b = 1$$

$$1 \times (W_1^T \phi(2,1) + b) = 1 \Rightarrow$$

$$1 \times ([0, W_2, W_3] \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} + b) = 1 \Rightarrow W_2 + W_3 + b = 1$$

$$1 \times ([0, W_2, W_3] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + b) = 1 \Rightarrow W_2 + W_3 + b = 1$$

$$2W_3 = 2 \Rightarrow W_3 = 1$$

$$W_2 = 0$$

پس  $\hat{W}_2 = (0, 0, 2)$ ،  $b = 1$ ، و اندازه نقطه بین بردارها  $\frac{r_1}{|W_2|} = 1$  بود و مارجین بین  $\frac{1}{2}$  می‌باشد

سوال ۵:

ا)

یعنی ۳ تا باید بیشتر از ۵ باشد؟ جواب صحیح بدهند برابر ۵ است:  $P=0.5$

$$\frac{5+1}{2} = 3$$

$$\binom{3}{5} P^3 (1-P)^2 + \binom{4}{5} P^4 (1-P) + P^5 = 0.5187$$

$$\frac{9+1}{2} = 5$$

ب) مثال ۱۱ باید ۵ تا یا بیشتر پاسخ درست بدهند (یا برعکس ۵ تا یا کمتر پاسخ درست بدهند)

$$\binom{5}{9} P^5 (1-P)^4 + \binom{6}{9} P^6 (1-P)^3 + \binom{7}{9} P^7 (1-P)^2 + \binom{8}{9} P^8 (1-P) + P^9 =$$

$$\frac{9 \times 8 \times 7 \times 6 \times 5}{4! \times 5!} P^5 (1-P)^4 + \frac{9 \times 8 \times 7 \times 6 \times 5}{3! \times 6!} P^6 (1-P)^3 + \frac{9 \times 8 \times 7 \times 6 \times 5}{2! \times 7!} P^7 (1-P)^2 + 9 \times P^8 (1-P) + P^9 = 0.524$$

ج) در صورتی که  $N \rightarrow \infty$  مقدار دقت به یک نزدیک می‌شود. اما در واقعیت نمی‌توانیم به این دقت برسیم چون خطای طبقه‌بندی از هم مستقل نیستند

د)

$$\binom{3}{5} 0.5^3 (0.5)^2 + \binom{4}{5} 0.5^4 \times 0.5 + 0.5^5 = 0.5$$

در صورتی که دقت مدل‌های ضعیف برابر با نیم بوده باشد (رندوم) مجموعه‌ای آن به صورت ensemble نیز همچنان به رندوم انتخاب نگاه خواهد کرد و نتیجه بهتری نخواهد بود

## سوال ۶:

(الف)

کرنل‌های linear، poly، rbf، sigmoid و precomputed در این کتابخانه موجود هستند که استفاده‌ی هر کدام به شرح زیر است:

Linear kernel:

کرل خطی زمانی استفاده می شود که داده ها به صورت خطی قابل تفکیک باشند، یعنی با استفاده از یک خط می توان آنها را جدا کرد. این یکی از رایج ترین کرل هایی است که مورد استفاده قرار می گیرد. کرل خطی معمولاً در مجموعه داده هایی با مقادیر زیادی ویژگی استفاده می شود، زیرا افزایش ابعاد در این مجموعه داده ها لزوماً تفکیک پذیری را بهبود نمی بخشد.

#### Poly kernel:

در یک کرل **poly**، داده ها با استفاده از یک تابع چند جمله ای در فضایی با ابعاد بالاتر نگاشت می شوند. حاصل ضرب نقطه ای نقاط داده در فضای اصلی و تابع چند جمله ای در فضای جدید گرفته می شود. کرل چند جمله ای اغلب در مسائل طبقه بندی SVM استفاده می شود که در آن داده ها به صورت خطی قابل تفکیک نیستند. با نگاشت داده ها در فضایی با ابعاد بالاتر، کرل چند جمله ای گاهی اوقات می تواند ابرصفحه ای را پیدا کند که کلاس ها را از هم جدا می کند. این کرل به شکل زیر می باشد:

$$K(x, y) = (x^T y + 1)^d$$

#### Rbf kernel:

زمانی که مجموعه داده به صورت خطی غیر قابل تفکیک باشد یا به عبارت دیگر مجموعه داده غیر خطی باشد، توصیه می شود از کرلی مانند RBF استفاده شود. این کمک می کند تا زمانی که هیچ دانش قبلی از داده ها وجود ندارد، جداسازی مناسب انجام شود. RBF از منحنی های نرمال در اطراف نقاط داده استفاده می کند و آنها را جمع می کند تا مرز تصمیم را بتوان با یک نوع شرایط توپولوژی تعریف کرد، مانند منحنی هایی که مجموع آنها بالاتر از مقدار ۰/۵ است. این کرل معمولاً به صورت دیفالت انتخاب می شود و خاصیت هایی مثل نرم (smooth)، ایستا و.. دارد که بهتر از سایر کرل ها عمل می کند.

$$K(x, y) = \exp(-||x - y||^2 / (2\sigma^2))$$

#### Sigmoid kernel:

کرل Sigmoid از نظر تئوری برای یک ماشین بردار پشتیبان (SVM) پیشنهاد شده است، زیرا از یک شبکه عصبی سرچشمه می گیرد، اما تاکنون به طور گسترده در عمل مورد استفاده قرار نگرفته است. کرل Sigmoid معمولاً مشکل ساز یا نامعتبر است زیرا شرط Mercer را در همه  $\kappa$  و  $\theta$  برآورده نمی کند.

$$K(x, y) = \tanh(\kappa x^T y + \theta)$$

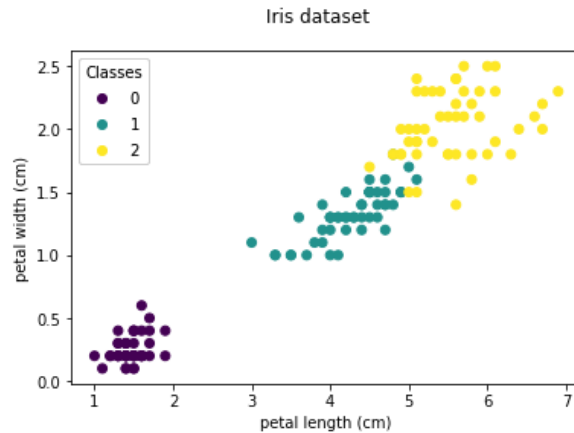
#### Precomputed kernel:

کرل های svm در این کتابخانه تا حد زیادی کند هستند. با استفاده از این کرل می توانیم کرل هایی که قبلاً حساب شده اند را برای آموزش و آزمون بدهیم تا هم از زمان و هم از حافظه برای آموزش مدل خود بهتر استفاده کنیم.

(ب)

داده های ما به صورت زیر پخش شده اند:





با توجه به شکل بالا می‌توانیم حدس بزنیم که کرنل خطی تا حد خوبی داده‌ها را جداسازی می‌کند.

حال داده‌ها را به دو دسته آموزش و آزمون تقسیم می‌کنیم (با نسبت ۰/۳۳) و برای هریک از کرنل‌های خواسته شده در صورت سوال اینکار را انجام می‌دهیم:

- SVM with Linear Kernel, one-vs-rest

The train accuracy for SVM with Linear Kernel, one-vs-rest= 0.91  
The test accuracy for SVM with Linear Kernel, one-vs-rest= 0.92

- SVM with Linear Kernel, one-vs-one

The train accuracy for SVM with Linear Kernel, one-vs-one= 0.95  
The test accuracy for SVM with Linear Kernel, one-vs-one= 1.0

- SVM with RBF Kernel, one-vs-rest

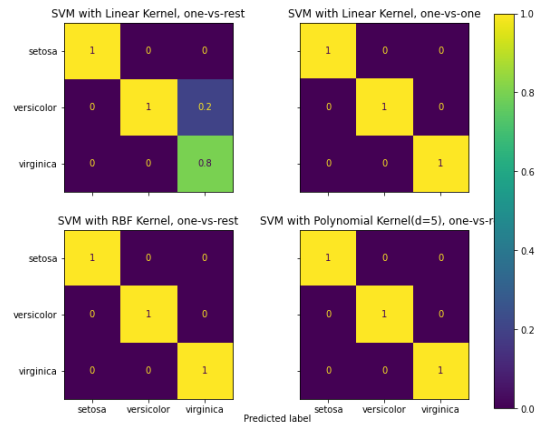
The train accuracy for SVM with RBF Kernel, one-vs-rest= 0.95  
The test accuracy for SVM with RBF Kernel, one-vs-rest= 1.0

- SVM with Polynomial Kernel (d=5), one-vs-rest

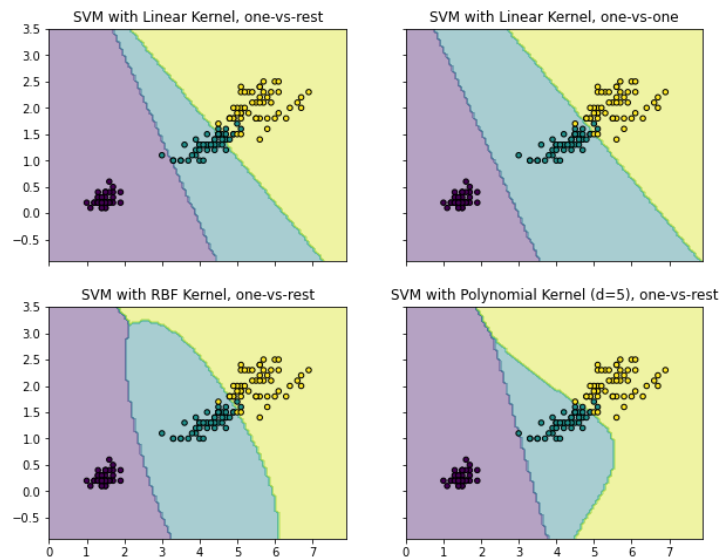
The train accuracy for SVM with Polynomial Kernel (d=5), one-vs-rest= 0.96  
The test accuracy for SVM with Polynomial Kernel (d=5), one-vs-rest= 1.0

دلیل اینکه دقت در داده‌های آموزش کمتر از داده‌های تست می‌باشد فقط به خاطر تصادفی بودن انتخاب این داده‌هاست و این دقت ممکن است تغییر کند.

ماتریس در هم‌ریختگی به شکل زیر خواهد بود:



و همچنین مرزهای تصمیم‌گیری به شکل‌های زیر خواهد بود:



مشخص است که مدل هر سه مدل خطی و poly و rbf نسبتاً خوب مرزهای تصمیم‌گیری را جداسازی کرده‌اند و با توجه به نتایج مدل خطی، نیاز به مدل‌های پیچیده‌تر همچون rbf نداریم. مدل poly نیز توانسته با یک شکل سهمی وار این نقاط را از هم جدا کند. ناحیه تصمیم‌گیری مدل one vs all نسبت به مدل one vs one کمی فشرده‌تر شده است. در این نوع تصمیم‌گیری چند کلاسه، هر بار یک لیبل مثبت را در برابر همه‌ی لیبل‌های منفی مقایسه می‌کنیم برای همین تعداد لیبل‌های منفی هر بار بیشتر از حالت عادی است که باعث فشرده‌تر شدن ناحیه شده است.

در اینجا نمونه‌ی آموزش برای یه نوع شافل دیگر از داده‌ها هم قرار می‌دهم:

- SVM with Linear Kernel, one-vs-rest

The train accuracy for SVM with Linear Kernel, one-vs-rest= 0.94

The test accuracy for SVM with Linear Kernel, one-vs-rest= 0.86

- SVM with Linear Kernel, one-vs-one

The train accuracy for SVM with Linear Kernel, one-vs-one= 0.98

The test accuracy for SVM with Linear Kernel, one-vs-one= 0.96

- SVM with RBF Kernel, one-vs-rest

The train accuracy for SVM with RBF Kernel, one-vs-rest= 0.97

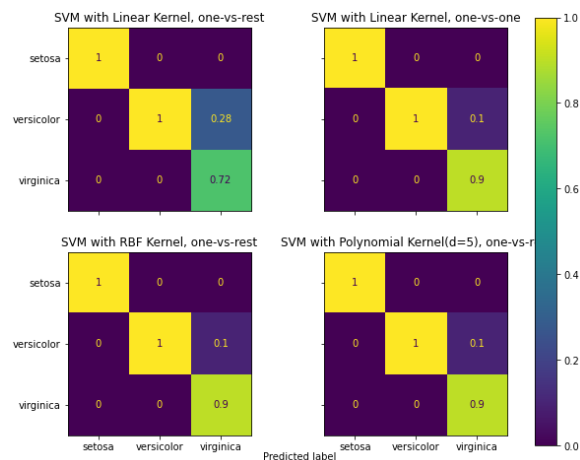
The test accuracy for SVM with RBF Kernel, one-vs-rest= 0.96

- SVM with Polynomial Kernel (d=5), one-vs-rest

The train accuracy for SVM with Polynomial Kernel (d=5), one-vs-rest= 0.95

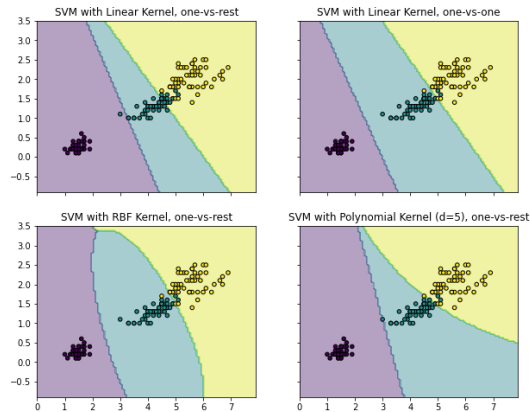
The test accuracy for SVM with Polynomial Kernel (d=5), one-vs-rest= 0.96

ماتریس درهم ریختگی:



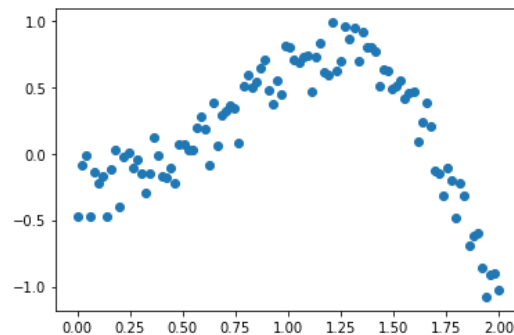
نتایج تقریباً همانند بالاست، فقط با توجه به داده‌ی تست جدید مدل در کلاس ۱ و ۲ (زرد و ابی) که داده‌ها نزدیک بهم هستند تعدادی داده را اشتباه طبقه بندی کرده است.

ناحیه تصمیم گیری:



## سوال ۷:

نمودار داده‌های ما به شکل زیر است:



حال برای کرنل خطی مقادیر  $C$  برابر با  $[1, 10, 100, 1000]$  برای مدل در نظر می‌گیریم.

همینطور برای کرنل چند جمله ای مقادیر زیر را در نظر می‌گیریم:

```
'C': [1, 10, 100, 1000], 'gamma': [0.0001, 0.001, 0.1, 1, 10],
'coef0': [0.1, 1, 5, 10]
```

Coef0 برای نامتقارن کردن کرنل استفاده می‌شود.

برای کرنل rbf مقادیر زیر را در نظر می‌گیریم:

```
'C': [1, 10, 100, 1000], 'gamma': [0.0001, 0.001, 0.1, 1, 10]
```

و مدل‌های خود را با کراس ولیدیشن ۵ آموزش می‌دهیم، که نتایج به صورت زیر خواهد شد:

```
The best parameters for linear are: {'C': 1, 'kernel': 'linear'}
The best parameters for poly are: {'C': 1, 'coef0': 0.1, 'gamma': 10,
'kernel': 'poly'}
The best parameters for rbf are: {'C': 10, 'gamma': 1, 'kernel': 'rbf'}
```

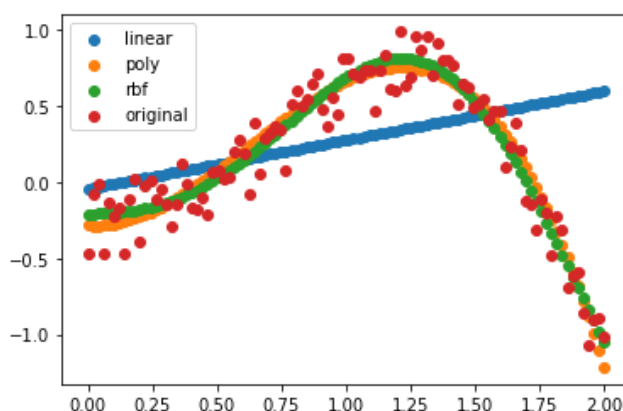
که نتایج مدل‌ها به صورت زیر خواهد بود:

The coefficient of determination for the linear classifier is: -  
 0.1765743645776896  
 The coefficient of determination for the poly classifier is:  
 0.9113147506100493  
 The coefficient of determination for the rbf classifier is:  
 0.9203043104715829

و همچنین خطا به صورت زیر است:

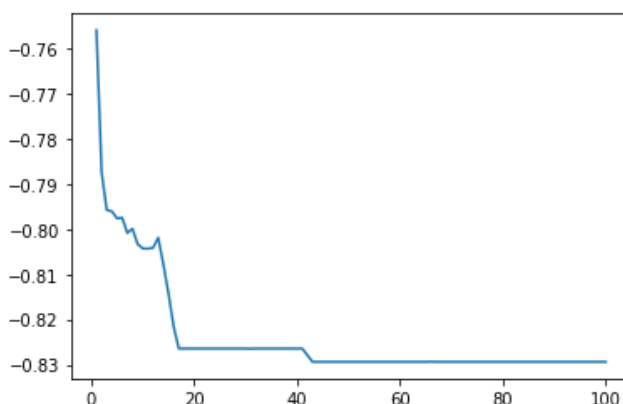
The mean\_squared\_error for the linear classifier is: 0.2817008417424394  
 The mean\_squared\_error for the poly classifier is: 0.021233429994248064  
 The mean\_squared\_error for the rbf classifier is: 0.019081108257409087

و نتایج در نمودار به صورت زیر خواهد بود:



با توجه به نمودار بالا و مقدار خطا و  $\text{coefficient of determination}$  مشخص است که کرنل  $\text{rbf}$  از دو کرنل دیگر بهتر عمل کرده است، که با توجه به شکل خاص  $\sin$  این پیش بینی میشود. همچنین کرنل  $\text{poly}$  از کرنل خطی بهتر عمل کرده و تقریباً به شکل  $\sin$  فیت شده است که این به خاطر ضریب  $\text{coef10}$  است که شکل منحنی را نامنتظران کرده است و توان سه قابلیت به فیت شدن به شکل بالا را داشت اما اگر بخش بزرگتر از  $\sin$  به عنوان داده‌ها انتخاب میشد،  $\text{poly}$  نیز بهتر نمیتوانست جواب بدهد. کرنل خطی نیز بدترین کرنل بوده و برای شکل  $\sin$  کافی نیست.

البته میتوان کارهای بالا را بدون استفاده از کتابخانه نیز انجام داد و نحوه کم شدن ارور را دید، برای مثال برای کرنل خطی بهترین  $c$  را با توجه به شکل زیر میتوان پیدا کرد:



که با توجه به نمودار بالا با افزایش  $c$  مقدار خطا در حال افزایش است پس  $c=1$  بهترین انتخاب است.