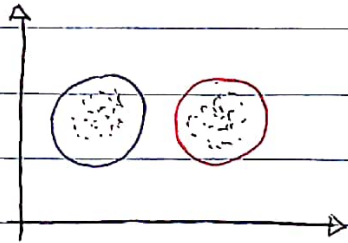
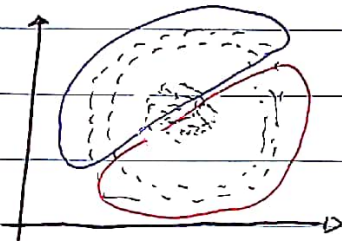


سوال ۱.۱

الف) کلاسترها به صورت زیر خواصند:



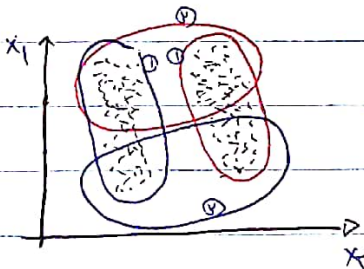
دو دسته نقطه به صورت جدا از هم به فضای شکل  
مایل دسته بندی شوند.



با توجه به تقارن شکل، اعداد clustering  
به دسته بندی به شکل مقایسه نمی از نقاط دو دسته  
به صورت شکل مانند شکل هستند.



به دلیل چگالی نقاط در سمت راست  
در آن ناحیه کلاستر تشکیل خواهد شد.  
و سایر نقاط که به این مرکز در ناحیه چپ نزدیک هستند  
نیز در این cluster قرار می دهند و شکل مانند شکل مایل  
ایجاد می شود.



با توجه به خواص انتخاب رده هم نقاط ادغام می شوند  
چونکه از حالات زیر یکی است اتفاق نیفتد.  
به یکی از حالت زیر cluster می شوند.

در صورت آخر feature scaling در راستای  $x_1$  باعث می شود که  
رابطه شود چپ نشیمن به شکل اول شود که نقاط به صورت کروی ترند و درشت  
حالت به ۲ در شکل آخر پیش نیاید.

سوال ۱.۲

در  $\leftarrow$  inertia فقط فاصله های درون کلاستر در نظر گرفته می شود، در صورتی که در silhouette score علاوه بر فواصل درون کلاستر، فاصله های بین کلاسترها نیز در نظر گرفته می شوند.

در inertia، بارش  $K$ ، خط به سمت تیرگی و کاهش میابد و این

خوبست زیرا به دنبال حالت بهینه  $K$  هستیم.

چنین silhouette score قابلیت تقسیم  $K$  را نیز فراهم کند.

اما در silhouette نمودار شکل خاص ندارد کامل معویه رفتار می کند که تصمیم گیری راحت تر باشد.

در ابتدا این نمودار خوب است که تمام cluster plot ها از خط  $\leftarrow$  باشند، که این موضوع در شکل های  $K=4$  و  $K=3$  رعایت نشود پس خوب نیست.

از بین نمودارهای باقی مانده  $K=5$  و  $K=4$ ، ما دانیم نمودار خوب است که اندازه کلاسترها به صورت یونیفرم باشند و هم اندازه پس نمودار سم با  $K=5$  مناسب ترین نمودار باشد در نتیجه تعداد خوشه ها = 5 انتخاب مناسب تر باشد.

سوال ۱.۳

active learning

حالتی است که در آن الگوریتم لرنینگ به صورت تعاملی و  $\leftarrow$  interactively

عمل می تواند  $\leftarrow$  فضات  $\leftarrow$  table جدید دهد در هر اجرا الگوریتم.

یک حالت خاص از semi-supervised محسوب می شود.

به این صورت که ابتدا خوشه بندی کنیم و به صورت  $\leftarrow$  interactively از

user می خواهیم به خوشه ها label دهد. اگر به خوشه ها label دهیم cluster

توانیم label با قطعیت بدهیم، به آن label می دهیم.

در ادامه با استفاده از data ها label شده مدل را train و

حالتی که ادامه می دهیم.

## Semi Supervised learning :

یک dataset نیمه نظارت شده که فقط بخشی از آن labeled و بقیه unlabeled است. در این روش، ما از داده‌های supervised و unsupervised برای آموزش استفاده می‌کنیم. این کار را با train کردن انجام می‌دهیم. فقط داده‌های supervised را train می‌کنیم و بقیه را به عنوان query در نظر می‌گیریم.

فرض کنید  $K$  خوشه داریم. query (همه داده‌های train کردن) که به این خوشه‌ها تعلق دارند. به این صورت عمل می‌کنیم که ابتدا به روش unsupervised  $K$  خوشه را پیدا می‌کنیم. سپس  $K$  خوشه را به روش supervised به این خوشه‌ها اختصاص می‌دهیم. در نهایت، به روش supervised  $K$  خوشه را به این خوشه‌ها اختصاص می‌دهیم. حال از هر خوشه یک نقطه به نام centroid (نقطه مرکزی خوشه) انتخاب می‌کنیم. و در نهایت به روش supervised train می‌کنیم.

## PCA whitening (نوع ۲)

در این روش، ما eigenvector ها را به روش supervised پیدا می‌کنیم. و سپس با این eigenvector ها dataset را به روش supervised train می‌کنیم. این کار را با eigenvector ها انجام می‌دهیم.  $X_{m \times n}$

ابتدا ماتریس  $X'$  را به صورت  $zero mean$  در می‌آوریم. حال به روش supervised  $X'$  را به این خوشه‌ها اختصاص می‌دهیم.

$$C = \frac{1}{m} X' X'^T$$

حال ماتریس کوواریانس را به صورت  $C$  قابل مشاهده است.

حال با استفاده از PCA (روش ~~ساده~~ <sup>تقریبی</sup>) eigenvector ها را به روش supervised پیدا می‌کنیم.

$$C = U \Sigma U^T$$

که  $U$  ماتریس eigenvector ها را به روش supervised پیدا می‌کنیم.

$$U = \begin{bmatrix} u_1 & \dots & u_n \\ 1 & & \end{bmatrix}$$





و ماتریس قطری  $\Sigma$  شامل eigen value ها بر روی قطر اصلی آن باشد.

بنابراین ماتریس  $\tilde{X}$  را به صورت  $\tilde{X} = U X'$  حساب میکنیم که در واقع  $\tilde{X}$  یک بردار  $k \times n$  است.  $\tilde{X}$  یک بردار  $k \times n$  است. eigen vector نامیده می شود.

حال برای PCA یعنی کاهش بعد ماتریس  $\tilde{X}$  و  $k$  عدد دلخواه را در نظر میگیریم. و مقدار  $k$  را مشخص می کنیم. که آن ها در نظر میگیریم. ماتریس  $\tilde{X}$  را به صورت زیر با بعد  $k \times n$  در نظر میگیریم.

$$\tilde{X} = \begin{bmatrix} \tilde{X}_{(1)} \\ \vdots \\ \tilde{X}_{(n)} \end{bmatrix} = \begin{bmatrix} \tilde{X}_{(1)} \\ \vdots \\ \tilde{X}_{(k)} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = D \quad \tilde{X} = \begin{bmatrix} \tilde{X}_{(1)} \\ \vdots \\ \tilde{X}_{(k)} \end{bmatrix}_{k \times n}$$

حال برای اندازه variance هر feature را می توانیم به دست آوریم. آن را تقسیم میکنیم. حال دیگر می بینیم  $X_{whitened}$  را به صورت زیر تعریف میکنیم.

$$X_{whitened} = \Sigma^{-1/2} (U^T X')$$

زیرا  $\Sigma^{-1/2}$  معکوس ماتریس قطری شامل مقادیر eigen value ها باشد. که به ماتریس  $\tilde{X}$  اعمال شود.

$$cov(X_{whitened}) = (\Sigma^{-1/2} U^T X') (\Sigma^{-1/2} U^T X')^T$$

حال ثابت می کنیم  $I = cov(X_{whitened})$  به این صورت.

$$X X^T = \Sigma^{-1/2} U^T X' X'^T U \Sigma^{-1/2}$$

$$= \Sigma^{-1/2} U^T \underbrace{X' X'^T}_I U \Sigma^{-1/2}$$

$$= \Sigma^{-1/2} \Sigma \Sigma^{-1/2} = I$$

از این نتیجه میگیریم که در نظر میگیریم.  $\Sigma^{-1/2}$  یکین می باشد.

برای به دست آوردن نقطه optimal در Ridge Regression از تابع  $loss$  آن مشتق می‌گیریم و برابر صفر قرار می‌دهیم.

$$\begin{aligned} L(w) &= \|xw - y\|_2^2 + \lambda \|w\|_2^2 \\ &= (xw - y)^T (xw - y) + \lambda w^T w \\ &= (w^T x^T - y^T) (xw - y) + \lambda w^T w \\ &= w^T x^T x w - w^T x^T y - y^T x w + y^T y + \lambda w^T w \end{aligned}$$

$$\frac{\partial L(w)}{\partial w} = x^T x w - x^T y + \lambda w$$

$$\begin{aligned} \frac{\partial L(w)}{\partial w} = 0 &\Rightarrow x^T y = x^T x w + \lambda w \\ (x^T x + \lambda I) w &= x^T y \\ w^* &= (x^T x + \lambda I)^{-1} x^T y \end{aligned}$$

در تابع نقطه optimal برابر خواص موجود است  $w^* = (x^T x + \lambda I)^{-1} x^T y$

شش داریم:  $L(w) = \|xw - y\|_2^2 + \lambda \|w\|_2^2$  می‌نویسیم

این  $L(w)$  را می‌توانیم به صورت  $\|B\|_2^2$  در نظر بگیریم به صورت

$$B = \begin{bmatrix} xw - y \\ \sqrt{\lambda} w \end{bmatrix}$$

باشد. حال خواص  $B$  را بنویسیم  $B = Aw - b$

$$B = \begin{bmatrix} x \\ \sqrt{\lambda} I \end{bmatrix} w - \begin{bmatrix} y \\ 0 \end{bmatrix} = Aw - b$$

که می‌خواهیم  $\|Aw - b\|_2^2$  را می‌نویسیم که  $LSE$   $\min_w$  باشد.  $A^T A w^* = A^T b$

پس داریم:



$$A^T A \omega^* = A^T b \quad A = \begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix} \quad b = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

$$A^T A = \begin{bmatrix} X^T & \sqrt{\lambda} I \end{bmatrix} \begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix} = X^T X + \lambda I$$

$$A^T b = \begin{bmatrix} X^T & \sqrt{\lambda} I \end{bmatrix} \begin{bmatrix} y \\ 0 \end{bmatrix} = X^T y$$

$$\Rightarrow A^T A \omega^* = A^T b \Rightarrow (X^T X + \lambda I) \omega^* = X^T y$$

$$\Rightarrow \omega^* = (X^T X + \lambda I)^{-1} X^T y$$

سوال (۴)

الف) dimensionality reduction باعث کاهش overfitting می‌شود.  
کاهش بعد باعث کاهش تعداد ویژگی‌ها می‌شود و کاهش تعداد ویژگی‌ها باعث  
کاهش بیش‌برازش می‌شود و overfitting را کاهش می‌دهد.  
بر این صورت که خصوصاً در PCA می‌تواند noise را از بین ببرد و مدل را بهتر کند  
درخت‌ها هم می‌تواند آنها را حذف کند.

ب) data augmentation methods

=> data augmentation for image classification  
یکی از ایجاد تغییر در تصویر است و می‌تواند به روش‌های مختلفی انجام شود.  
از دیگر استفاده‌ها generative adversarial net برای تولید تصویر مصنوعی است.

=> data augmentation for speech recognition  
تولید داده مصنوعی از MFCC ها و تغییراتی که بر روی آن‌ها اعمال می‌شود.  
تولید از شبکه‌ها از طریق انتقال یادگیری در داده‌های مصنوعی از طریق CNN در سطح  
کالندر RNN این‌ها است.

=> data augmentation for signal processing  
تولید generative adversarial network (GCGAN) برای  
تولید سیگنال‌ها و اکثر سیگنال‌های مصنوعی استفاده کرده و مطابق  
سیگنال‌ها در زمان باشد توسط تبدیل به یادگیری است.

استفاده از convolutional Neural Network بر روی سیگنال‌های احساسات است.  
بر EEG در سیگنال‌های احساسات می‌تواند.

مدل Open AEC-2 قادر به تولید سیگنال‌های بیولوژیکی مصنوعی مانند EEG  
و EMG می‌باشد.



✓



مثال ۵) مسئله یادداشت که  $n$  اشیاء داریم  $x, y \in \mathbb{R}^n$

۱.۵) ابتدا با توجه به اینکه  $\cos(\angle(x, y)) = \frac{\langle x, y \rangle}{\|x\| \|y\|}$  می توانیم تابع کرنل را به صورت زیر بازنویس کنیم

$$K(x, y) = (1 + \cos(\angle(x, y)))^r = 1 + 2\cos(\angle(x, y)) + \cos^2(\angle(x, y))$$

$$= 1 + 2 \frac{\langle x, y \rangle}{\|x\| \|y\|} + \frac{\langle x, y \rangle \cdot \langle x, y \rangle}{\|x\|^2 \|y\|^2}$$

$$= 1 + 2 \frac{\sum_{i=1}^n x_i y_i}{\|x\| \|y\|} + \frac{(\sum_{i=1}^n x_i y_i) (\sum_{j=1}^n x_j y_j)}{\|x\|^2 \|y\|^2}$$

Feature space  $\phi$ ، اگر بخواهیم به صورت زیر تعریف کنیم

$$\phi(x) = \begin{bmatrix} 1 \\ \sqrt{\frac{x_1}{\|x\|}} \\ \vdots \\ \sqrt{\frac{x_n}{\|x\|}} \\ \frac{x_1 x_1}{\|x\|^2} \\ \frac{x_1 x_2}{\|x\|^2} \\ \vdots \\ \frac{x_n x_1}{\|x\|^2} \end{bmatrix}$$

$\left. \begin{matrix} 1 \\ \sqrt{\frac{x_1}{\|x\|}} \\ \vdots \\ \sqrt{\frac{x_n}{\|x\|}} \end{matrix} \right\} = D$   
 $\left. \begin{matrix} \frac{x_1 x_1}{\|x\|^2} \\ \frac{x_1 x_2}{\|x\|^2} \\ \vdots \\ \frac{x_n x_1}{\|x\|^2} \end{matrix} \right\} = D \frac{x_i x_j}{\|x\|^2}$   
 $\Rightarrow \frac{x_i x_j}{\|x\|^2}$

$$\langle \phi(x) \cdot \phi(y) \rangle = 1 + \sum_{i=1}^n \sqrt{\frac{x_i}{\|x\|}} \sqrt{\frac{y_i}{\|y\|}} + \sum_{i=1}^n \sum_{j=1}^n \frac{x_i x_j}{\|x\|^2} \frac{y_i y_j}{\|y\|^2}$$

$$= 1 + 2 \frac{\sum_{i=1}^n x_i y_i}{\|x\| \|y\|} + \frac{\sum_i \sum_j (x_i y_i) (x_j y_j)}{\|x\|^2 \|y\|^2}$$

$$= 1 + 2 \frac{\langle x, y \rangle}{\|x\| \|y\|} + \frac{(\sum x_i y_i) (\sum x_j y_j)}{\|x\|^2 \|y\|^2}$$

$$= 1 + 2 \frac{\langle x, y \rangle}{\|x\| \|y\|} + \frac{\langle x, y \rangle \langle x, y \rangle}{\|x\|^2 \|y\|^2}$$

$$= K(x, y) \Rightarrow K(x, y) = \langle \phi(x) \cdot \phi(y) \rangle \quad \checkmark$$





(۵.۲.۱) ماتریس  $K$  در فضا  $\mathcal{H}$  مثبت معین است.  $\forall v \in \mathbb{R}^n$  ;  $v^T K v \geq 0$

به این صورت می‌توان نوشت:

$$v^T K v = [v_1 \dots v_m] \begin{bmatrix} c_1 & \dots & c_m \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix}$$

$$v^T K v = \begin{bmatrix} \langle v, c_1 \rangle & \dots & \langle v, c_m \rangle \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix}$$

$$v^T K v = \sum_{i=1}^m \langle v, c_i \rangle v_i$$

$$= \sum_{i=1}^m \left( \sum_{j=1}^m v_j K_{ji} \right) v_i$$

$$= \sum_{i=1}^m \sum_{j=1}^m v_i v_j K_{ji}$$

$$= \sum_{i=1}^m \sum_{j=1}^m v_i v_j \phi(x^{(j)})^T \phi(x^{(i)}) = \phi(x^{(j)})^T \phi(x^{(i)})$$

$$= \left\langle \left( \sum_{j=1}^m v_j \phi(x^{(j)}) \right), \left( \sum_{i=1}^m v_i \phi(x^{(i)}) \right) \right\rangle$$

$$= \left\| \sum_{i=1}^m v_i \phi(x^{(i)}) \right\|_2^2 \geq 0$$

$$\forall v \in \mathbb{R}^n ; v^T K v \geq 0$$

ماتریس  $K$  مثبت معین است.

$$\Rightarrow K \text{ مثبت معین است}$$



مثال ۵.۲.۲) خواهم نشان دهم اگر  $K$  یک ماتریس نیمه مثبت معین باشد.  
آنگاه میتوان نوشت  $K(x, y) = \langle \phi(x), \phi(y) \rangle$

اگر  $K$  ماتریس نیمه مثبت باشد آنگاه مقادیر منفی نیست و میتوان آن را مقدر کرد  
و بنویسیم  $K = P D P^T$  نوشت که در آن  $D$  ماتریس مقدر با مقادیر غیر  
صفری مقدر شد  $P$  ماتریس متعام باشد به طریقی دیگر

$$O = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_m \end{bmatrix} \quad , \quad P = \begin{bmatrix} | & | & & | \\ v_1 & v_2 & \dots & v_m \\ | & | & & | \end{bmatrix}$$

نام  $\phi(x)$  را تعریف کنیم طوری که  $K(x, y) = \langle \phi(x), \phi(y) \rangle$

ابتدا خواهیم نوشت  $K$  را به واسطه  $P D P^T$  بنویسیم

$$K = \begin{bmatrix} v_1 & \dots & v_m \end{bmatrix} \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_m \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_m^T \end{bmatrix}$$

$$K = \begin{bmatrix} \lambda_1 v_1 & \dots & \lambda_m v_m \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_m^T \end{bmatrix}$$

↓ ضرب در  $\lambda$    
  $K_{ij} = \langle r_i, c_j \rangle$

$$r_i = \begin{bmatrix} \lambda_1 v_i^{(1)} & \dots & \lambda_m v_i^{(m)} \end{bmatrix} \quad , \quad c_j = \begin{bmatrix} v_j^{(1)} \\ \vdots \\ v_j^{(m)} \end{bmatrix}$$

$$K_{ij} = \sum_{\ell=1}^m \lambda_{\ell} v_i^{(\ell)} v_j^{(\ell)} \quad (*)$$

حال کافی است  $\phi(x_i)$  را به صورت زیر تعریف کنیم:

$$\phi(x_i) = \begin{bmatrix} \sqrt{\lambda_1} v_i^{(1)} \\ \vdots \\ \sqrt{\lambda_m} v_i^{(m)} \end{bmatrix}$$

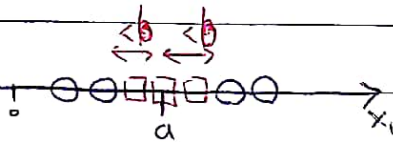
در این صورت داریم:

$$\langle \phi(x_i), \phi(x_j) \rangle = \sum_{\ell=1}^m \sqrt{\lambda_{\ell}} v_i^{(\ell)} \sqrt{\lambda_{\ell}} v_j^{(\ell)} = K_{ij} \quad (*)$$

$$\Rightarrow \langle \phi(x_i), \phi(x_j) \rangle = K_{ij} \quad \checkmark$$

(7.013)

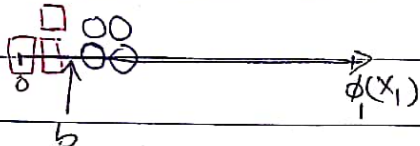
(4.1)



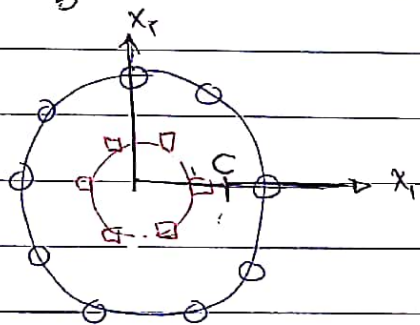
$\Phi(x)$  را به صورت زیر تعریف کنیم

$$\Phi_1(x) = |x - a|$$

حال کانتر است یک threshold ما  $b$  قرار بدهیم تا به راحتی cluster خود را در این صورت خواهیم داشت:



که داده ها که از مرکز دورند



ما فاصله نقاط از مرکز داده ها را cluster کنیم به این صورت که:

$$\Phi_2(x_1, x_2) = (x_1^2 + x_2^2)$$

حال کانتر است یک threshold ما  $c$  در نظر بگیریم. اگر

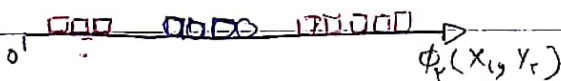
$$\Phi_2(x_1, x_2) < c$$
 داده ها که در یک مرکز قرار

$$\Phi_2(x_1, x_2) > c$$
 و داده ها که در یک مرکز قرار

در این حالت ترکیب از داده ها بالا باشد به این صورت که این داده ها

این داده ها  $\Phi_2$  و داده ها باشد به این صورت که این داده ها

از فاصله از مرکز است که



حال با این  $\Phi_2$  به این صورت که این داده ها

پس با این  $\Phi_2$  در این حالت به هم  $\Phi_1(\Phi_2(x_1, x_2))$  خواهد بود و یک threshold

کاف است داده ها که از مرکز دورند





$$(x_1 - a)^T + (x_2 - b)^T - r^T = 0 \quad (4.2)$$

$$x_1^T - 2ax_1 + a^T + x_2^T - 2bx_2 + b^T - r^T = 0$$

$$(x_1^T) + (x_2^T) - 2a(x_1) - 2b(x_2) + (b^T + a^T - r^T) = 0$$

پس در feature space  $n$  به صورت  $(x_1^T, x_2^T, x_1, x_2)$  تعریف شود

مکانی در  $n$  به صورت یک معادله یک خط در فضای  $(x_1^T, x_2^T, x_1, x_2)$

مانند باجه  $(b^T + a^T - r^T)$  و  $(1, 1, -2a, -2b)$  و  $(1, 1, -2a, -2b)$  باشد

پس در این فضای  $n$  linearly separable است

(4.3)

مردانم در کدن خط در SVM، این دو صورت زیر باشد

$$\min \frac{1}{2} \|w\|^2 + C \sum (y^{(i)} (w^T x^{(i)} + b))$$

که در آن  $C$  در واقع hyperparameter است که کنترل کننده این است

چنان اندازه margin و دقت دسته بندی را با هم

هم  $C$  در برگیرنده یعنی به این که دسته ها درست cluster شده باشند  
در دسته درست اینصورت و در فاصله دسته از آن (میزان دقت دسته بندی)

قدرت دسته باشد. این به این بست به margin size هر دو

و هم  $C$  کوچکتر به بیشترین حاشیه این به این بست به margin size هر دو

است به این که در نمونه ها کمتر و اثر آن ها نیز در آن کمتر شود

linear

$$\begin{aligned} C=10 \Rightarrow \text{F} \\ C=1 \Rightarrow \text{B} \\ C=1 \Rightarrow \text{C} \end{aligned}$$

بنابراین دهات ها کردند خطی شکل با  $C$  بزرگ  $C$  در  $C$  بزرگ

شکل  $C$  کوچکتر حاشیه در دسته بزرگتر  $C$  یعنی  $C$  بزرگتر است

شکل  $C$  بزرگتر حاشیه را دارد در دسته  $C$  کوچکتر  $C$  یعنی  $C$  بزرگتر است

در آخر شکل  $C$  با حاشیه توسط خط مستقیم  $C$  متوسط یعنی  $C=1$  باشد

در بین RBF دام  $\lambda$  hyperparameter می‌باشد و در جزی جزو پارامترهای مدل است  
 به این صورت که هر چه  $\lambda$  بزرگتر باشد، مدل منظم‌تر (smoother) می‌شود.

RBF

$$\lambda = 0 \rightarrow \infty$$

$$\lambda = 1 \rightarrow D$$

$$\lambda = \infty \rightarrow 0$$

بین قدرت دام در سطح RBF kernel و  $(a, d, e)$  این شکلهای

بروز دایره  $d$  و  $e$  smooth تر می‌شوند و به سمت  $\lambda = 0$  می‌روند  
 بروز دایره  $a$  کمتر می‌شود و به سمت  $\lambda = \infty$  می‌روند  
 و در آخر بروز دایره  $e$  از یک smoothness در دایره  $d$  و  $e$  به سمت  $\lambda = 1$  می‌باشد

با  $\lambda$  بزرگتر می‌شود و به سمت  $\lambda = 1$  می‌باشد که به سمت margin می‌رود.

CA 013

اگر در رابطه با تاثیر اندازه training set بر bias و variance  
 هم اندازه training set بیشتر باشد، باین سبب می‌شود  
 مدل دقیق‌تر شود و در واقع اطلاعات جدید اضافه شده باعث ضعیف‌تر شدن  
 مدل می‌شود که در نتیجه باین سبب واریانس بیشتر می‌شود.

و در نتیجه اندازه  $k$  کمتر باین سبب واریانس کمتر می‌شود.

با افزایش اندازه training set بیش شباهت بین هوش مصنوعی  
 و بیش overfitting می‌شود و باین سبب کم واریانس زیاد می‌شود.

بجای  $k$  fold هر نقطه یکبار تست می‌شود و  $k$  دفعه train  
 می‌شود. در روش MCCV تقسیم بندی‌های بیشتری استفاده می‌شوند  
 البته برای این تقسیم بندی‌ها استفاده نمی‌شوند.

میانگین ریشه اشتباه در  $k$  fold با  $k$  بزرگتر می‌شود و باین سبب باین سبب زیاد  
 در  $k$  بزرگتر می‌شود و باین سبب باین سبب زیاد.

