

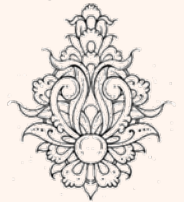
یادگیری ماشین  
(۰۱-۸۰۵-۱۱-۱۳)  
فصل سیزدهم



دانشگاه شهید بهشتی  
دانشکده‌ی مهندسی برق و کامپیوتر  
پاییز ۱۳۹۳  
احمد محمودی ازناوه

# فهرست مطالب

- SVM
  - تاریخچه
  - معرفی
  - داده‌های جدایی‌پذیر خطی
- Soft Margin
- مجموعه‌های جدایی‌ناپذیر خطی
  - نگاشت به فضای با ابعاد بالا
  - Inner product kernel
- مثال XOR
- Matlab در SVM

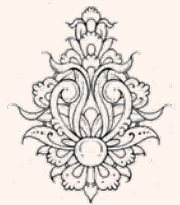


# تاریخچه

- نسخه‌ی اولیه‌ی SVM توسط آقای Vladimir Vapnik ارائه شد.
- Vapnik با همکاری خانم Corinna Cortes استاندارد کنونی SVM را در سال ۱۹۹۳ پایه‌ریزی کرده و در سال ۱۹۹۵ منتشر نمودند.

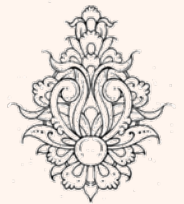
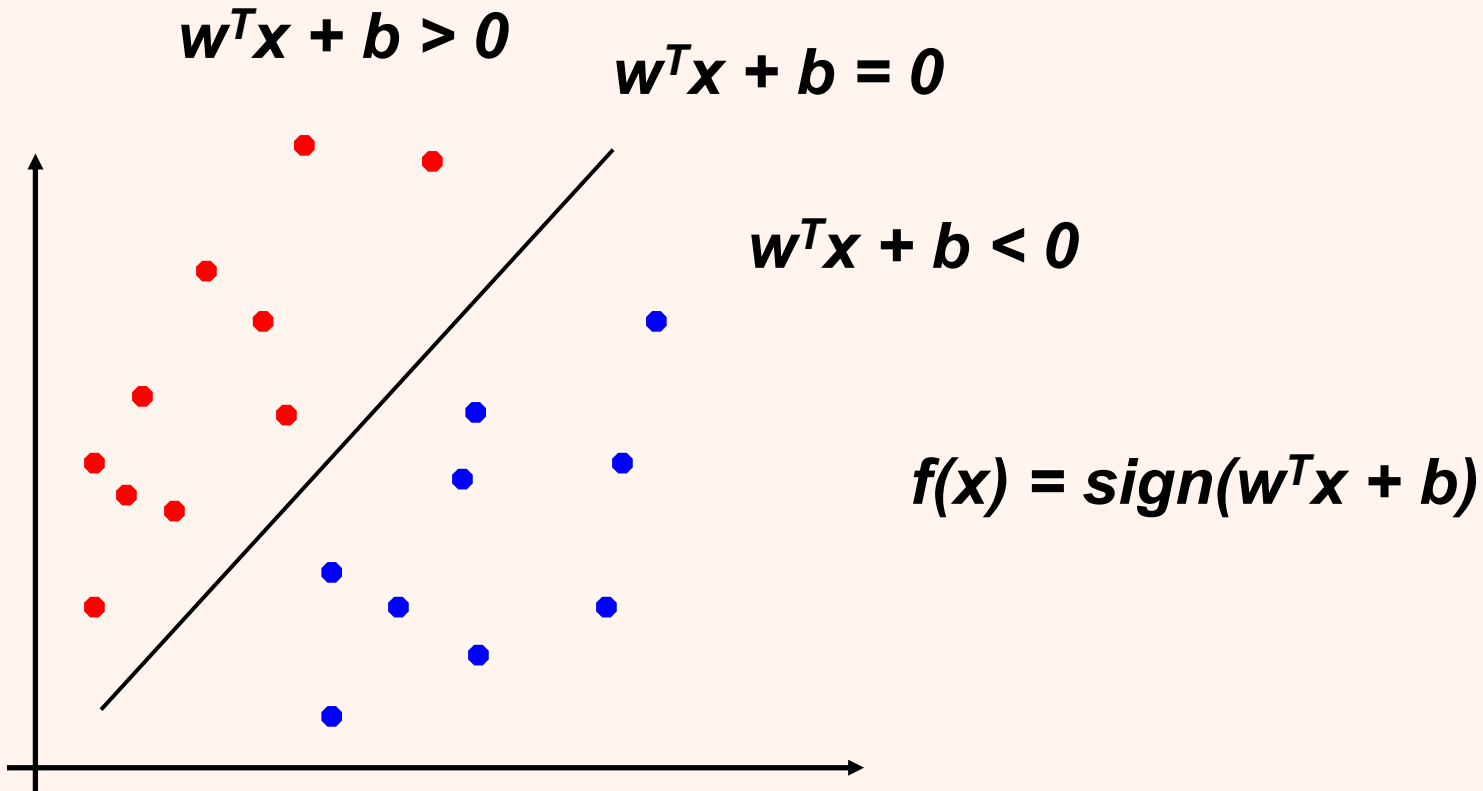


**Cortes, C. and V. Vapnik (1995). "Support-vector networks." Machine Learning 20(3): 273-297.**



# معرفی

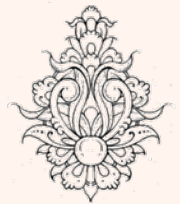
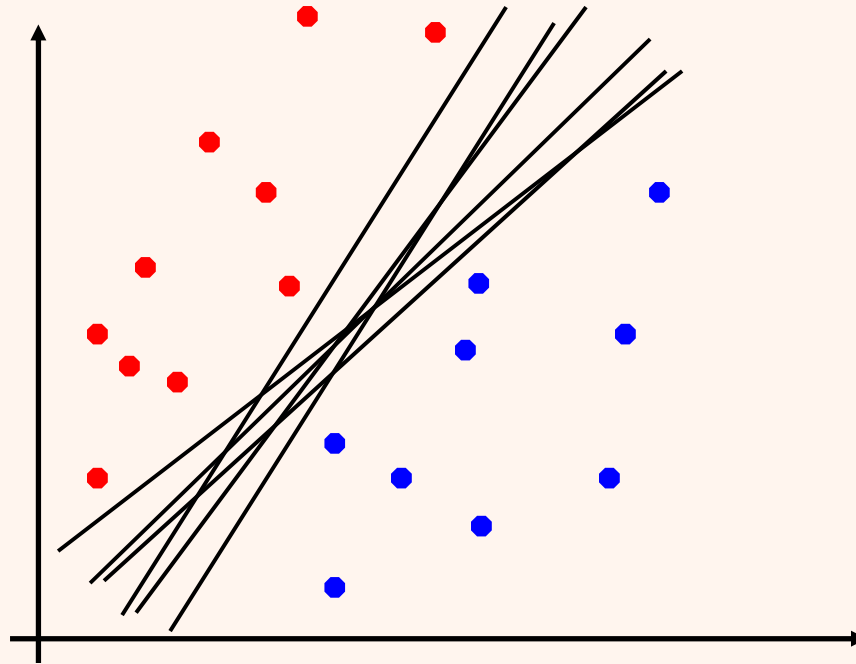
- یک جداکنندهی خطی را می‌توان همانند شکل زیر در نظر گرفت.



# مرز بهینه

## • سوال

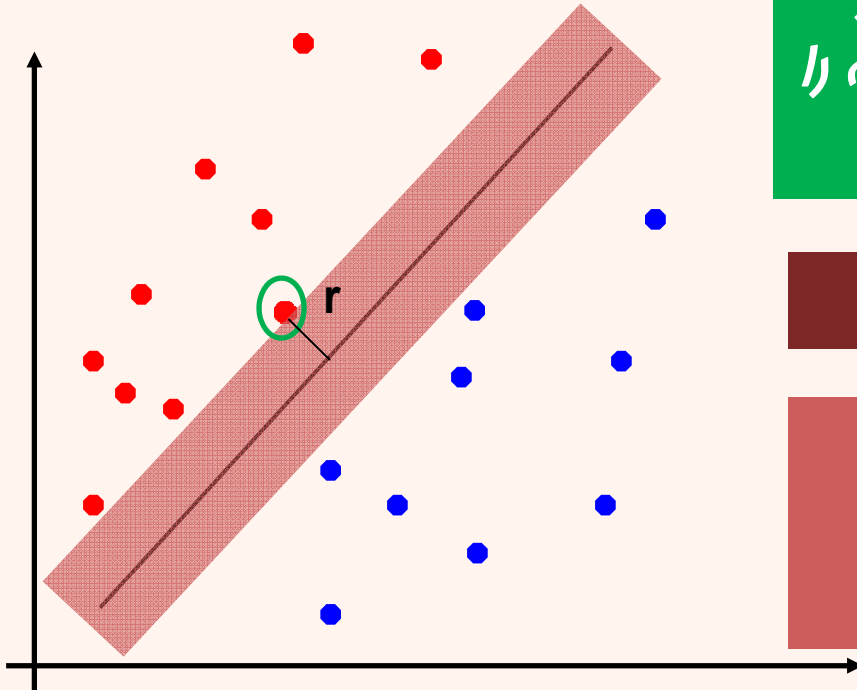
– کدام یک از مرزها، مرزی بهینه برای جداسازی است؟



# مرز جداسازی

- می‌خواهیم به گونه‌ای بهترین مرز جداسازی را به دست آوریم.

Margin of separation



فرض کنیم نزدیک‌ترین نقطه به مرز جداسازی در نظر گرفته شده و فاصله را  $r$  بنامیم.

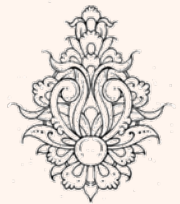
هدف ماکزیم نمودن  $r$  است.

یک ماشیه مشخص می‌کنیم هر مرزی که ماشیهی پهن‌تری را نتیجه دهد، بهتر است.

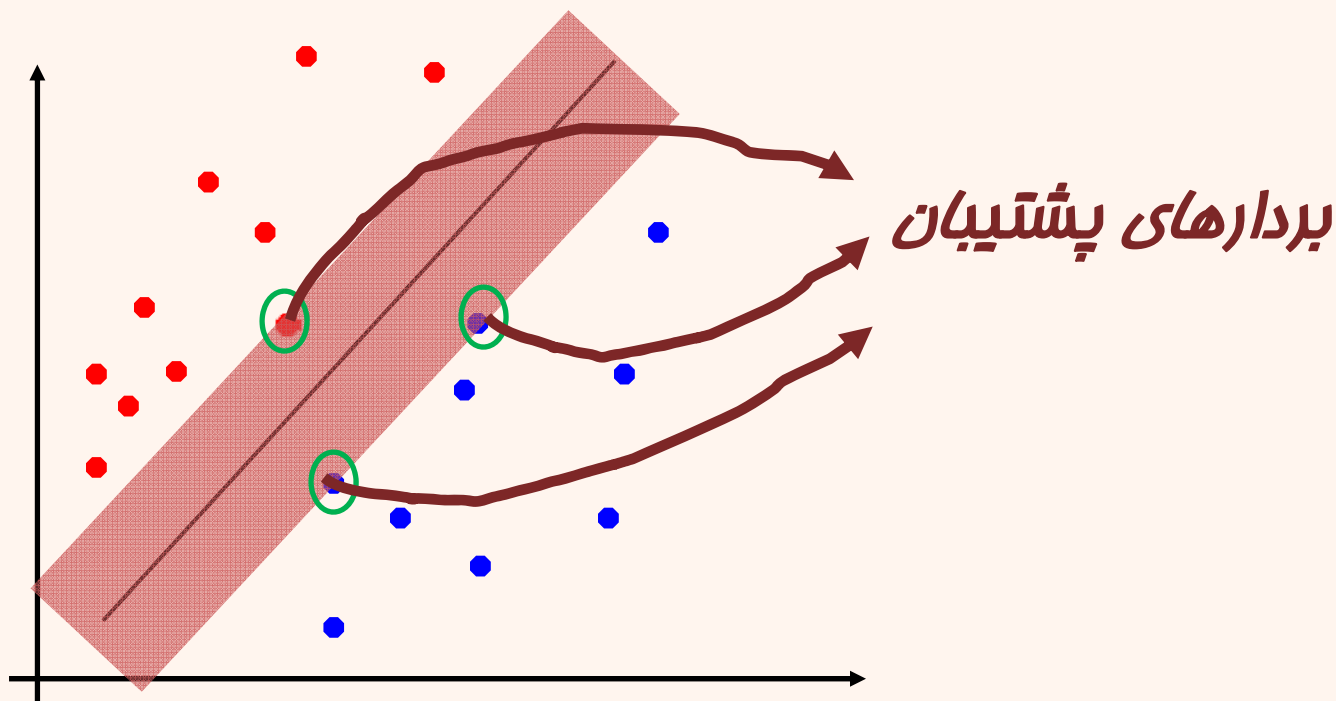
دانشگاه  
سپهر  
پهشتی

# ماشینی ماکزیمم

- ماکزیمم نمودن ماشینه (Margin) ایده‌ی خوبی است جهت جداسازی خطی، این شیوه را **LSVM** یا **Linear SVM** می‌نامند.
- در این حالت نمونه‌هایی که به روی مرز هستند، از اهمیت ویژه‌ای برخوردارند.
- بدین وسیله می‌توان از نمونه‌های دیگر صرف‌نظر کرد و تنها به نمونه‌های مهم روی مرز پرداخت.



- به نمونه‌های روی مرز «بردار پشتیبان» گویند.



Optimal hyperplane





# مرز جداسازی

- برای معادله‌ی مرز جداسازی داشته‌یم:

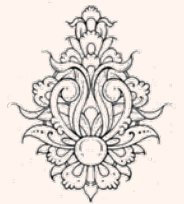
$$\mathbf{W}^T \mathbf{X} + b = 0$$

$$(\mathbf{x}_i, d_i = +1) \quad \mathbf{W}^T \mathbf{x}_i + b > 0$$

$$(\mathbf{x}_i, d_i = -1) \quad \mathbf{W}^T \mathbf{x}_i + b < 0$$

- فرض کنیم مرز بهینه توسط  $\mathbf{W}_{op}$  و  $b_{op}$  مشخص شود.

- فرض: فرض کنیم نزدیک‌ترین نقطه به مرز جداسازی را در نظر گرفته، فاصله را « $r$ » بنامیم.



# مرز جداسازی (ادامه...)

## • هدف

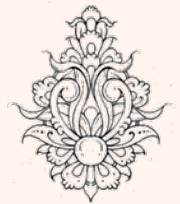
- ماکزیمم نمودن فاصله یا همان  $\rho$  است
- برای نقاط روی مرز جداسازی بهینه داریم:

$$\mathbf{W}_{op}^T \mathbf{X} + b_{op} = 0$$

- برای نقاط خارج از مرز داریم:

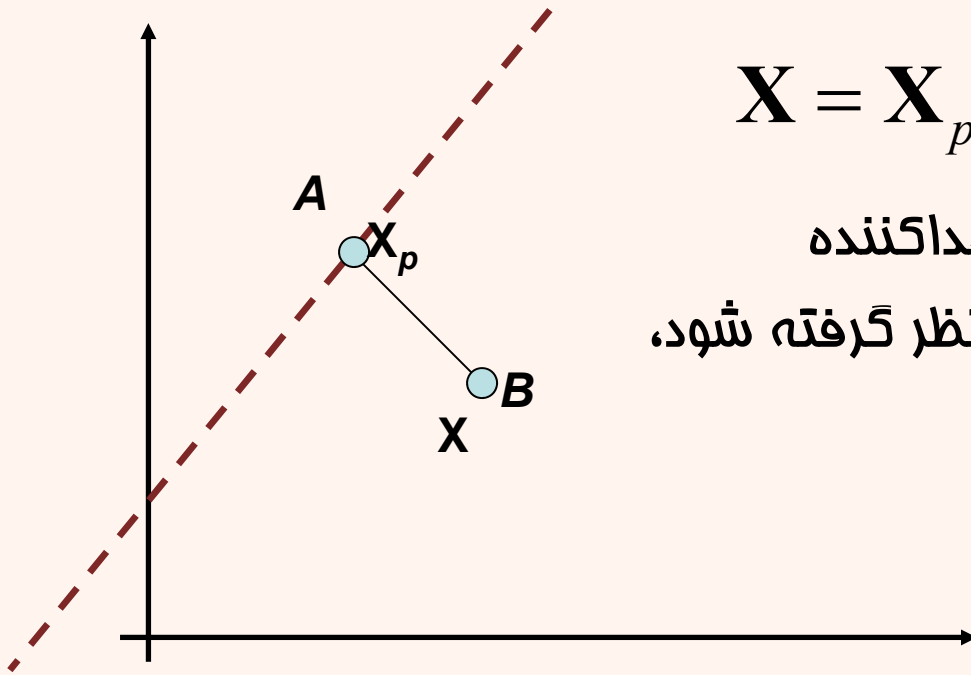
$$g(\mathbf{X}) = \mathbf{W}_{op}^T \mathbf{X} + b_{op}$$

- $g(x)$  می‌تواند مثبت یا منفی باشد.



# مرز جداسازی (ادامه...)

- در صورتی که  $X$  بردار پشتیبان باشد، طبق شکل زیر خواهیم داشت:

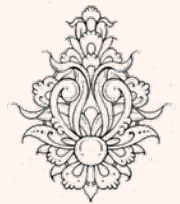


$$\mathbf{X} = \mathbf{X}_p + \overrightarrow{AB}$$

- $AB$  در جهت عمود بر مرز جداکننده
- اگر اندازهی بردار  $AB=r$  در نظر گرفته شود، خواهیم داشت:

$$\mathbf{X} = \mathbf{X}_p + r \frac{\mathbf{W}_{op}}{\|\mathbf{W}_{op}\|}$$

$$\overrightarrow{AB} = r \frac{\mathbf{W}_{op}}{\|\mathbf{W}_{op}\|}$$



# مرز جداسازی (ادامه...)

$$g(\mathbf{X}) = \mathbf{W}_{op}^T \mathbf{X} + b_{op}$$

• داشته‌یم:

$$\mathbf{X} = \mathbf{X}_p + r \frac{\mathbf{W}_{op}}{\|\mathbf{W}_{op}\|}$$

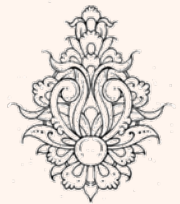
$$g(\mathbf{X}) = \mathbf{W}_{op}^T \left[ \mathbf{X}_p + r \frac{\mathbf{W}_{op}}{\|\mathbf{W}_{op}\|} \right] + b_{op}$$

$$g(\mathbf{X}) = \underbrace{\mathbf{W}_{op}^T \mathbf{X}_p + b_{op}}_{\text{روی مرز پس برابر با صفر}} + r \frac{\mathbf{W}_{op}}{\|\mathbf{W}_{op}\|} \mathbf{W}_{op}^T$$

روی مرز پس برابر با صفر

$$g(\mathbf{X}) = r \frac{\|\mathbf{W}_{op}\|^2}{\|\mathbf{W}_{op}\|}$$

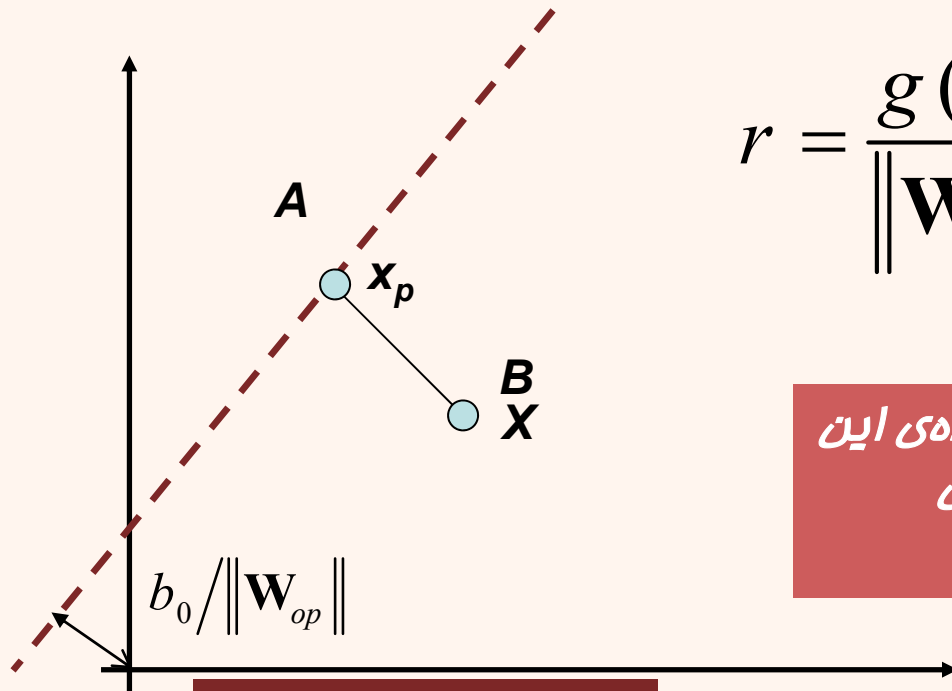
$$g(\mathbf{X}) = r \|\mathbf{W}_{op}\|$$



# مرز جداسازی (ادامه...)

$$g(\mathbf{X}) = r \|\mathbf{W}_{op}\|$$

$$r = \frac{g(\mathbf{X})}{\|\mathbf{W}_{op}\|}$$



فاصله از مبدا مقلصات

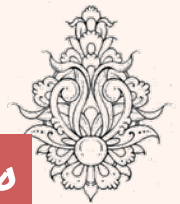
هدف ما کمینه نمودن  $r$  است.

در این حالت سمت شرایطی می‌باید  $W$  کمینه گردد.

$\mathbf{X}=0$

$$r = \frac{g(\mathbf{X})}{\|\mathbf{W}_{op}\|} = \frac{b_{op}}{\|\mathbf{W}_{op}\|}$$

مثبت یا منفی بودن  $b_{op}$  نشان دهنده این است که مبدأ در کدام سمت خط مرزی هستیم.



# مرز جداسازی (ادامه...)

مکان یافتن  $W_{op}$  و  $b_{op}$  است

• جداساز خطی را به صورت زیر در نظر می‌گیریم:

$$(X_i, +1) \quad W_{op}^T X_i + b_{op} \geq 1 \quad \text{for } d_i = +1$$

$$(X_i, -1) \quad W_{op}^T X_i + b_{op} \leq -1 \quad \text{for } d_i = -1$$

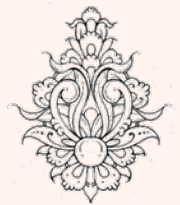
به صورت کلی داریم:

$$d_i (W_{op}^T X + b_{op}) \geq 1$$

• رابطه‌ی بالا برای تمامی الگوهای آموزشی برقرار است.

• و در نتیجه برای بردارهای پشتیبان

$$g(X_s) = W_{op}^T X_s + b_{op} = \pm 1$$



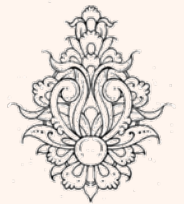
# مرز جداسازی (ادامه...)

$$g(\mathbf{X}_s) = \mathbf{W}_{op}^T \mathbf{X}_s + b_{op} = \pm 1$$

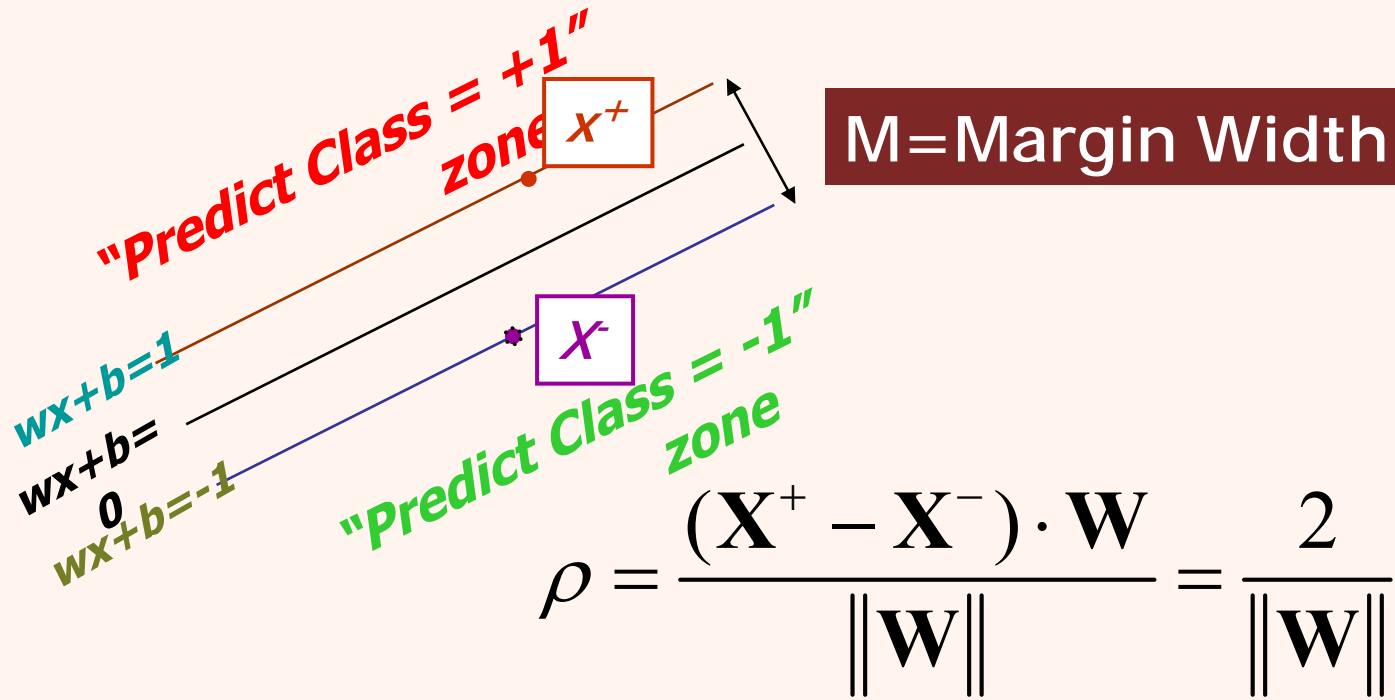
$$r = \frac{g(\mathbf{X}_s)}{\|\mathbf{W}_{op}\|} = \begin{cases} \frac{1}{\|\mathbf{W}_{op}\|} \\ -\frac{1}{\|\mathbf{W}_{op}\|} \end{cases}$$

- در نتیجه فاصله‌ی دو بردار پشتیبان در دو طرف مرز:

$$\rho = 2r = \frac{2}{\|\mathbf{W}_{op}\|}$$

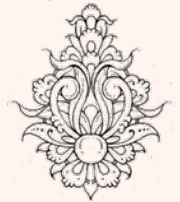


# جدایی پذیر خطی



می دانیم:

- $W \cdot X^+ + b = +1$
- $W \cdot X^- + b = -1$
- $W \cdot (X^+ - X^-) = 2$





# جدایی پذیر خطی

$$\rho = 2r = \frac{\|\pm 2\|}{\|\mathbf{W}_{op}\|}$$

- با توجه به دو رابطه‌ی
- به این نتیجه می‌رسیم که  $\mathbf{W}_{op}$  می‌باید مینیمم گردد.

• این مسأله معادل مینیمم کردن

$$\Phi(\mathbf{W}) = \frac{1}{2} \mathbf{W}^T \mathbf{W}$$

–  $\Phi$  یک تابع محدب (Convex Function) است.

– طبق رابطه‌ی  $d_i(\mathbf{W}_{op}^T \mathbf{X} + b_{op}) \geq 1$  برای  $N$  الگوی آموزشی شرط زیر می‌باید برقرار باشد:

$$\sum_{i=1}^N [d_i(\mathbf{W}_{op}^T \mathbf{X} + b_{op}) - 1]$$

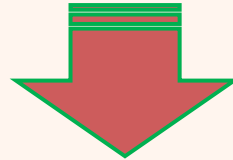
این میزان بزرگتر یا مساوی صفر است



وزن‌ها و بایاس را به گونه‌ای بیابید که:

$$\rho = \frac{2}{\|\mathbf{W}\|} \text{ is maximized}$$

and for all  $(\mathbf{X}_i, y_i), i=1..n$  :  $y_i(\mathbf{W}^T \mathbf{X}_i + b) \geq 1$



وزن‌ها و بایاس را به گونه‌ای بیابید که:

$$\Phi(\mathbf{W}) = 1/2 \|\mathbf{W}\|^2 = 1/2 \mathbf{W}^T \mathbf{W} \text{ is minimized}$$

and for all  $(\mathbf{X}_i, y_i), i=1..n$  :  $y_i (\mathbf{W}^T \mathbf{X}_i + b) \geq 1$



# جدایی پذیر خطی

- رابطه‌ی لاگرانژ زیر تعریف می‌شود به گونه‌ای که هر دو قید ذکر شده را پوشش دهد:

$$j(\mathbf{W}, b, \alpha) = \frac{1}{2} \mathbf{W}^T \mathbf{W} - \sum_{i=1}^N \alpha [d_i (\mathbf{W}^T \mathbf{X}_i + b) - 1]$$

Lagrange multiplier(nonnegative)

- برای به دست آوردن  $\mathbf{W}_{op}$  و  $b_{op}$  نسبت به هر دو مشتق می‌گیریم.

$$\frac{\partial j}{\partial \mathbf{W}} = 0 \Rightarrow \mathbf{W} - \sum_{i=1}^N \alpha_i d_i \mathbf{X}_i = 0 \Rightarrow$$

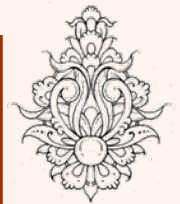
$$\mathbf{W}_{op} = \sum_{i=1}^N \alpha_i d_i \mathbf{X}_i$$

$$\frac{\partial j}{\partial b} = 0 \Rightarrow -\sum_{i=1}^N \alpha_i d_i = 0 \Rightarrow$$

$$\sum_{i=1}^N \alpha_i d_i = 0$$

$b_{op}$  به دست  
نمی‌آید

باید قید می‌دهد



# جدایی پذیر خطی

- دو شرط برای  $\alpha$  خواهیم داشت:

$$\left\{ \begin{array}{l} \sum_{i=1}^N \alpha_i d_i = 0 \\ \alpha_i [d_i (\mathbf{W}^T \mathbf{X}_i + b) - 1] = 0 \end{array} \right.$$

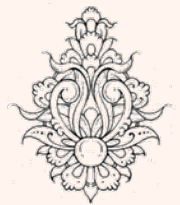
صفر

غیر صفر

Kuhn-Tucker condition of optimization theory

- به ازای هر  $\alpha_i$  برای الگوهای آموزشی متناظر با SVها رابطه‌ی زیر برقرار است.

$$d_i (\mathbf{W}_{op}^T \mathbf{X}_i + b_{op}) - 1 = 0$$



# یافتن رویه بهینه

$$j(\mathbf{W}, b, \alpha) = \frac{1}{2} \|\mathbf{W}\|^2 - \sum_{i=1}^N \alpha_i [d_i (\mathbf{W}^T \mathbf{X}_i + b) - 1]$$

$$j(\mathbf{W}, b, \alpha) = \frac{1}{2} \mathbf{W}^T \mathbf{W} - \sum_{i=1}^N \alpha_i d_i \mathbf{W}^T \mathbf{X}_i - \sum_{i=1}^N \alpha_i d_i b + \sum_{i=1}^N \alpha_i$$

• برای مقادیر بهینه داشتیم:

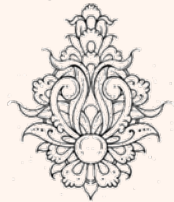
$$\mathbf{W}_{op} = \sum_{i=1}^N \alpha_i d_i \mathbf{X}_i$$

$$\sum_{i=1}^N \alpha_i d_i = 0$$

Duality theorem

• پس خواهیم داشت:

$$j(\mathbf{w}_{op}, b_{op}, \alpha) = \frac{1}{2} \mathbf{w}_{op}^T \mathbf{w}_{op} - \mathbf{w}_{op}^T \mathbf{w}_{op} + 0 + \sum_{i=1}^N \alpha_i$$

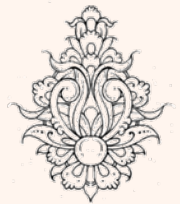


# Dual Problem

$$j(\mathbf{W}_{op}, b_{op}, \alpha) = \frac{1}{2} \mathbf{W}_{op}^T \mathbf{W}_{op} - \mathbf{W}_{op}^T \mathbf{W}_{op} + 0 + \sum_{i=1}^N \alpha_i$$

$$\left\{ \begin{array}{l} = \sum_{i=1}^N \alpha_i - \frac{1}{2} \mathbf{W}_{op}^T \mathbf{W}_{op} = Q(\alpha) \\ \sum_{i=1}^N \alpha_i d_i = 0 \end{array} \right.$$

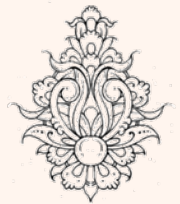
مینیمم نمودن  $W$  همانند ماکزیمم نمودن  $Q$  است  
زیرا در  $W_{op}$  مقدار  $W$  کمترین میزان است و در این  
صورت است که کل عبارت ماکزیمم می‌شود.



# یافتن رویه‌ی بهینه

$$\begin{aligned} Q(\alpha) &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \mathbf{w}_{op}^T \mathbf{w}_{op} \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \left[ \sum_{i=1}^N \alpha_i d_i \mathbf{X}_i \right]^T \left[ \sum_{j=1}^N \alpha_j d_j \mathbf{X}_j \right] \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i d_i \alpha_j d_j \mathbf{X}_i^T \mathbf{X}_j \end{aligned}$$

$$\left\{ \begin{aligned} &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i d_i \alpha_j d_j \mathbf{X}_i^T \mathbf{X}_j \\ &\sum_{i=1}^N \alpha_i d_i = 0 \\ &\alpha_i \geq 0 \text{ for } i=0,1,\dots,N \end{aligned} \right.$$



$\alpha_i$  ها وابسته به الگوهای  $\mathbf{X}_i$ ،  $\mathbf{X}_j$  و خروجی‌های مرتبط است

# یافتن رویه بهینه

$$= \sum_{i=1}^N \alpha_i - \frac{1}{2} \left[ \sum_{i=1}^N \alpha_i d_i \mathbf{X}_i \right]^T \left[ \sum_{j=1}^N \alpha_j d_j \mathbf{X}_j \right]$$

جهت محاسبه  $\alpha$  نسبت به  $\alpha_k$  مشتق گرفته برابر با صفر قرار می‌دهیم

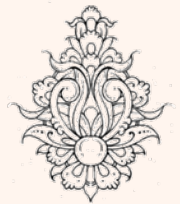
$$\frac{\partial Q(\alpha)}{\partial \alpha_k} = 1 - \sum_{\substack{i=1 \\ i \neq k}}^N \alpha_i d_i d_k \mathbf{X}_i^T \mathbf{X}_k - \alpha_k d_k^2 \mathbf{X}_k^T \mathbf{X}_k = 0$$

$N$  معادله و  $N$  مجهول

$$M_{i,j} = \mathbf{X}_i^T \mathbf{X}_j$$

ضرب داخلی

$$\frac{\partial Q(\alpha)}{\partial \alpha_k} = 1 - d_k \sum_{i=1}^N \alpha_i d_i M_{i,k} = 0$$







# یافتن رویه بهینه

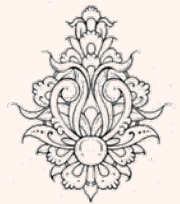
$$\frac{\partial Q(\alpha)}{\partial \alpha_k} = 1 - d_k \sum_{i=1}^N \alpha_i d_i M_{i,k} = 0$$

• پس از به دست آوردن  $\alpha$  خواهیم داشت:

$$w_{op} = \sum_{i=1}^N \alpha_i d_i x_i$$

$\mathbf{x} = \mathbf{x}^{\text{support vector}}$    $w_{op}^T \mathbf{x}^s + b_{op} = \pm 1$

  $b_{op} = \pm 1 - w_{op}^T \mathbf{x}^s$



# یافتن رویه‌ی بهینه

$$W = \sum \alpha_i y_i X_i \quad b = y_k - W^T X_k \text{ for any } X_k \text{ such that } \alpha_k \neq 0$$

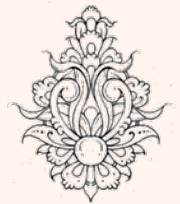
هر  $\alpha_i$  مخالف صفر، نشان‌دهنده‌ی این است که  $X_i$  متناظرش یک بردار پشتیبان است.  
در این حالت تابع جداکننده همانند زیر است:

$$f(X) = \sum \alpha_i d_i \underbrace{X_i^T X}_{} + b$$

ضرب داخلی دو بردار

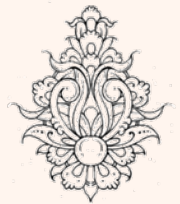
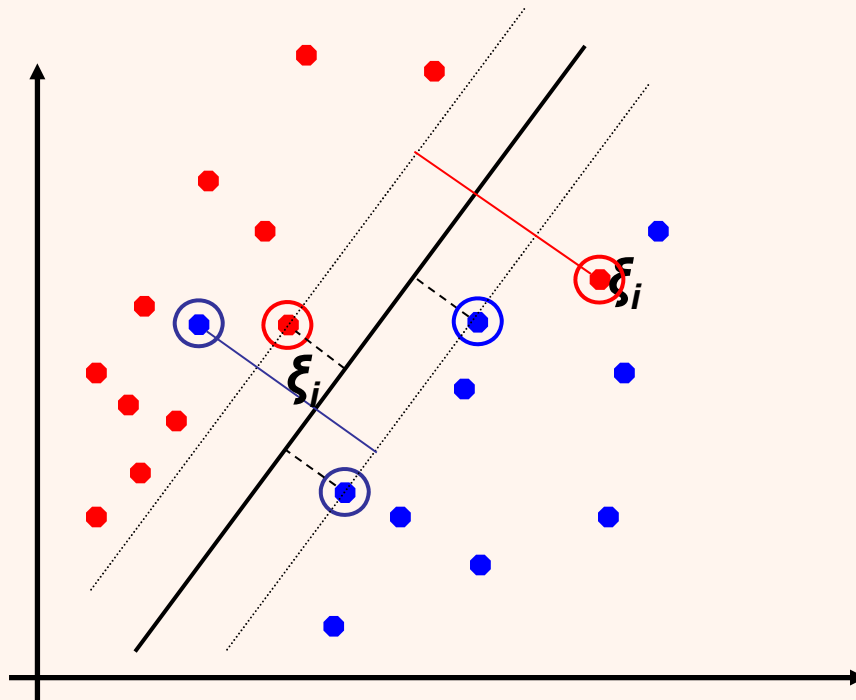
توجه:

حل مساله بهینه‌سازی وابسته به محاسبه ضرب داخلی بین تمامی نمونه‌های آموزشی است.



# Soft Margin

- SVM برای داده‌های جدایی‌پذیر خطی مورد بررسی قرار گرفت.
- حال اگر مجموعه‌ی داده‌های آموزش قابلیت جداسازی را نداشته باشند، چه خواهد شد؟ به بیان بهتر صحبت در مورد مسائل جدایی‌پذیر است که با نویز همراه هستند.



# Soft Margin

- مسأله‌ی **Hard Margin** را تبدیل به حل مسأله‌ی **Soft Margin** می‌شود.

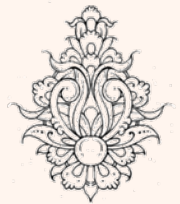
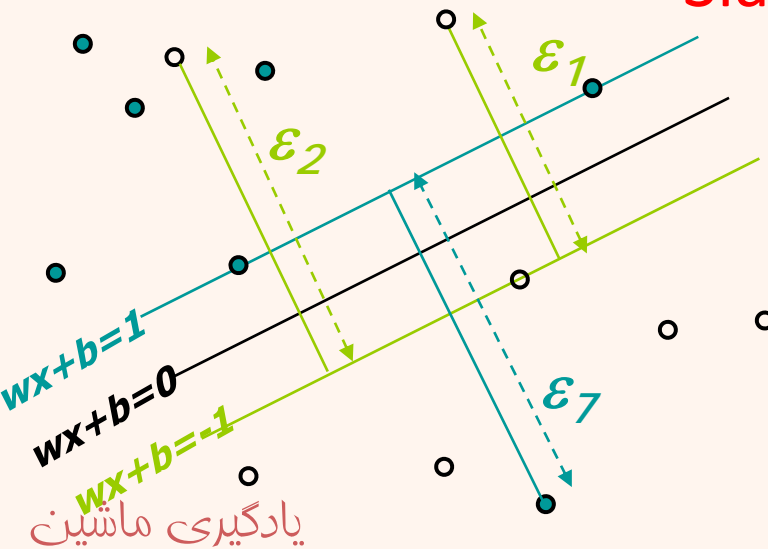
- ماشینی جداسازی soft گفته می‌شود، در صورتی که برای برخی داده‌ها شرط زیر نقض شود:

$$d_i(\mathbf{W}_{op}^T \mathbf{X} + b_{op}) \geq 1$$

- با اضافه کردن یک **Slack Variable**

- مسأله را بار دیگر بررسی می‌کنیم.

- این متغیر میزان انحراف از شرط فوق را نشان می‌دهد.



# XOR Problem مثال

$$N = 4$$

$$X_1 = [-1 \ -1] \rightarrow d_1 = -1$$

$$X_2 = [-1 \ 1] \rightarrow d_2 = +1$$

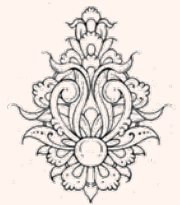
$$X_3 = [1 \ -1] \rightarrow d_3 = +1$$

$$X_4 = [1 \ 1] \rightarrow d_4 = -1$$

$$K(\mathbf{X}, \mathbf{X}_i) = \Phi^T(\mathbf{X}) \cdot \Phi(\mathbf{X}_i)$$

$$K(\mathbf{X}, \mathbf{X}_i) = (1 + \mathbf{X}^T \mathbf{X}_i)^2$$

• نمونه‌های آموزشی دو بعدی هستند.



# XOR Problem

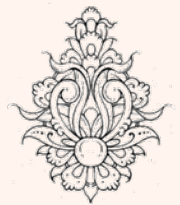
$$\mathbf{X}_i = [x_{i1} \ x_{i2}]$$

$$\mathbf{X} = [x_1 \ x_2]$$

$$K(\mathbf{X}, \mathbf{X}_i) = (1 + \mathbf{X}^T \mathbf{X}_i)^2$$

$$= (1 + [x_{i1} \ x_{i2}] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix})^2 = (1 + x_{i1}x_1 + x_{i2}x_2)^2$$
$$= 1 + x_1^2 x_{i1}^2 + 2x_1 x_2 x_{i1} x_{i2} + x_{i2}^2 x_2^2 + 2x_1 x_{i1} + 2x_2 x_{i2}$$

- حال اگر بخواهیم پاسخ به دست آمده را با ضرب داخلی دو بردار  $\phi(\mathbf{X})$  و  $\phi(\mathbf{X}_i)$  نشان دهیم خواهیم داشت:

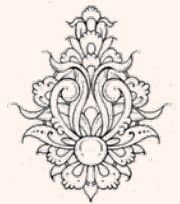


# XOR Problem

$$= 1 + x_1^2 x_{i1}^2 + 2x_1 x_2 x_{i1} x_{i2} + x_{i2}^2 x_2^2 + 2x_1 x_{i1} + 2x_2 x_{i2}$$

$$\boldsymbol{\varphi}(\mathbf{x}) = [1, x_1^2, \sqrt{2}x_1 x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2]^T$$

$$\boldsymbol{\varphi}(x_i) = [1 + x_{i1}^2, \sqrt{2}x_{i1} x_{i2}, x_{i2}^2, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}]^T \quad i=1,2,3,4$$

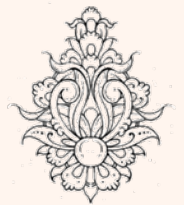


# XOR Problem

$$\boldsymbol{\varphi}(\mathbf{x}) = [1, x_1^2, \sqrt{2}x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2]^T$$

$$\boldsymbol{\varphi}(x_i) = [1 + x_{i1}^2, \sqrt{2}x_{i1}x_{i2}, x_{i2}^2, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}]^T \quad i=1,2,3,4$$

$$\begin{array}{l} X_1 = [-1 \ -1] \\ X_2 = [-1 \ 1] \\ X_3 = [1 \ -1] \\ X_4 = [1 \ 1] \end{array} \left[ \begin{array}{cccccc} 1 & 1 & \sqrt{2} & 1 & -\sqrt{2} & -\sqrt{2} \\ 1 & 1 & -\sqrt{2} & 1 & -\sqrt{2} & \sqrt{2} \\ 1 & 1 & -\sqrt{2} & 1 & \sqrt{2} & -\sqrt{2} \\ 1 & 1 & \sqrt{2} & 1 & \sqrt{2} & \sqrt{2} \end{array} \right]$$

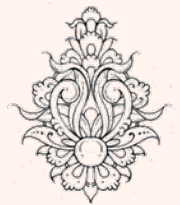




# XOR Problem

$$K_{4 \times 4} = \begin{bmatrix} 1 & 1 & \sqrt{2} & 1 & -\sqrt{2} & -\sqrt{2} \\ 1 & 1 & -\sqrt{2} & 1 & -\sqrt{2} & \sqrt{2} \\ 1 & 1 & -\sqrt{2} & 1 & \sqrt{2} & -\sqrt{2} \\ 1 & 1 & \sqrt{2} & 1 & \sqrt{2} & \sqrt{2} \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ \sqrt{2} & -\sqrt{2} & -\sqrt{2} & \sqrt{2} \\ 1 & 1 & 1 & 1 \\ -\sqrt{2} & -\sqrt{2} & \sqrt{2} & \sqrt{2} \\ -\sqrt{2} & \sqrt{2} & -\sqrt{2} & \sqrt{2} \end{bmatrix}$$

$$K_{4 \times 4} = \begin{bmatrix} 9 & 1 & 1 & 1 \\ 1 & 9 & 1 & 1 \\ 1 & 1 & 9 & 1 \\ 1 & 1 & 1 & 9 \end{bmatrix}$$



# XOR Problem

$$N = 4$$

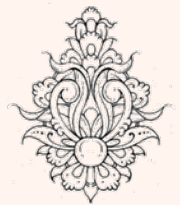
$$Q(\alpha) = \sum_{i=1}^4 \alpha_i - \frac{1}{2} \sum_{i=1}^4 \sum_{j=1}^4 \alpha_i \alpha_j d_i d_j K(\mathbf{X}_i, \mathbf{X}_j)$$

$$Q(\alpha) = \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4$$

$$- \frac{1}{2} (9\alpha_1^2 + 9\alpha_2^2 + 9\alpha_3^2 + 9\alpha_4^2 - 2\alpha_1\alpha_2 - 2\alpha_1\alpha_3 + 2\alpha_1\alpha_4 + 2\alpha_2\alpha_3 - 2\alpha_2\alpha_4 - 2\alpha_3\alpha_4)$$

- برای به دست آوردن  $\alpha$  ها مشتق گرفته برابر با صفر قرار می‌دهیم:

$$\frac{\partial Q(\alpha)}{\partial \alpha_1} = 0 \Rightarrow 1 - 9\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 = 0$$



# XOR Problem

$$\frac{\partial Q(\alpha)}{\partial \alpha_1} = 0 \Rightarrow 1 - 9\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 = 0$$

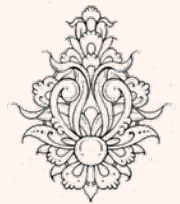
$$\frac{\partial Q(\alpha)}{\partial \alpha_2} = 0 \Rightarrow 1 + \alpha_1 - 9\alpha_2 - \alpha_3 + \alpha_4 = 0$$

$$\frac{\partial Q(\alpha)}{\partial \alpha_3} = 0 \Rightarrow 1 + \alpha_1 - \alpha_2 - 9\alpha_3 - \alpha_4 = 0$$

$$\frac{\partial Q(\alpha)}{\partial \alpha_4} = 0 \Rightarrow 1 - \alpha_1 + \alpha_2 + \alpha_3 - 9\alpha_4 = 0$$

$$\alpha_i = \frac{1}{8}$$

$$Q(\alpha) = \frac{1}{4}$$



• پس از محاسبه  $\alpha$  ها  $W_{opt}$  را محاسبه می‌کنیم:



# XOR Problem

• جهت محاسبه‌ی اندازه‌ی وزن بهینه داریم:

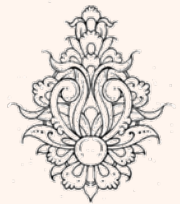
$$\frac{1}{2} \|\mathbf{w}_{opt}\|^2 = \frac{1}{4} \quad \Rightarrow \quad \|\mathbf{w}_{opt}\| = \frac{1}{\sqrt{2}}$$

• داشتیم:

$$\mathbf{w}_{opt} = \sum_{i=1}^N \alpha_i \cdot d_i(\varphi_j(\mathbf{x}_i))$$

$$\mathbf{w}_o = \frac{1}{8} [-\varphi(\mathbf{x}_1) + \varphi(\mathbf{x}_2) + \varphi(\mathbf{x}_3) - \varphi(\mathbf{x}_4)]$$

$$= \frac{1}{8} \left[ - \begin{bmatrix} 1 \\ 1 \\ \sqrt{2} \\ 1 \\ -\sqrt{2} \\ -\sqrt{2} \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ -\sqrt{2} \\ 1 \\ -\sqrt{2} \\ \sqrt{2} \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ -\sqrt{2} \\ 1 \\ \sqrt{2} \\ -\sqrt{2} \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \\ \sqrt{2} \end{bmatrix} \right] = \begin{bmatrix} 0 \\ 0 \\ -1/\sqrt{2} \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

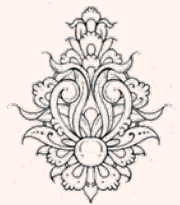


# XOR Problem

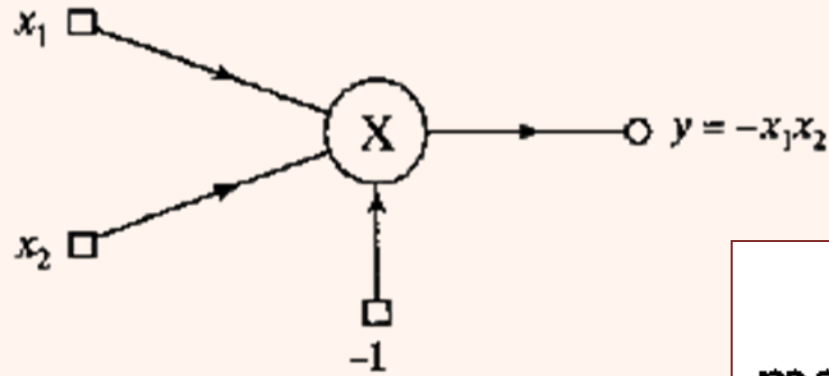
- رویه‌ی بهینه به وسیله‌ی رابطی زیر محاسبه می‌شود:

$$\mathbf{W}_{opt}^T \varphi(\mathbf{X}) = 0$$

$$\begin{bmatrix} 0 & 0 & \frac{-1}{\sqrt{2}} & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \end{bmatrix} = 0 \quad \Rightarrow \quad -x_1x_2 = 0$$

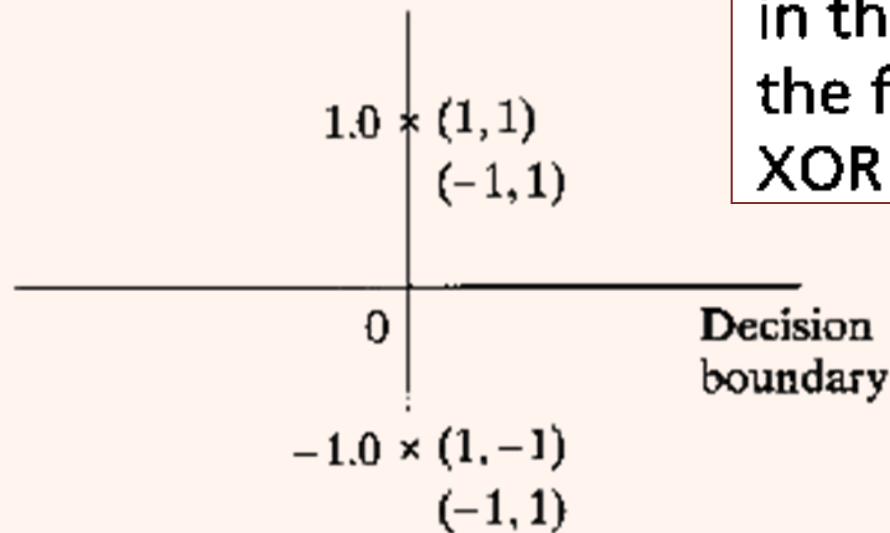


# XOR Problem

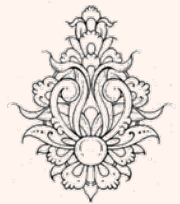


(a)

(a) Polynomial machine for solving the XOR problem. (b) Induced images in the feature space due to the four data points of the XOR problem.



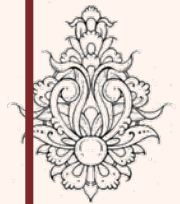
(b)



# XOR Problem

```
X=[-1 -1 1 1;  
    -1 1 -1 1];  
d=[-1 1 1 -1];  
n=4;  
K=zeros(n,n);  
for i=1:n  
    for j=1:n  
        xi1=X(1,i);  
        xi2=X(2,i);  
        xj1=X(1,j);  
        xj2=X(2,j);  
        fi=[1 xi1^2 xi2^2  
sqrt(2)*xi1*xi2 sqrt(2)*xi1  
sqrt(2)*xi2 ];  
        fj=[1 xj1^2 xj2^2  
sqrt(2)*xj1*xj2 sqrt(2)*xj1  
sqrt(2)*xj2 ];  
        K(i,j)=fi*fj';%Kernel  
    end  
end
```

```
end  
KD=zeros(n,n);  
for i=1:n  
    for j=1:n  
  
KD(i,j)=K(i,j)*d(i)*d(j);  
    end  
end  
  
alfa=inv(KD)*ones(n,1);
```



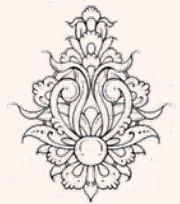
• می‌توان گفت برای به دست آوردن  $\alpha$

$$Q(\alpha) = \sum_{i=1}^{m_1} \alpha_i - \frac{1}{2} \sum_{i=1}^{m_1} \sum_{j=1}^{m_1} \alpha_i \alpha_j d_i d_j K(x_i, x_j)$$

$$\frac{\partial Q(\alpha)}{\partial \alpha_j} = 1 - \sum_{i=1}^{m_1} \alpha_j d_i d_j K(x_i, x_j) = 0 \quad j=1, 2, \dots, m_1$$

$$K'(i, j)$$

$$\alpha = K^{-1} \cdot \begin{bmatrix} 1 \\ \cdot \\ \cdot \\ 1 \end{bmatrix}$$



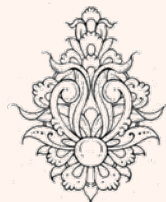


```

f=zeros(4,6);
for i=1:n
    xi1=X(1,i);
    xi2=X(2,i);
    f(i,:)=alfa(i)*d(i)*[1 xi1^2 xi2^2
sqrt(2)*xi1*xi2 sqrt(2)*xi1 sqrt(2)*xi2 ];
end
W=sum(f);

O=zeros(1,4);
for i=1:n
    xi1=X(1,i);
    xi2=X(2,i);
    O(i)=W*[1 xi1^2 xi2^2 sqrt(2)*xi1*xi2
sqrt(2)*xi1 sqrt(2)*xi2 ]';
end

```

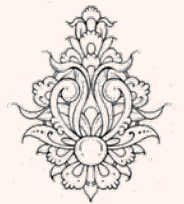


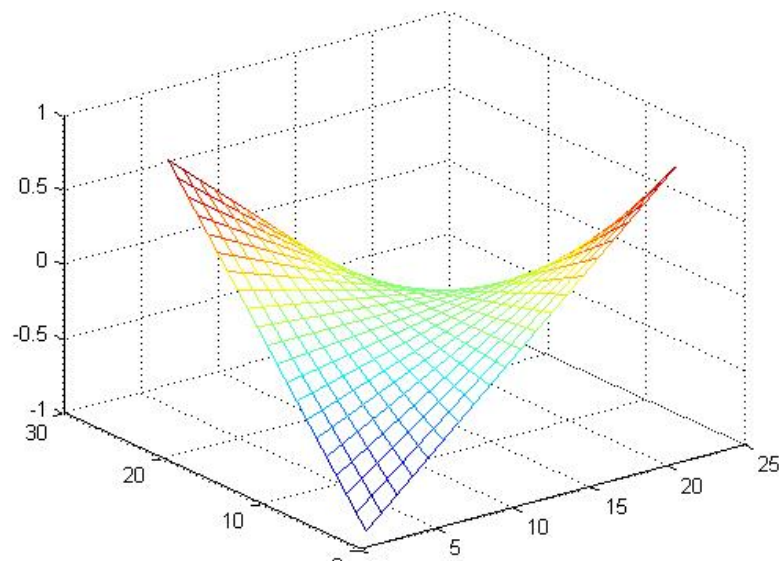
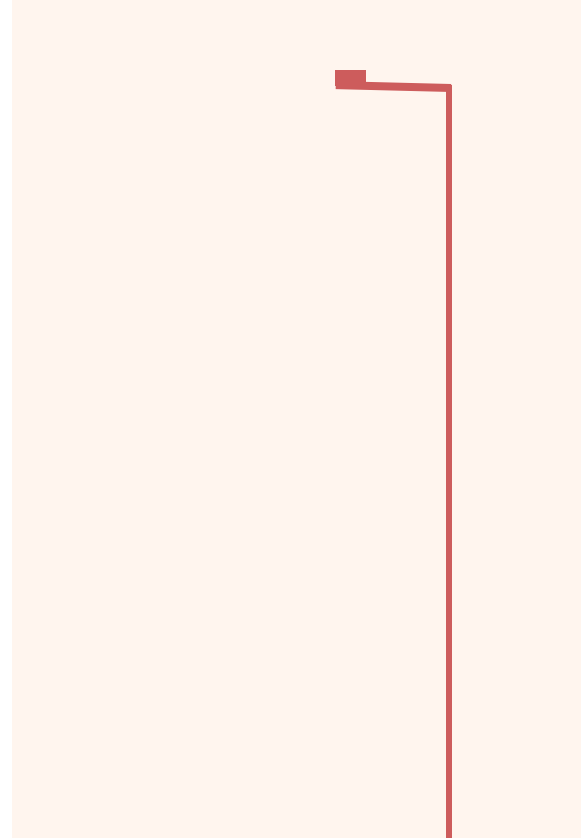
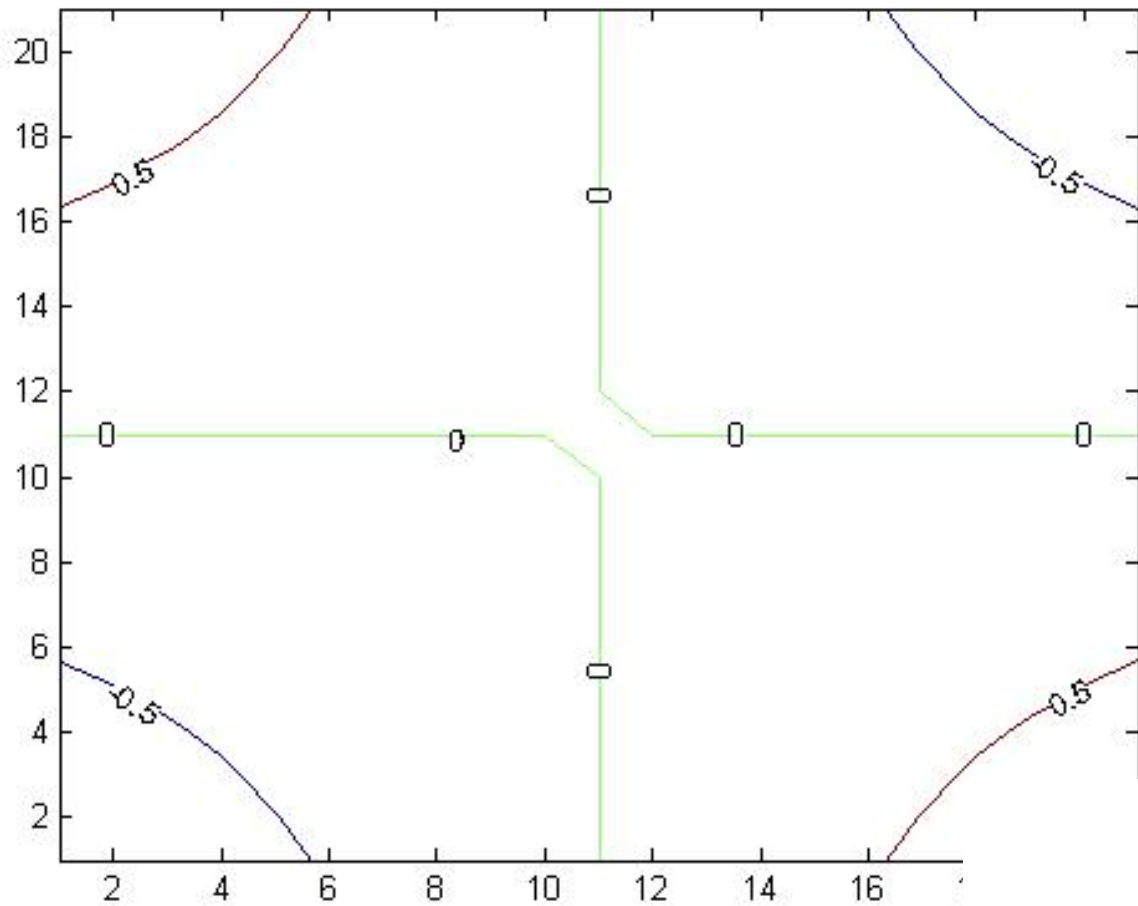
```

x1=-1:0.1:1;
x2=-1:0.1:1;
for i=1:21
    for j=1:21
        xi1=x1(i);
        xi2=x2(j);
        M(i,j)=W*[1 xi1^2 xi2^2
sqrt(2)*xi1*xi2 sqrt(2)*xi1 sqrt(2)*xi2
]';
    end
end

mesh(M);
figure;
[c,h]=contour(M,[-1 -0.5 0 0.5 1]);
clabel(c,h);

```

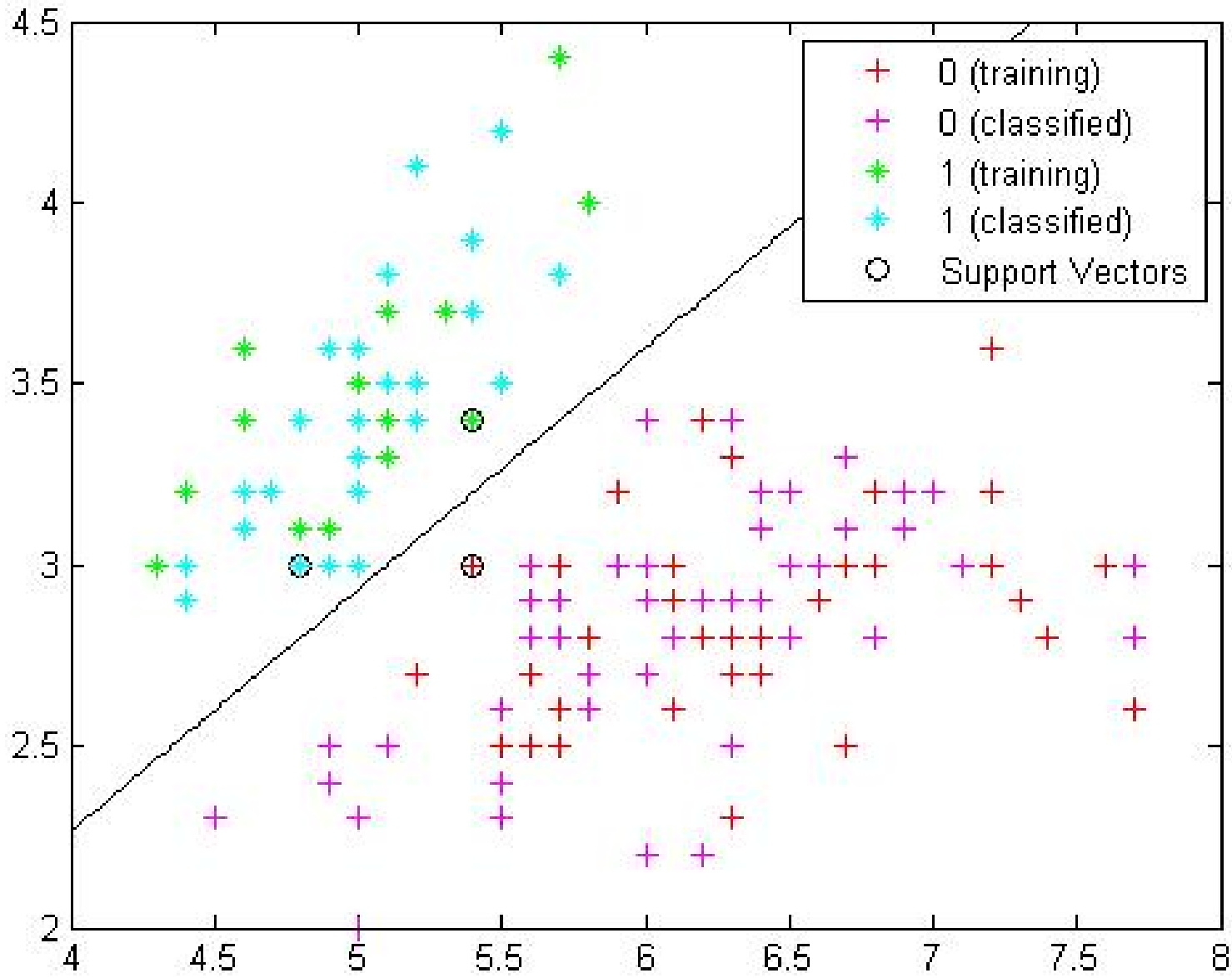




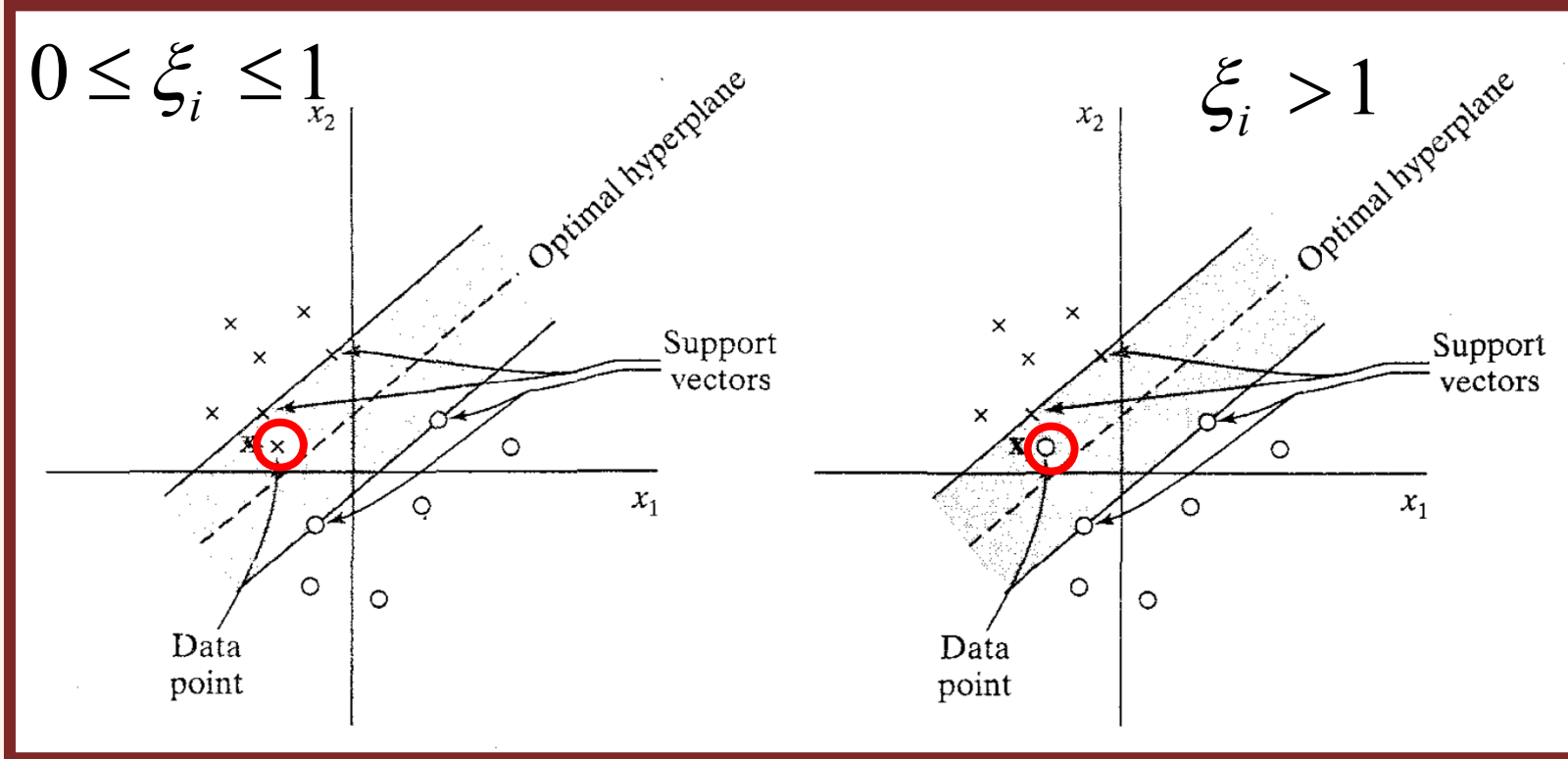
```
load fisheriris% load dataset
data = [meas(:,1), meas(:,2)];%Create data, a two-column matrix
groups = ismember(species,'setosa');%divide data into two groups: Setosa
and non-Setosa.
[train, test] = crossvalind('holdOut',groups);%test and train randomly
cp = classperf(groups);%Evaluate performance of classifier
svmStruct = svmtrain(data(train,:),groups(train),'showplot',true);%train
an SVM classifier using a linear kernel function and plot the grouped
data
title(sprintf('Kernel Function: %s',...
              func2str(svmStruct.KernelFunction)),...
       'interpreter','none');
classes = svmclassify(svmStruct,data(test,:), 'showplot',true);
classperf(cp,classes,test);
p.CorrectRate

figure
svmStruct = svmtrain(data(train,:),groups(train),...
                    'showplot',true,'boxconstraint',1e6);
classes = svmclassify(svmStruct,data(test,:), 'showplot',true);
classperf(cp,classes,test);
cp.CorrectRate
```



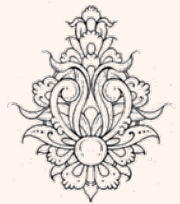


# Soft Margin Classification



- دو حالت ممکن است رخ دهد:
- داده‌ی در کلاس درست ولی در ماشیه قرار گیرد.
- داده‌ی آموزشی به اشتباه طبقه‌بندی شود.

$$d_i(\mathbf{W}_{op}^T \mathbf{X} + b_{op}) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N$$

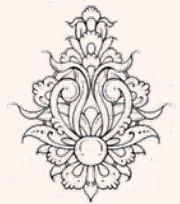


# Soft Margin Classification

$$d_i(\mathbf{W}_{op}^T \mathbf{X} + b_{op}) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N$$

- در این حالت برداهای پشتیبان آن‌هایی هستند که در رابطه‌ی تساوی در عبارت بالا صدق می‌کنند، حتی با وجود  $\xi > 0$
- در صورتی که داده‌های نویزی از مجموعه خارج شود، رویه‌ی جداکننده تخییر خواهد کرد.
- هدف یافتن «رویه‌ای جداکننده» است که در آن خطای طبقه‌بندی نادرست در آن مینیمم شود:

$$\Phi(\xi) = \sum_{i=1}^N I(\xi_i - 1) \quad I(\xi) = \begin{cases} 0 & \text{if } \xi \leq 0 \\ 1 & \text{if } \xi > 0 \end{cases}$$





# Soft Margin Classification

- با توجه به این که کمینه کردن چنین تابعی یک مسأله‌ی بهینه‌سازی **nonconvex** است و در رده‌ی NP-complete قرار می‌گیرد، آن را با تابع زیر تقریب می‌زنیم:

$$\Phi(\xi) = \sum_{i=1}^N \xi_i$$

- و در کل هدف می‌نیمیم کردن عبارت زیر است:

$$\Phi(\mathbf{W}, \xi) = \frac{1}{2} \mathbf{W}^T \mathbf{W} + C \sum_{k=1}^R \xi_k$$

regularization parameter

این پارامتر نوعی مصالحه بین پیچیدگی ماشین و خطا برقرار می‌کند. هرچه  $C$  به صفر نزدیک‌تر باشد به این معناست که نمونه‌هایی در حاشیه قرار می‌گیرند، اهمیت کمتری دارند و در نتیجه حاشیه بزرگ‌تر می‌شود. و هرچه بزرگ‌تر باشد، ما به حالت **hard margin** نزدیک‌تر می‌شود.



# Soft Margin

- برای Hard Margin داشتیم:

Find  $\mathbf{W}$  and  $b$  such that

$\Phi(\mathbf{W}) = \frac{1}{2} \mathbf{W}^T \mathbf{W}$  is minimized and for all  $\{(X_i, y_i)\}$   
 $y_i (\mathbf{W}^T X_i + b) \geq 1$

- با اضافه کردن Slack Variable داریم:

Find  $\mathbf{W}$  and  $b$  such that

$\Phi(\mathbf{W}) = \frac{1}{2} \mathbf{W}^T \mathbf{W} + C \sum \xi_i$  is minimized and for all  $\{(X_i, y_i)\}$   
 $y_i (\mathbf{W}^T X_i + b) \geq 1 - \xi_i$  and  $\xi_i \geq 0$  for all  $i$



# Soft Margin Classification

هدف یافتن ضرایب لاگراشر بیشینه در عبارت زیر است:

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \mathbf{X}_i^T \mathbf{X}_j$$

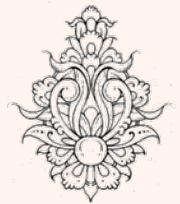
با در نظر گرفتن محدود زیر

$$\sum_{i=1}^N \alpha_i d_i = 0$$

$$0 \leq \alpha \leq C$$

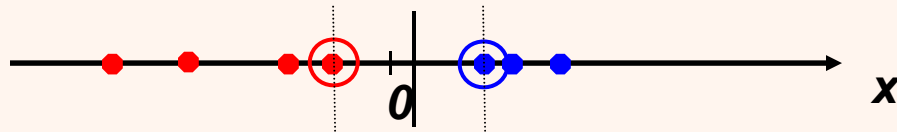
بقیه مراحل مانند حالت قبل خواهد بود:

$$W_{op} = \sum_{i=1}^N \alpha_i d_i \mathbf{X}_i$$

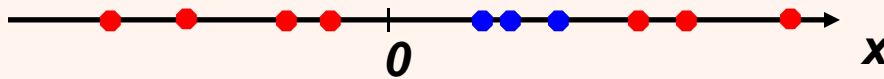


# SVM غیرخطی

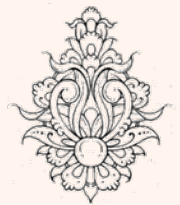
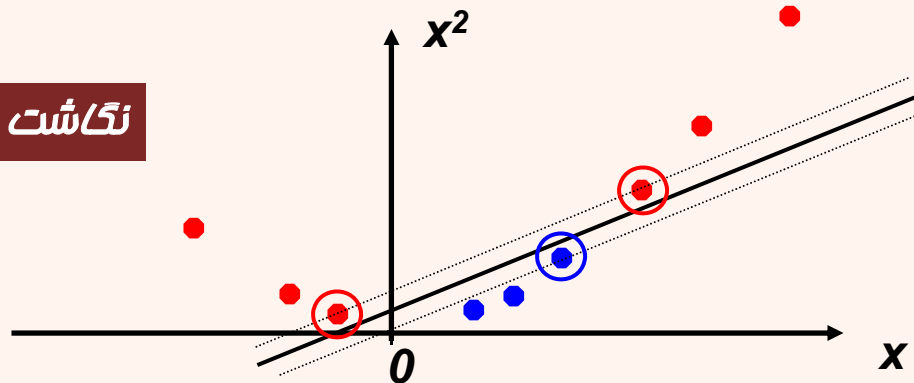
- برای داده‌هایی که قابلیت جداسازی خطی دارند، عملکرد سیستم بسیار خوب است.



- اگر داده‌ها به صورت‌های زیر باشد، مسأله چگونه حل می‌شود؟

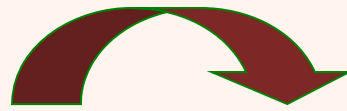


نگاشت به یک فضای *High Dimension*

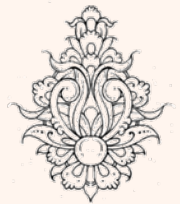
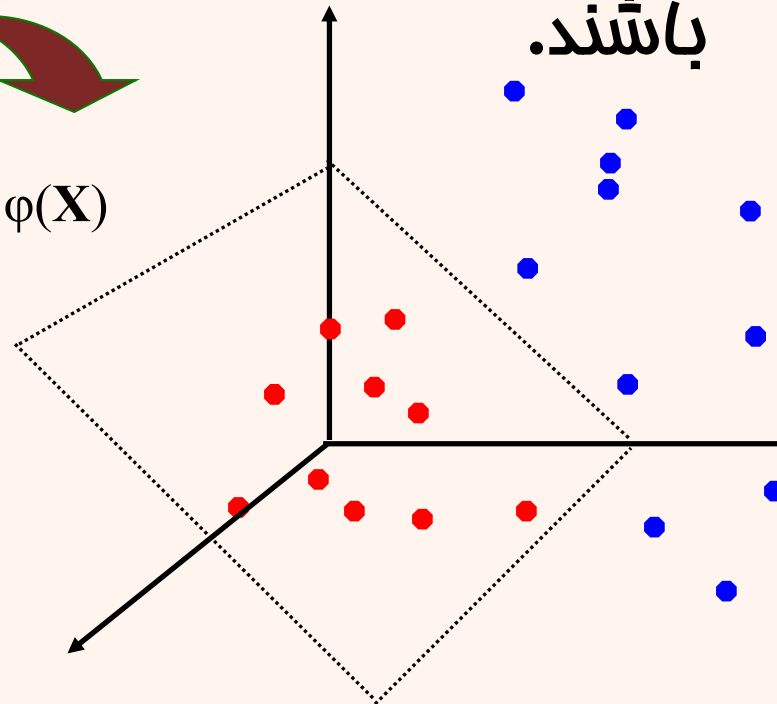
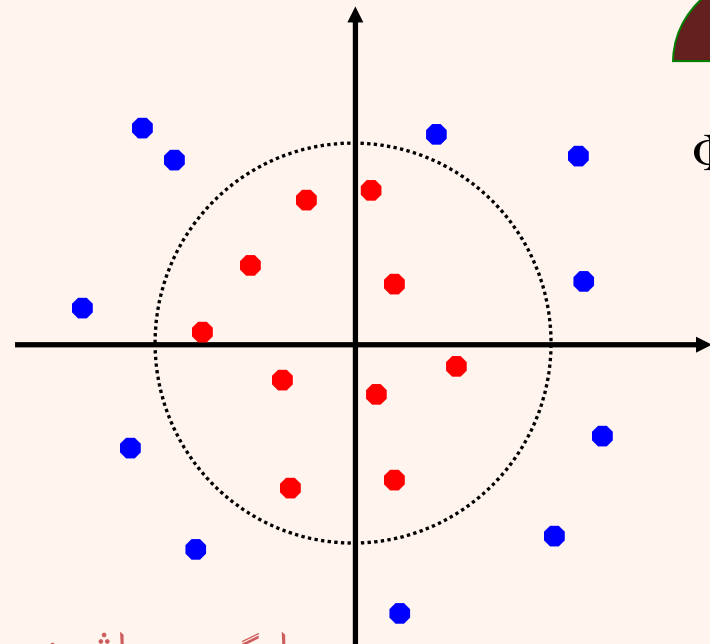


# نگاشت به فضای بالاتر

- همواره فضای ورودی می‌تواند به فضایی با ابعاد بالاتر نگاشت گردد.
- این نگاشت می‌تواند به صورتی باشد که در این فضای جدید ورودی‌ها قابلیت جداسازی داشته باشند.

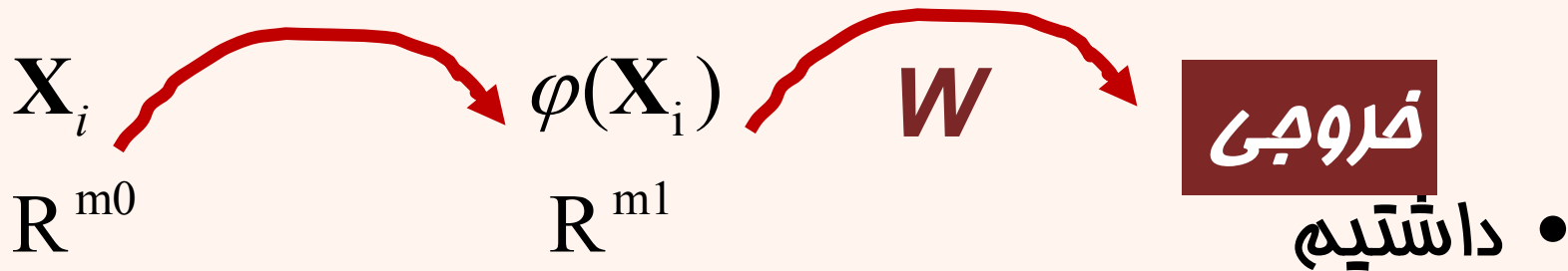


$$\Phi: \mathbf{X} \rightarrow \phi(\mathbf{X})$$



دانشگاه  
تهران  
پیشرو

# نگاشت به فضای بالاتر



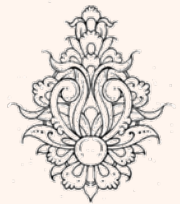
$$\mathbf{W}^T \mathbf{X} + b = 0$$

- هنگامی که ورودی ها به فضای دیگری نگاشت شوند برای نگاشت جدید خواهیم داشت:

$$\varphi(\mathbf{X}) = [\varphi_1(\mathbf{X}), \varphi_2(\mathbf{X}), \dots, \varphi_{m_1}(\mathbf{X})]^T$$

- در این حالت هدف یافتن رویه‌ی جداسازی است به‌گونه‌ای که:

$$\sum_{j=1}^{m_1} w_j \varphi_j(\mathbf{X}) + b = 0$$



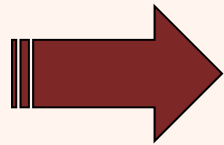
# نگاشت به فضای بالاتر

$$\sum_{j=1}^{m1} w_j \varphi_j(\mathbf{X}) + b = 0$$

• با فرض  $\varphi_0(\mathbf{X}) = 1$

• خواهیم داشت:

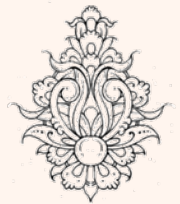
$$\sum_{j=0}^{m1} w_j \varphi_j(\mathbf{X}) = 0$$



$$\mathbf{W}^T \Phi(\mathbf{X}) = 0$$

$$\Phi(\mathbf{X}) = [1, \Phi(\mathbf{X})]^T$$

$$\mathbf{W} = [b, w_1, w_2, \dots, w_{m1}]^T$$



# نگاشت به فضای بالاتر

- در این مرحله تمامی شروط و قیودی که برای جداسازی خطی در نظر گرفتیم وجود دارد تنها به ازای  $x_i$  ها  $\varphi_i(x_i)$  در نظر گرفته می‌شود:

$$d_i \left[ \sum_{j=1}^{m_1} w_j \varphi_j(\mathbf{X}_i) - 1 \right] \geq 0$$

$$\mathbf{W}_{opt} = \sum_{i=1}^N \alpha_i \cdot d_i (\varphi_j(\mathbf{X}_i))$$

اسکالر  $\searrow$   $m_1 \times 1$

$$\mathbf{W}_{opt}^T \Phi(\mathbf{X}) = 0 \quad \Rightarrow \quad \sum_{i=1}^N \alpha_i \cdot d_i \Phi^T(\mathbf{X}_i) \Phi(\mathbf{X}) = 0$$





# نگاشت به فضای بالاتر

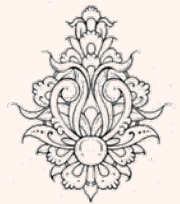
$$\sum_{i=1}^N \alpha_i \cdot d_i \Phi^T(\mathbf{X}_i) \Phi(\mathbf{X}) = 0$$

$$K(\mathbf{X}_i, \mathbf{X}_j) = \varphi(\mathbf{X}_i)^T \varphi(\mathbf{X}_j)$$

$$\sum_{i=1}^N \alpha_i \cdot d_i K(\mathbf{X}_i, \mathbf{X}) = 0$$

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(\mathbf{X}_i, \mathbf{X}_j)$$

تابع kernel، تابعی است که معادل ضرب داخلی در بردار  
خصیصه است.



# مثال

$$\mathbf{x}=[x_1 \ x_2];$$

$$K(\mathbf{x}_i, \mathbf{x}_j)=(1 + \mathbf{x}_i^T \mathbf{x}_j)^2,$$

$$K(\mathbf{x}_i, \mathbf{x}_j)= \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j):$$

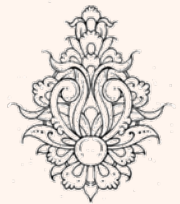
$$\begin{aligned} K(\mathbf{x}_i, \mathbf{x}_j) &= (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 = 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2x_{i1} x_{j1} + 2x_{i2} x_{j2} \\ &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] \\ &= \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j), \end{aligned}$$

$$\text{where } \boldsymbol{\varphi}(\mathbf{x}) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2]$$

## Mercer's theorem:

Every semi-positive definite symmetric function is a kernel

$K(\mathbf{x}_1, \mathbf{x}_1)$	$K(\mathbf{x}_1, \mathbf{x}_2)$	$K(\mathbf{x}_1, \mathbf{x}_3)$	...	$K(\mathbf{x}_1, \mathbf{x}_n)$
$K(\mathbf{x}_2, \mathbf{x}_1)$	$K(\mathbf{x}_2, \mathbf{x}_2)$	$K(\mathbf{x}_2, \mathbf{x}_3)$		$K(\mathbf{x}_2, \mathbf{x}_n)$
...	...	...	...	...
$K(\mathbf{x}_n, \mathbf{x}_1)$	$K(\mathbf{x}_n, \mathbf{x}_2)$	$K(\mathbf{x}_n, \mathbf{x}_3)$	...	$K(\mathbf{x}_n, \mathbf{x}_n)$



# نگاشت به فضای بالاتر

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \underbrace{\Phi(\mathbf{X}_i) \Phi(\mathbf{X}_j)}_{K(\mathbf{X}_i, \mathbf{X}_j)}$$

$$K_{N \times N} = \left\{ K(\mathbf{X}_i, \mathbf{X}_j) \right\}_{i,j=1}^N$$

ماتریس متقارن

هدف یافتن ضرایب لاگراثر بیشینه در عبارت زیر است:

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(\mathbf{X}_i, \mathbf{X}_j)$$

$$\sum_{i=1}^N \alpha_i d_i = 0$$

با در نظر گرفتن قيود زیر

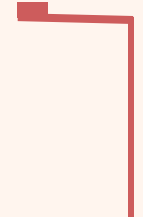
kernel trick

در صورت یافتن تابع kernel مناسب بدون این که در لایه مایل در فضایی با ابعاد بالاتر شویم، تنها از نتیجه این نگاشت

بهره می‌بریم.

$$g(\mathbf{x}) = \sum \alpha_i d_i K(\mathbf{X}_i, \mathbf{X}) + b$$





**TABLE 6.1** Summary of Inner-Product Kernels

Type of support vector machine	Inner product kernel $K(\mathbf{x}, \mathbf{x}_i), i = 1, 2, \dots, N$	Comments
Polynomial learning machine	$(\mathbf{x}^T \mathbf{x}_i + 1)^p$	Power $p$ is specified <i>a priori</i> by the user
Radial-basis function network	$\exp\left(-\frac{1}{2\sigma^2} \ \mathbf{x} - \mathbf{x}_i\ ^2\right)$	The width $\sigma^2$ , common to all the kernels, is specified <i>a priori</i> by the user
Two-layer perceptron	$\tanh(\beta_0 \mathbf{x}^T \mathbf{x}_i + \beta_1)$	Mercer's theorem is satisfied only for some values of $\beta_0$ and $\beta_1$

# ساختار SVM

