

# ماشین بردار پشتیبان (Support Vector Machine)

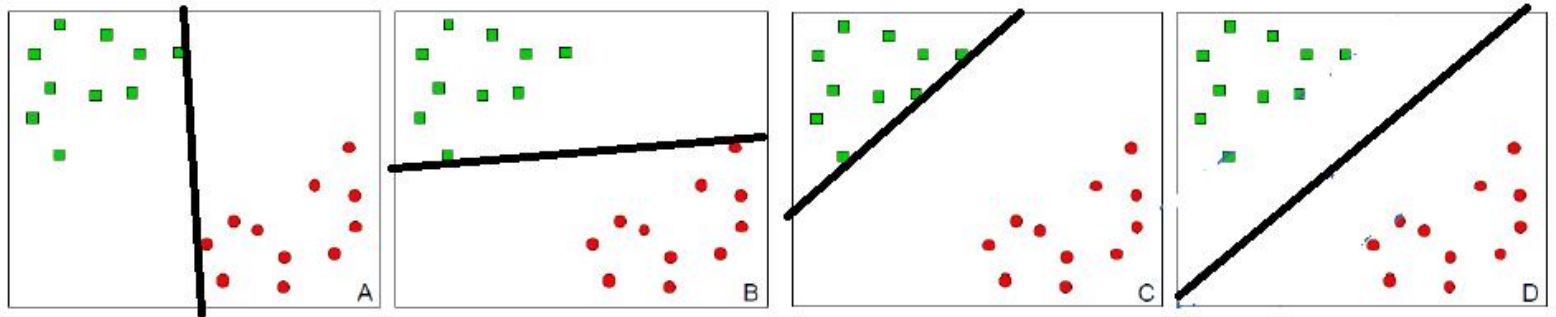
❖ یادگیری با سرپرست (SVM)

❖ در برخی از مسائل غیرخطی پیچیده کاربرد آن از رگرسیون لجستیک و شبکه عصبی بهتر است.

❖ کاربرد SVM عمدتاً در مواردی است که داده ها تفکیک پذیر خطی نباشند.

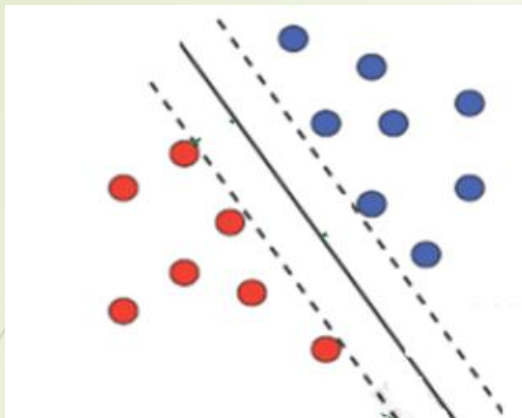
❖ SVM فضای مسئله را تغییر می دهد که نمونه های جدید تفکیک پذیر خطی شوند.

❖ در SVM دسته ها طوری تفکیک می شوند که فاصله بین آنها حداکثر شود

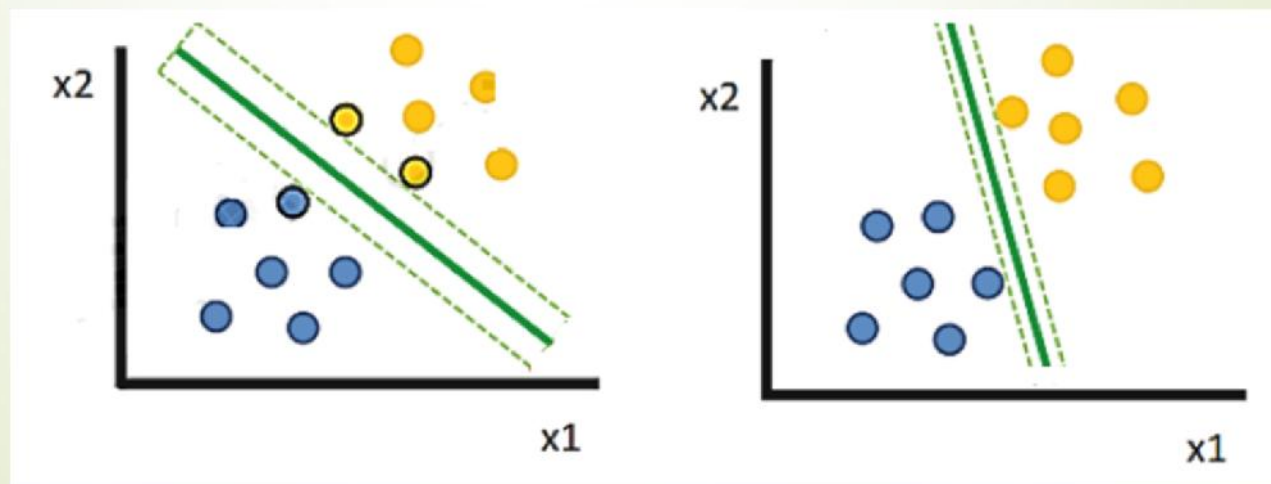


➤ SVM خط راستی که بیشترین حاشیه را با داده دسته ها دارد.

➤ فاصله خط جدا کننده تا تا نمونه دو دسته برابر است.



➤ به SVM دسته بندی با حاشیه پهن نیز می گویند چون خطای تعمیم را کاهش می دهد.



## طراحی تابع هزینه SVM

► SVM با اصلاحات در الگوریتم لجستیک حاصل می شود.

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad 0 \leq h_{\theta}(x) \leq 1$$

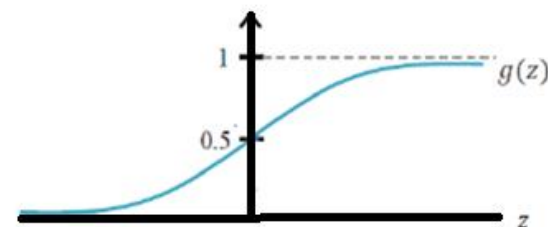
$$\text{if } h_{\theta}(x) \geq 0.5 : y = 1$$

$$\text{if } h_{\theta}(x) < 0.5 : y = 0$$

$$\text{if } \theta^T x \geq 0 : y = 1$$

$$\text{if } \theta^T x < 0 : y = 0$$

$$g(z) = \frac{1}{1 + e^{-z}}$$



Sigmoid function

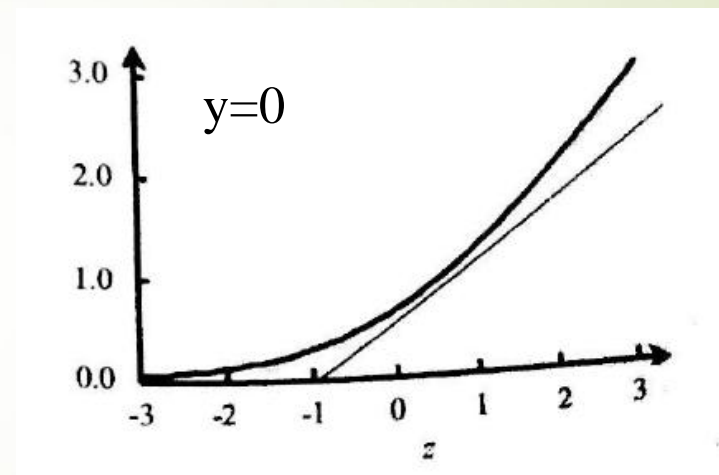
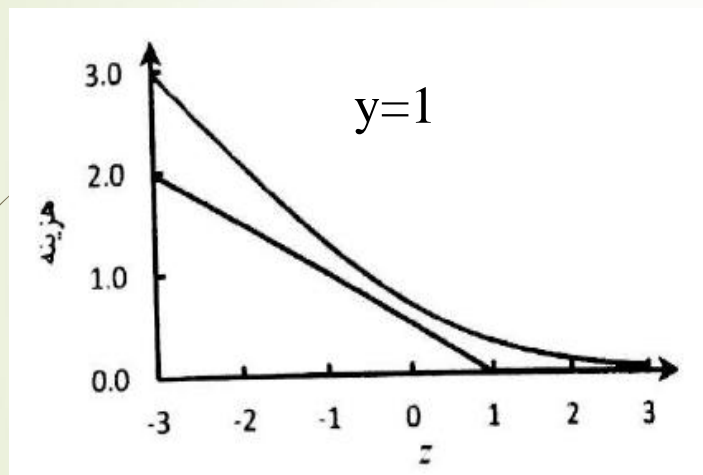
## تابع هزینه

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & y = 1 \\ -\log(1 - h_{\theta}(x)) & y = 0 \end{cases}$$

$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[ -y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

$$J(\theta) = -\left[\frac{1}{m} \sum_{i=1}^m y \log\left(\frac{1}{1 + e^{\theta^T x}}\right) + (1 - y) \log\left(1 - \frac{1}{1 + e^{\theta^T x}}\right)\right]$$



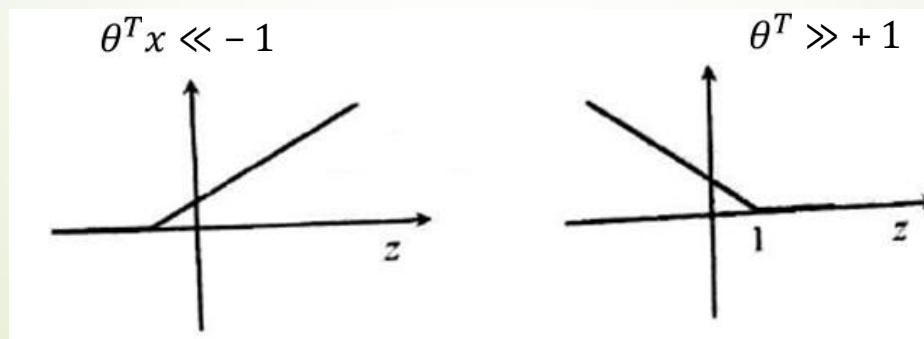
ایده ساخت SVM این است که بجای تابع هزینه لگاریتمی از این خطوط به عنوان تابع هزینه استفاده شود.

$$\text{Min } J(\theta) = \min \left( \left[ \frac{-1}{m} \sum_{i=1}^m y \log \left( \frac{1}{1 + e^{\theta^T x}} \right) + (1 - y) \log \left( 1 - \frac{1}{1 + e^{\theta^T x}} \right) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta^2 \right)$$

$$\min \left( \sum_{i=1}^m y \log \left( \frac{1}{1 + e^{\theta^T x}} \right) + (1 - y) \log \left( 1 - \frac{1}{1 + e^{\theta^T x}} \right) + \frac{\lambda}{2} \sum_{j=1}^n \theta^2 \right)$$

$$\min_{\theta} (A + \lambda B)$$

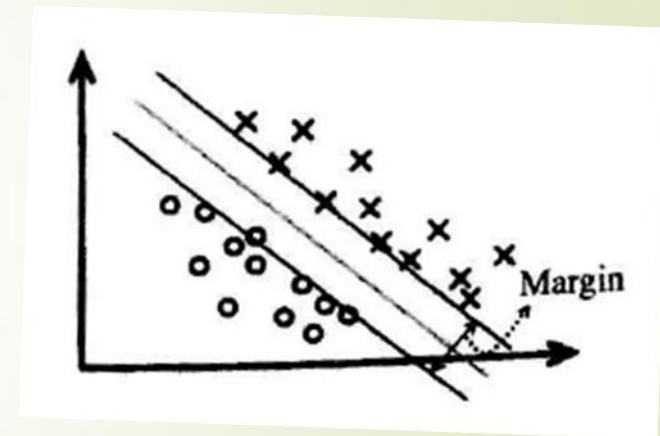
$$\min_{\theta} (CA + B) \quad C = \frac{1}{\lambda}$$



- برای تعیین مرز تفکیک تا جایی که امکان دارد مقدار تابع هزینه باید حداقل شود.
- اگر مقدار  $C$  کم شود  $A$  به سمت صفر میل کرده و فقط باید مقدار  $B$  حداقل شود.

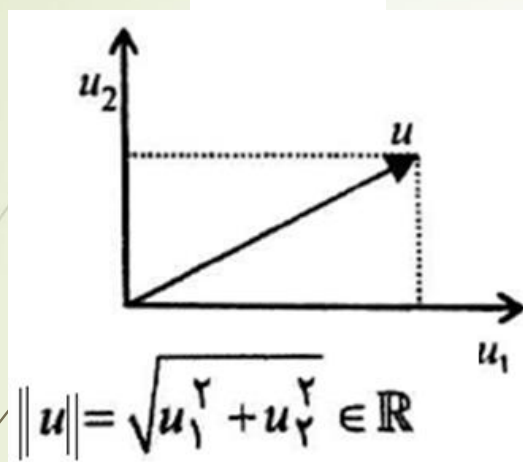
$$\min_{\theta} (CA + B) \quad C = \frac{1}{\lambda}$$

$$\begin{aligned} \min B &= \min \frac{\lambda}{2} \sum_{j=1}^n \theta^2 = \min \frac{1}{2} \|\theta\|^2 \\ \text{s.t.} \quad &\theta^T x^i \geq 1 \quad \text{if } y^i = 1 \\ &\theta^T x^i \leq -1 \quad \text{if } y^i = 0 \end{aligned}$$

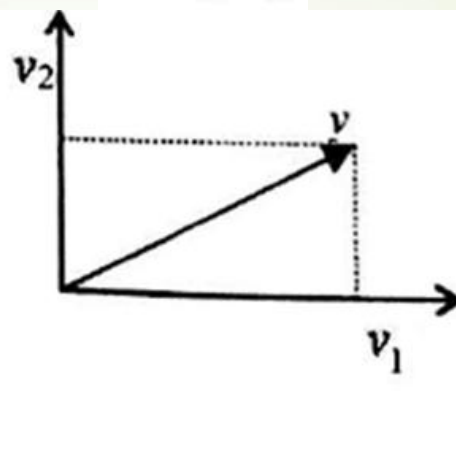


# بردارها

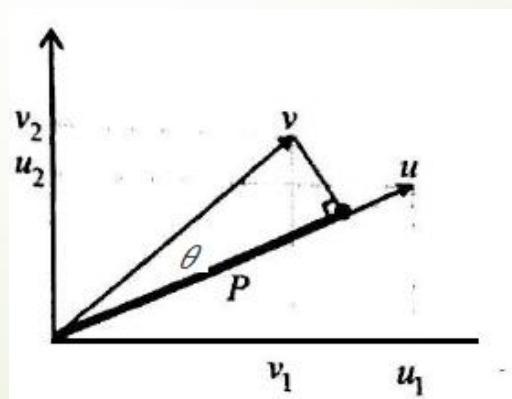
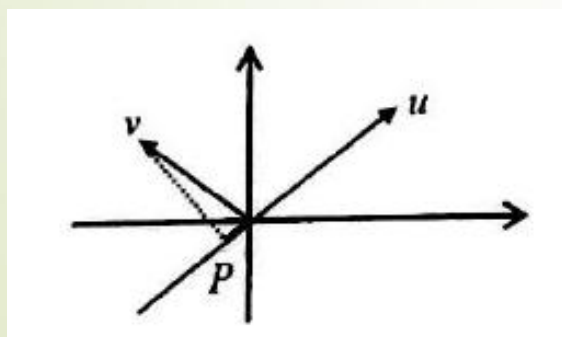
$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$



$$v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$



$$u^T v = \begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = u_1 v_1 + u_2 v_2$$

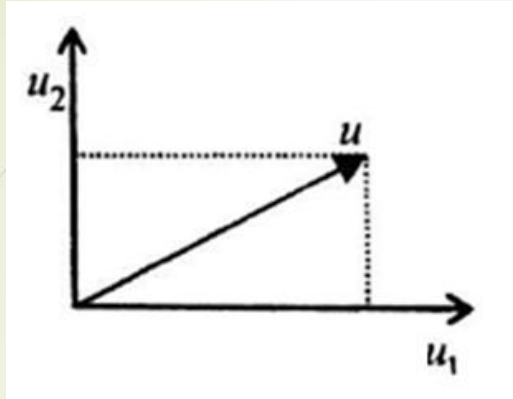


$$\cos(\theta) = \frac{\|P\|}{\|V\|}$$

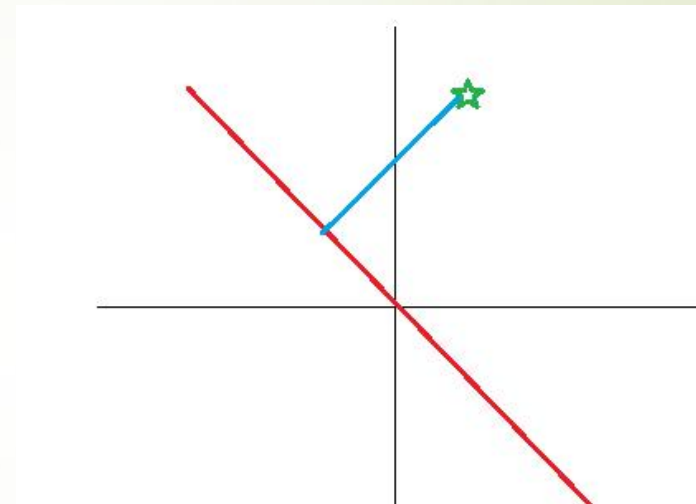
$$\|P\| = \|V\| \cos(\theta) = \frac{U \cdot V}{\|U\|}$$



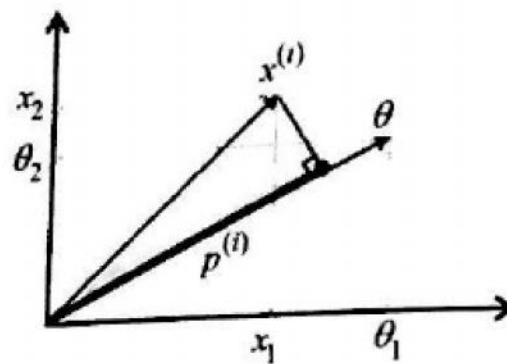
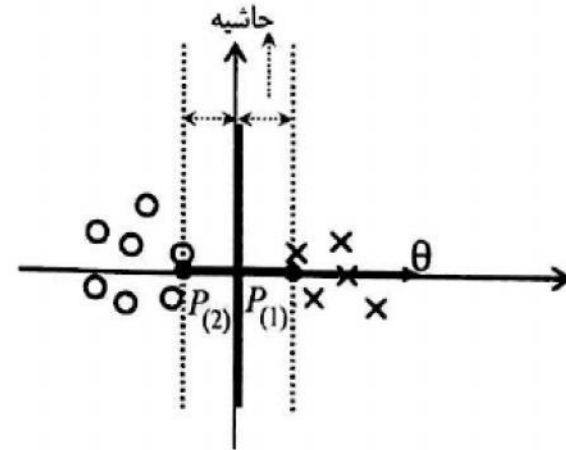
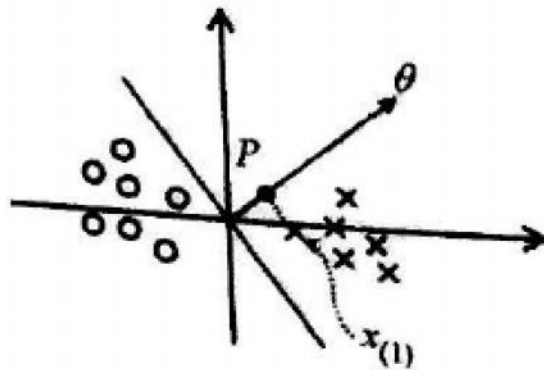
## جهت بردار و فاصله نقطه از خط



$$w = \left( \frac{u_1}{||u||}, \frac{u_2}{||u||} \right)$$



$$\frac{\theta^T X_A}{||\theta||}$$



$$\min \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2$$

$$\theta^T x^{(i)} = P^{(i)} \cdot \theta = \theta_1 x_1 + \theta_2 x_2$$

## ضرایب لاگرانژ

برای پیدا کردن مینیمم یا ماکزیمم یک تابع با توجه به محدودیت‌ها، از لاگرانژ استفاده می‌کنیم.

optimization problem

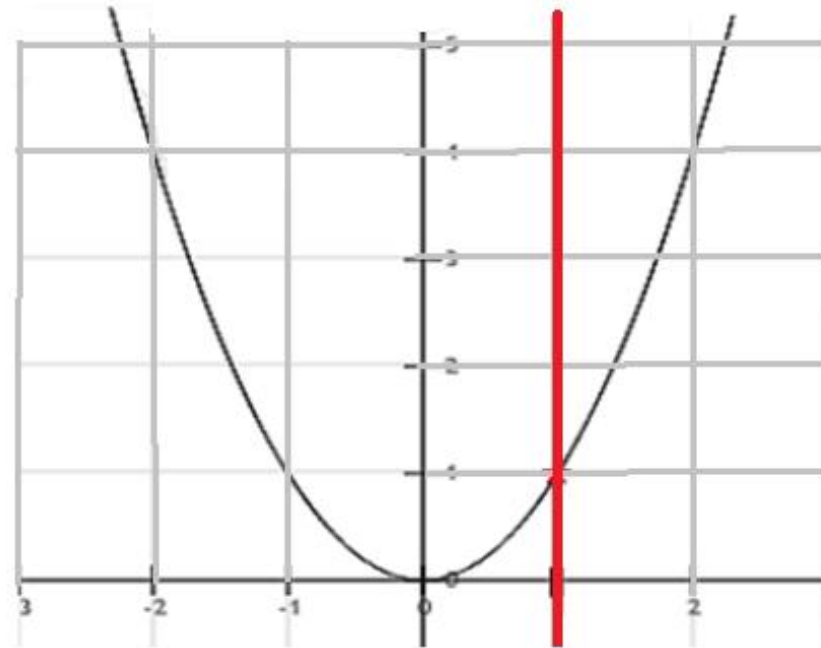
$$\begin{array}{l} \min f(x, y) \\ \text{subject to } g(x, y) = 0. \end{array}$$

$$\mathcal{L}(x, y, \lambda) = f(x, y) - \lambda \cdot g(x, y)$$

$$\nabla_{x,y,\lambda} \mathcal{L}(x, y, \lambda) = 0$$

$$\begin{array}{l} L(x, \lambda) = f(x) - \sum_i \lambda_i g_i(x) \\ \nabla L(x, \lambda) = 0 \end{array}$$

## مثال



$$f(x) = x^2$$
$$g(x) : x = 1$$

## مثال

$$f(x, y) = x^2 + y^2$$

$$\begin{cases} g_1(x, y) = x + 1 = 0 \\ g_2(x, y) = y + 1 = 0 \end{cases}$$

## اگر شروط نامساوی باشند؟

$$\min f(x, y)$$

$$g(x) \geq 0 \Rightarrow \lambda \geq 0$$

$$g(x) \leq 0 \Rightarrow \lambda \leq 0$$

مثال:

$$f(x, y) = x^3 + y^2$$

$$g(x, y) = x^2 - 1 \geq 0$$

$$f(x, y) = x^3 + y^3$$

$$g_1(x, y) = x^2 - 1 \geq 0$$

$$g_2(x, y) = y^2 - 1 \geq 0$$

## بهینه سازی SVM

$$\min \frac{1}{2} ||\theta||^2$$

$$\text{s.t.} \quad \begin{aligned} \theta^T x^i &\geq 1 & \text{if } y^i = 1 \\ \theta^T x^i &\leq -1 & \text{if } y^i = 0 \end{aligned}$$

$$\text{s.t.} \quad \begin{aligned} \theta^T x^i &\geq 1 & \text{if } y^i = 1 \\ \theta^T x^i &\leq -1 & \text{if } y^i = -1 \end{aligned}$$

$$\text{s.t.} \quad y^i(\theta^T x^i) \geq 1 = y^i(\theta^T x^i) - 1 \geq 0 = y^i(\theta_0 + x_1^i \theta_1 + \dots + x_n^i \theta_n) - 1 \geq 0$$

$$L(\theta, \theta_0, \lambda) = \frac{1}{2} ||\theta^2|| - \sum_{i=1}^m \lambda_i (y^i(\theta^T x^i) - 1) \quad \lambda_i \geq 0$$

## بهینه سازی SVM

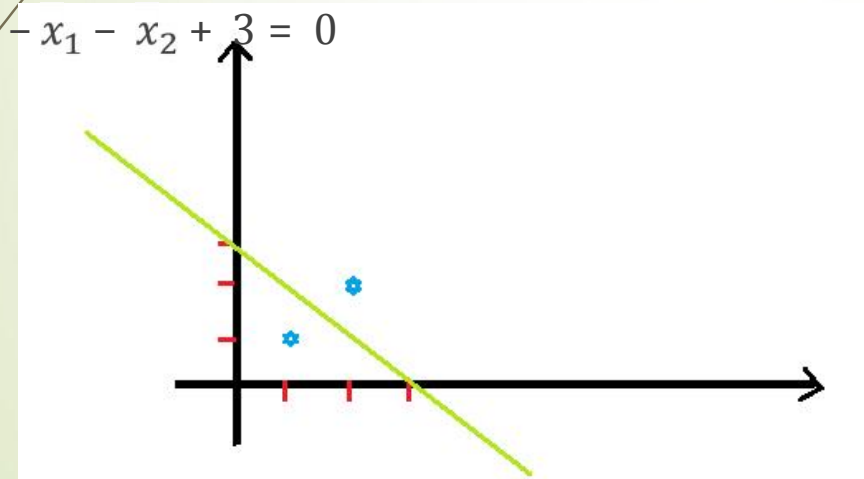
$$L(\theta, \theta_0, \lambda) = \frac{1}{2} ||\theta^2|| - \sum_{i=1}^m \lambda_i (y^i (\theta^T x^i) + \sum_{i=1}^m \lambda_i \quad \lambda_i \geq 0$$

مثال:

داده ها:

$(1,1) \rightarrow +1$

$(2,2) \rightarrow -1$





## دوگان لاگرانژ

$$L(\theta, \theta_0, \lambda) = \frac{1}{2} \|\theta^2\| - \sum_{i=1}^m \lambda_i (y^i (\theta^T x^i) + \sum_{i=1}^m \lambda_i \quad \lambda_i \geq 0$$

$$\theta = \sum_{i=1}^m \lambda_i y^i x^i, \quad \sum_{i=1}^m \lambda_i y^i = 0$$

$$L_d = \sum_{i=1}^m \lambda_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y_i y_j (x_i \cdot x_j)$$

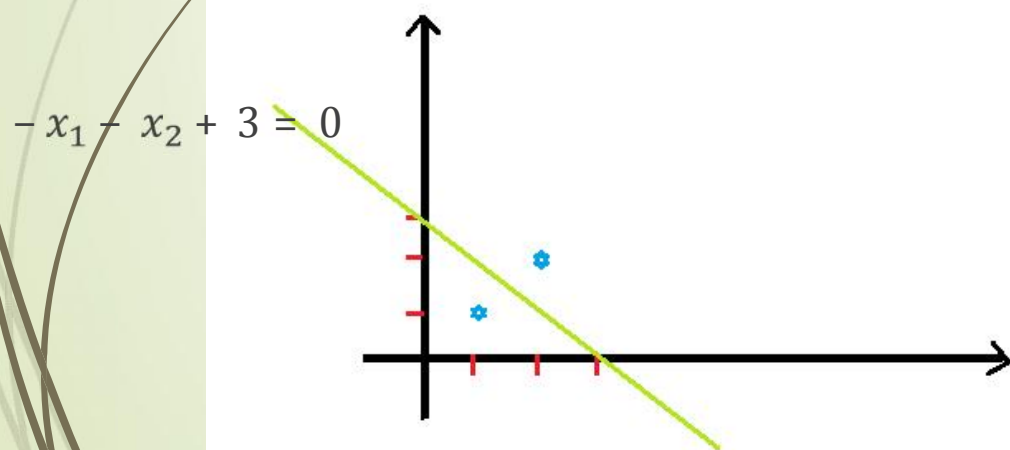
$$L_d = \sum_{i=1}^m \lambda_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y_i y_j (x_i \cdot x_j)$$

مثال:

داده ها:

$(1,1) \rightarrow +1$

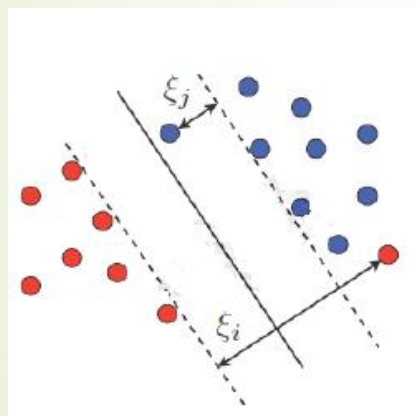
$(2,2) \rightarrow -1$



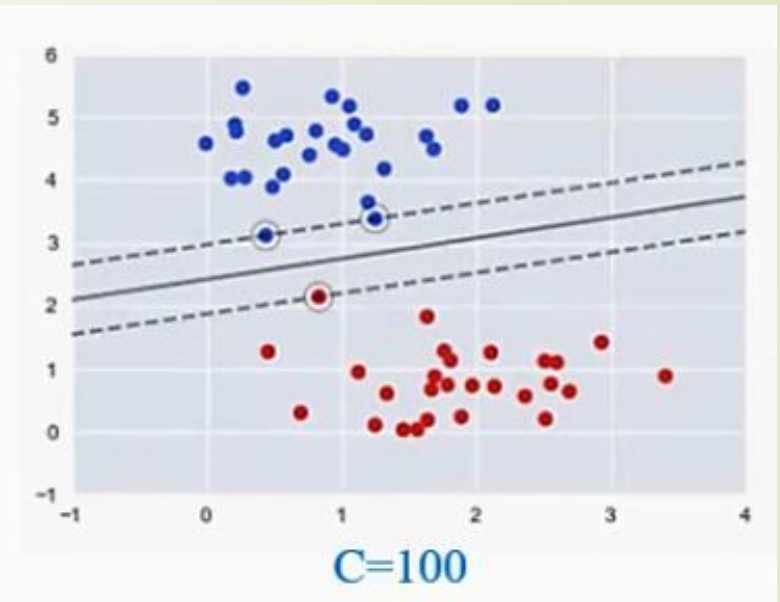
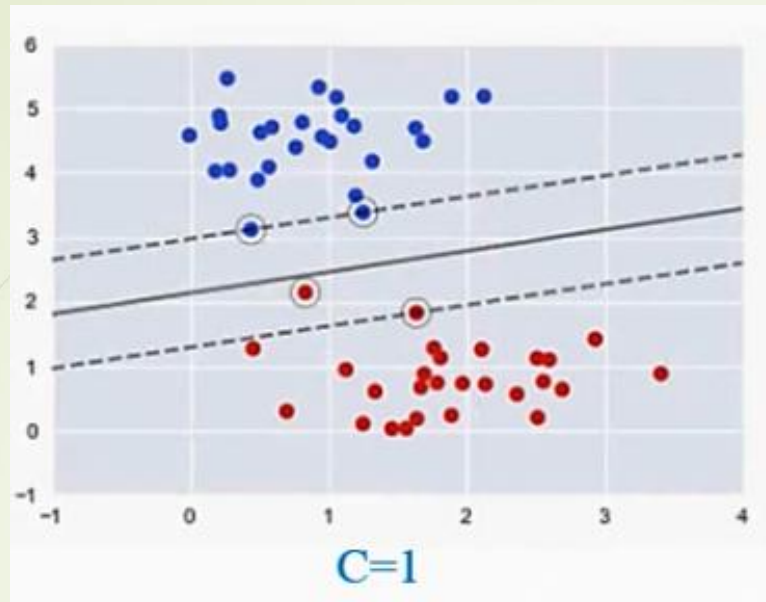
## ماشین بردار پشتیبان با حاشیه نرم

در SVM قبلی داده ها کاملاً به دو دسته تفکیک می شوند.

در SVM نرم تعدادی از داده ها می توانند مرزها را رعایت نکنند.



$$\min \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \varepsilon_i$$
$$y^i(\theta^T x^i) \geq 1 - \varepsilon_i$$
$$\varepsilon_i \geq 0$$



$$\min \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \varepsilon_i$$

$$S.t. \quad y^i(\theta^T x^i) \geq 1 - \varepsilon_i$$

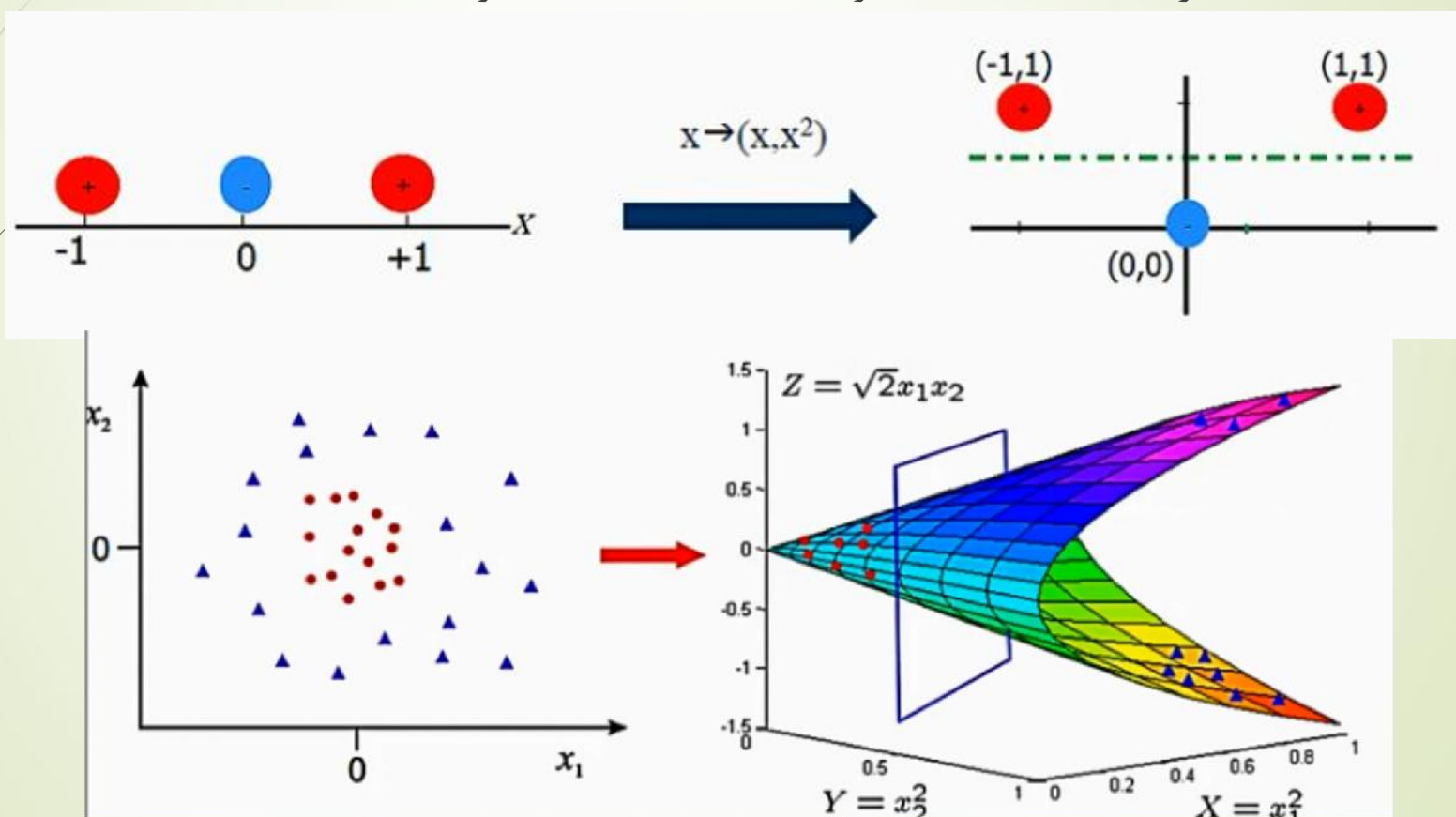
$$\varepsilon_i \geq 0$$

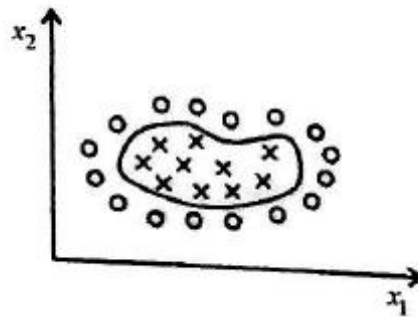
$$L(\theta, \theta_0, \lambda, \varepsilon) = \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \varepsilon_i - \sum_{i=1}^m \lambda_i (y^i(\theta^T x^i) - 1 + \varepsilon_i) - \sum_{i=1}^m \mu_i \varepsilon_i \quad \lambda_i, \varepsilon_i \geq 0$$

## کرنل (Kernel)

در واقعیت داده ها به صورت خطی تفکیک پذیر نیستند.

نگاشت داده ها از یک فضا با ابعاد کمتر به یک فضا با ابعاد بیشتر





$$\theta_0 + \theta_1 x_1 + \theta_r x_r + \theta_r x_1 x_r + \theta_r x_1^r + \theta_\delta x_r^r + \dots \geq .$$

$$h_\theta(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \theta_r x_r + \theta_r x_1 x_r + \dots \geq . \\ 0 & \text{if } \theta_0 + \theta_1 x_1 + \theta_r x_r + \theta_r x_1 x_r + \dots < . \end{cases}$$

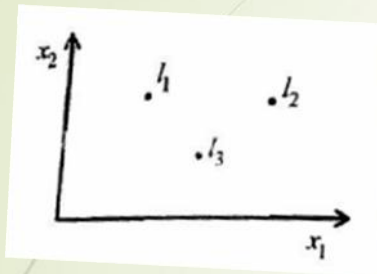
$$\theta_0 + \theta_1 f_1 + \theta_r f_r + \theta_r f_r + \dots$$

$$f_1 = x_1, f_r = x_r, f_r = x_1 x_r, f_r = x_1^r, f_\delta = x_r^r, \dots$$

می توان ویژگی های  $f$  را به صورت دیگر نیز تعیین کرد؟

براساس شباهت با یکسری نقاط مشخص شده در فضای داده ویژگی های جدید را بدست آورد.

فرض داده ها به جای دو ویژگی  $(X_1, X_2)$  با سه ویژگی جایگزین شوند  $(l_1, l_2, l_3)$ .



$F_1$  میزان شباهت داده ها با شاخص  $l_1$

$F_2$  میزان شباهت داده ها با شاخص  $l_2$

$F_3$  میزان شباهت داده ها با شاخص  $l_3$

$$k(x_i, x_j) = (x_i \cdot x_j + c)^d \quad c > 0, d \in \mathbb{N}$$

از فرمول های مختلف می توان برای سنجش شباهت

نمونه ها استفاده کرد که به آن **کرنل** گویند.

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x}' - \mathbf{x}\|^2}{2\sigma^2}\right)$$

$$K(\mathbf{x}, \mathbf{x}') = \tanh(a(\mathbf{x} \cdot \mathbf{x}') + b) \quad a, b > 0$$

## کرنل گوسی (Gaussian Kernel)

استفاده از فاصله اقلیدسی در تابع کرنل گوسی

برای هر نمونه مانند  $x$ :

$$f_1 = \text{sim}(x, l^1) = \exp\left(-\frac{\|x - l^1\|^2}{2\sigma^2}\right) = e^{\left(-\frac{\|x - l^1\|^2}{2\sigma^2}\right)}$$

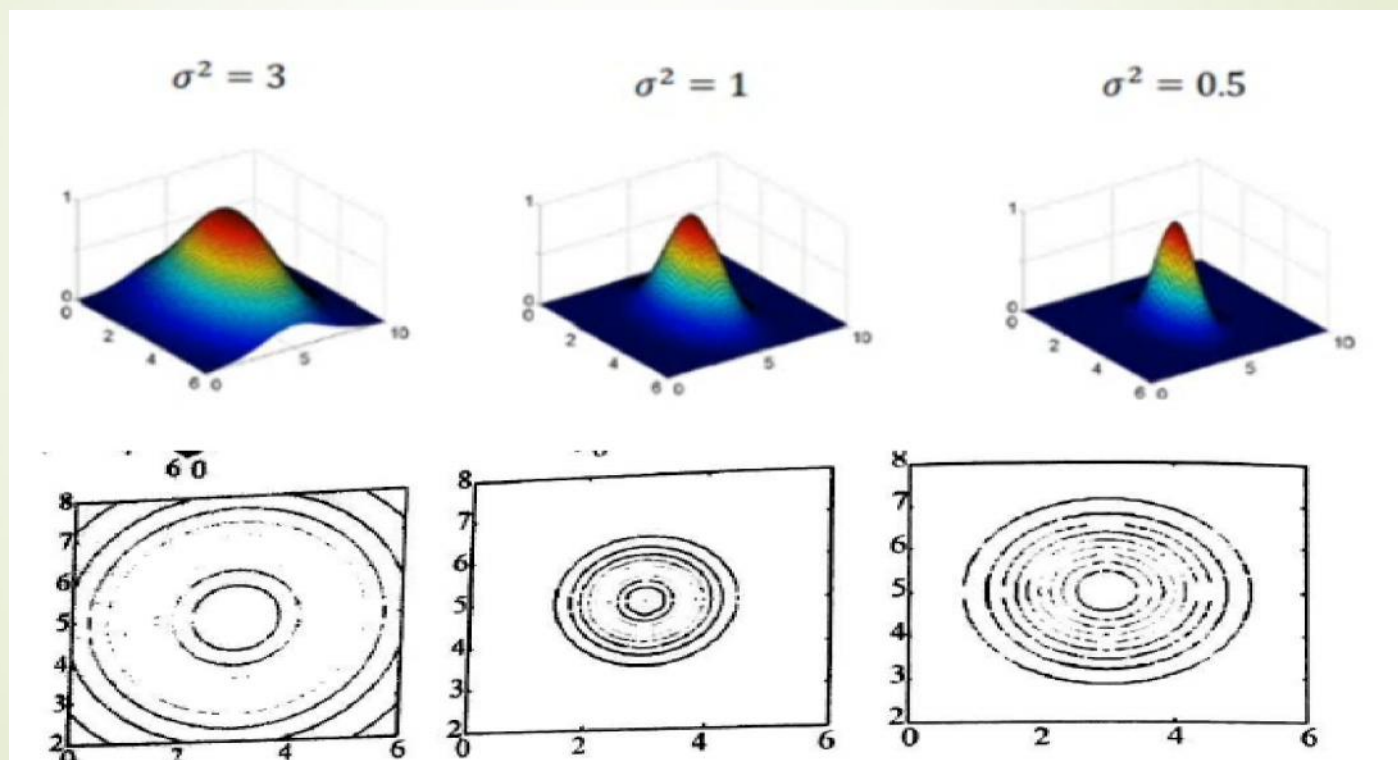
$$f_2 = \text{sim}(x, l^2) = \exp\left(-\frac{\|x - l^2\|^2}{2\sigma^2}\right) = e^{\left(-\frac{\|x - l^2\|^2}{2\sigma^2}\right)}$$

$$f_3 = \text{sim}(x, l^3) = \exp\left(-\frac{\|x - l^3\|^2}{2\sigma^2}\right) = e^{\left(-\frac{\|x - l^3\|^2}{2\sigma^2}\right)}$$

اگر  $x$  به  $l$  نزدیک باشد  $f$  متناظر تقریباً یک می شود و اگر دور باشد تقریباً صفر خواهد شد.



- هر چه  $\sigma$  کوچکتر باشد بایاس کمتر، واریانس بیشتر و خروجی سریعتر کاهش می یابد.
- هر چه  $\sigma$  بزرگتر باشد بایاس بیشتر، واریانس کمتر و خروجی کندتر کاهش می یابد.



## تعیین مراکز کرنل

تعیین  $m$  نقطه شاخص دلخواه

در نظر گرفتن هر داده به عنوان یک نقطه شاخص

$$x^1 = l^1, x^2 = l^2, \dots, x^m = l^m$$

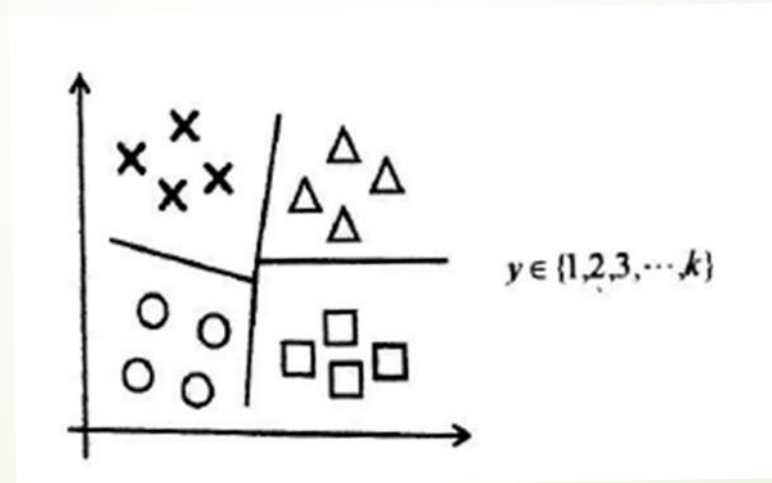
$$f_j^i = \text{sim}(x^i, l^j) = \exp\left(-\frac{\|x^i - l^j\|^2}{2\sigma^2}\right) = e^{\left(-\frac{\|x^i - l^j\|^2}{2\sigma^2}\right)}$$

مثال:

$$(2, 3) \rightarrow +1, (1, 1) \rightarrow +1, (2, 2) \rightarrow -1$$

## دسته بندی چند دسته ای

روش یک دسته در مقابل بقیه



## مقایسه SVM و الگوریتم لجستیک

➤ اگر ویژگی ها زیاده از تعداد داده ها باشد ( $n > m$ ) می توان از الگوریتم لجستیک یا SVM بدون کرنل استفاده کرد (ایمیل ها).

➤ اگر ویژگی ها کم و تعداد داده ها متوسط باشد استفاده از SVM با کرنل گوسی توصیه می شود.

➤ اگر ویژگی ها کم و تعداد داده ها زیاد باشد می توان از لجستیک یا SVM بدون کرنل استفاده کرد (کرنل بار محاسبات را زیاد می کند).

➤ در موارد فوق می توان از شبکه عصبی استفاده کرد فقط آموزش آن کندتر است.

➤ SVM محدب است و همیشه جواب بهینه و یا نزدیک بهینه تولید می کند.