

ماشین‌های بردار پشتیبان

سید ناصر رضوی www.snrazavi.ir

۱۳۹۶

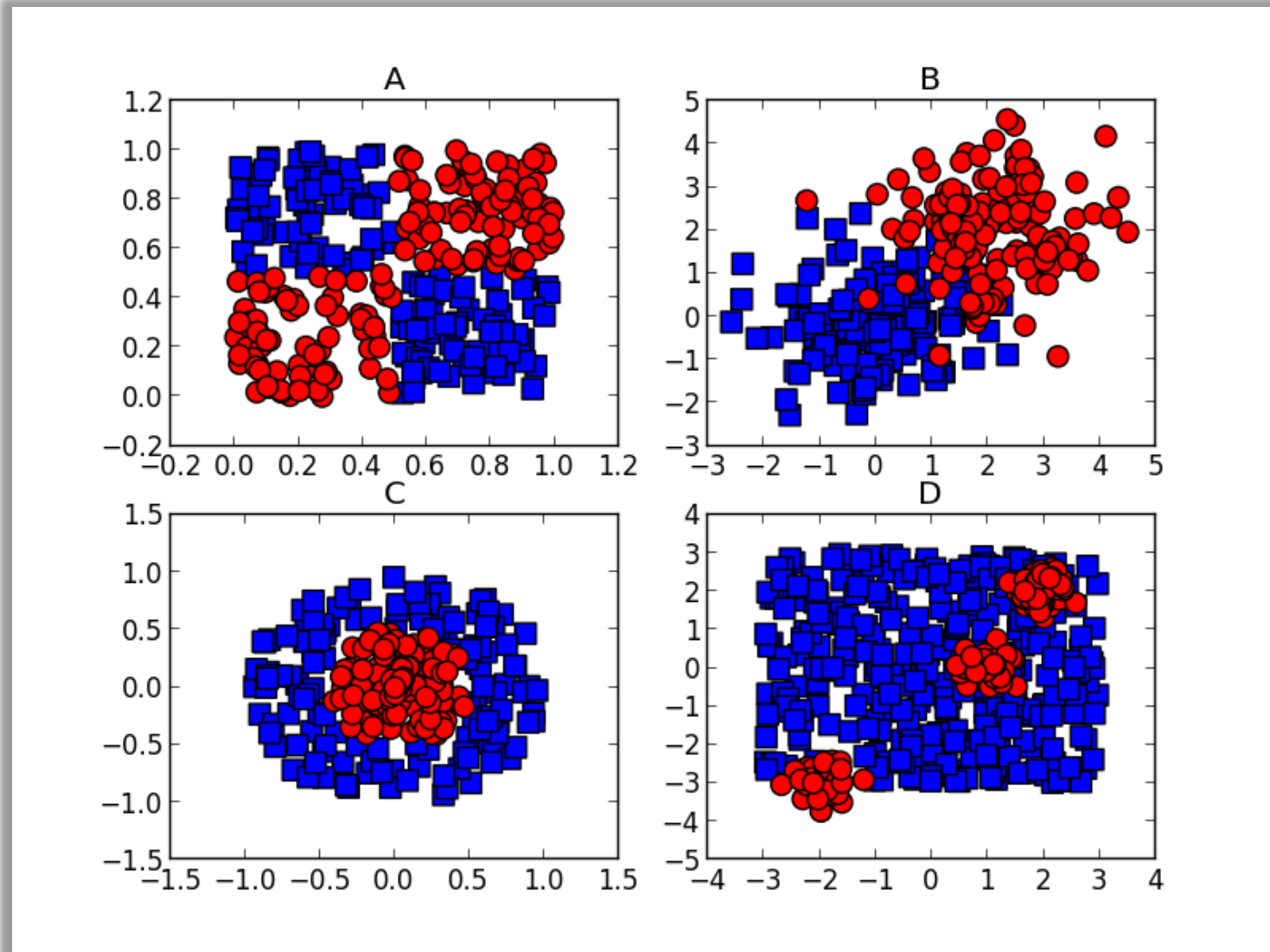
فهرست مطالب

- مفاهیم پایه‌ای
- مسائل SVM: مسئله اصلی و مسئله دوگان
- آموزش SVM های خطی و غیرخطی
- انتخاب پارامترها و تابع کرنل
- کلاس بندی چند کلاسی
- بحث و نتیجه گیری

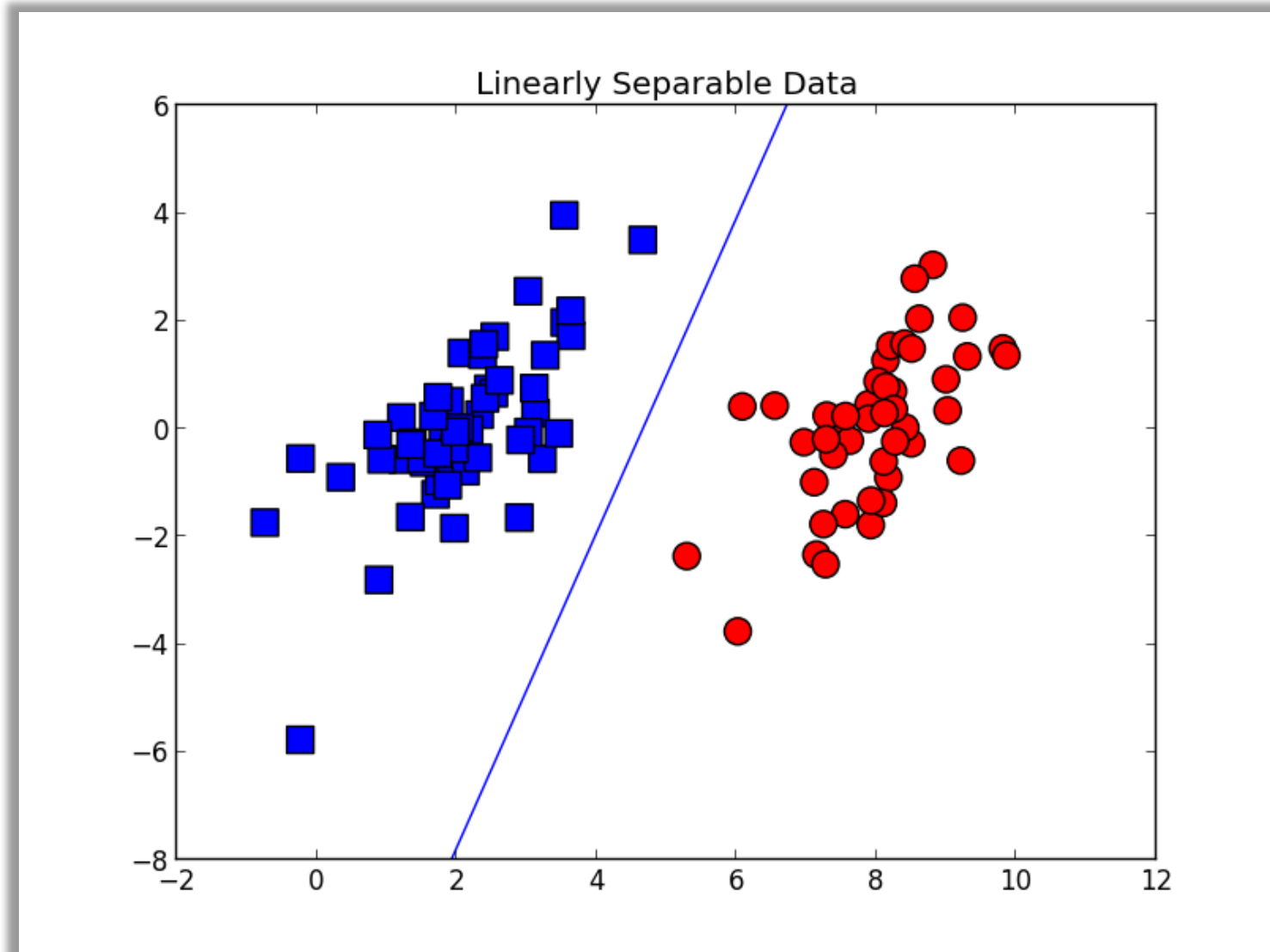
□ ماشین‌های بردار پشتیبان [اوپنیک، ۱۹۹۲]

- یکی از پرطرفدارترین الگوریتم‌های یادگیری ماشین!
- جداسازی بهتر داده‌ها نسبت به سایر روش‌های یادگیری ماشین (مسائل کلاس‌بندی)!
- استفاده از آن نسبتاً آسان است!
- استفاده از **ترفند کرنل**:
- کلاس‌بندی، رگرسیون، تخمین توزیع، کلاس‌بندی تک کلاسی و ...

داده‌های تفکیک ناپذیر خطی



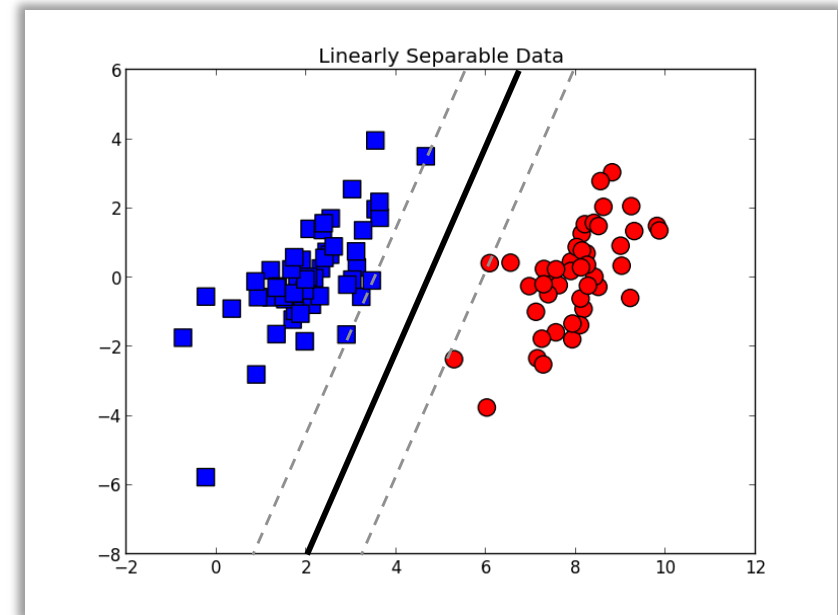
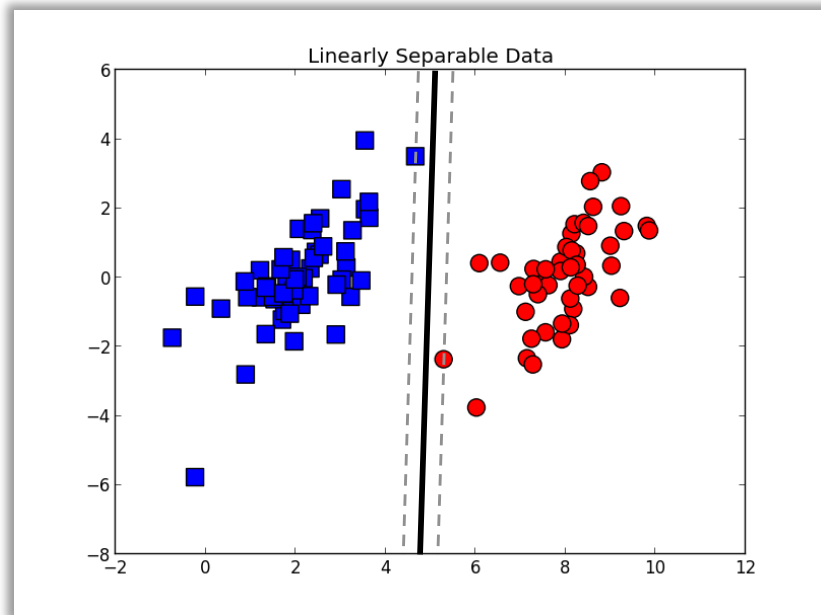
داده‌های تفکیک‌پذیر خطی



مرز تصمیم‌گیری بهینه

۶

□ س. کدام مرز تصمیم‌گیری بهتر است؟

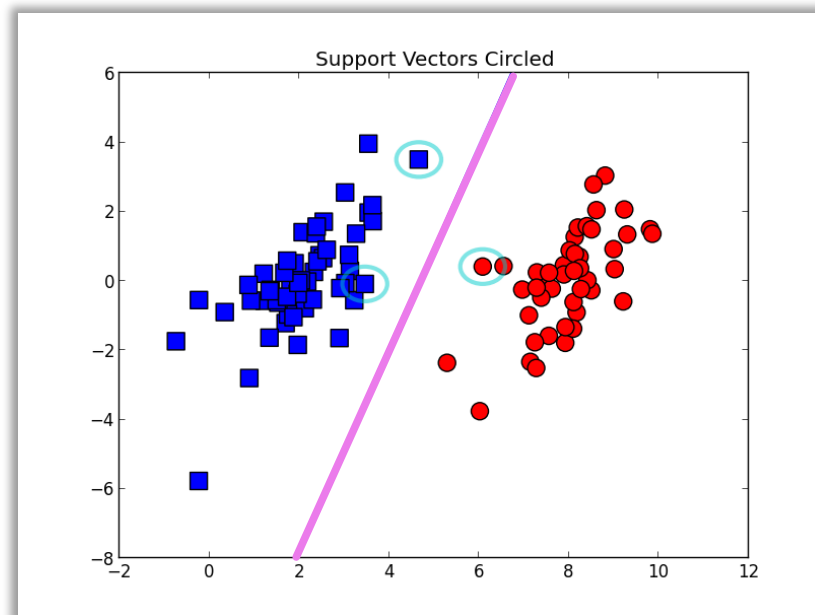


□ راه‌حل بیشترین حاشیه. بیشترین پایداری در برابر تخریب داده‌ها. [افزایش قابلیت تعمیم]

بردارهای پشتیبان

□ بردار پشتیبان. نزدیکترین نقاط به مرز تصمیم‌گیری.

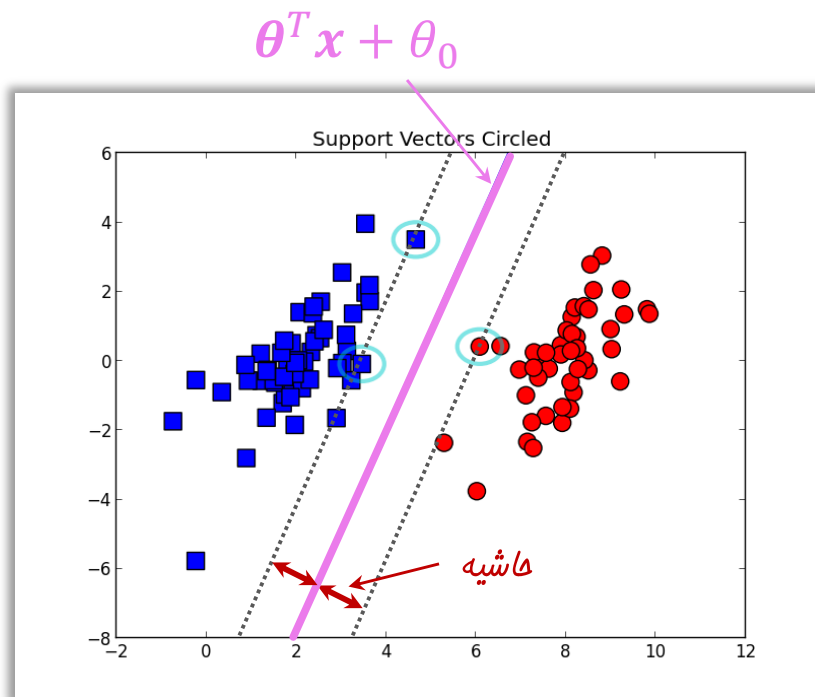
□ هدف. بیشینه کردن فاصله بردارهای پشتیبان از مرز تصمیم‌گیری.



ماشین‌های بردار پشتیبان: کلاس‌بند با بزرگ‌ترین حاشیه

□ حاشیه. فاصله بردارهای پشتیبان تا مرز تصمیم‌گیری.

□ هدف. بیشینه کردن فاصله بردارهای پشتیبان از مرز تصمیم‌گیری.



$$\frac{|\theta^T x + \theta_0|}{\|\theta\|} \geq \rho$$

مرز تصمیم‌گیری بهینه: نمادها

۹

□ نمونه‌های آموزشی.

$$X = (\mathbf{x}^t, y^t), \quad y^t = \begin{cases} +1 & \text{if } \mathbf{x}^t \in C_1 \\ -1 & \text{if } \mathbf{x}^t \in C_2 \end{cases}$$

□ هدف. یافتن بردار θ و مقدار θ_0 به طوری که:

$$\theta^T \mathbf{x}^t + \theta_0 \geq +1 \quad \text{for } y^t = +1$$

$$\theta^T \mathbf{x}^t + \theta_0 \leq -1 \quad \text{for } y^t = -1$$



$$y^t (\theta^T \mathbf{x}^t + \theta_0) \geq +1$$

تابع هدف

۱۰

□ هدف. بیشینه کردن فاصله بردارهای پشتیبان از مرز تصمیم‌گیری.

□ فاصله داده x از مرز تصمیم‌گیری:

$$\frac{|\boldsymbol{\theta}^T \mathbf{x} + \theta_0|}{\|\boldsymbol{\theta}\|} \geq \rho \Rightarrow |\boldsymbol{\theta}^T \mathbf{x} + \theta_0| \geq \rho \|\boldsymbol{\theta}\|$$

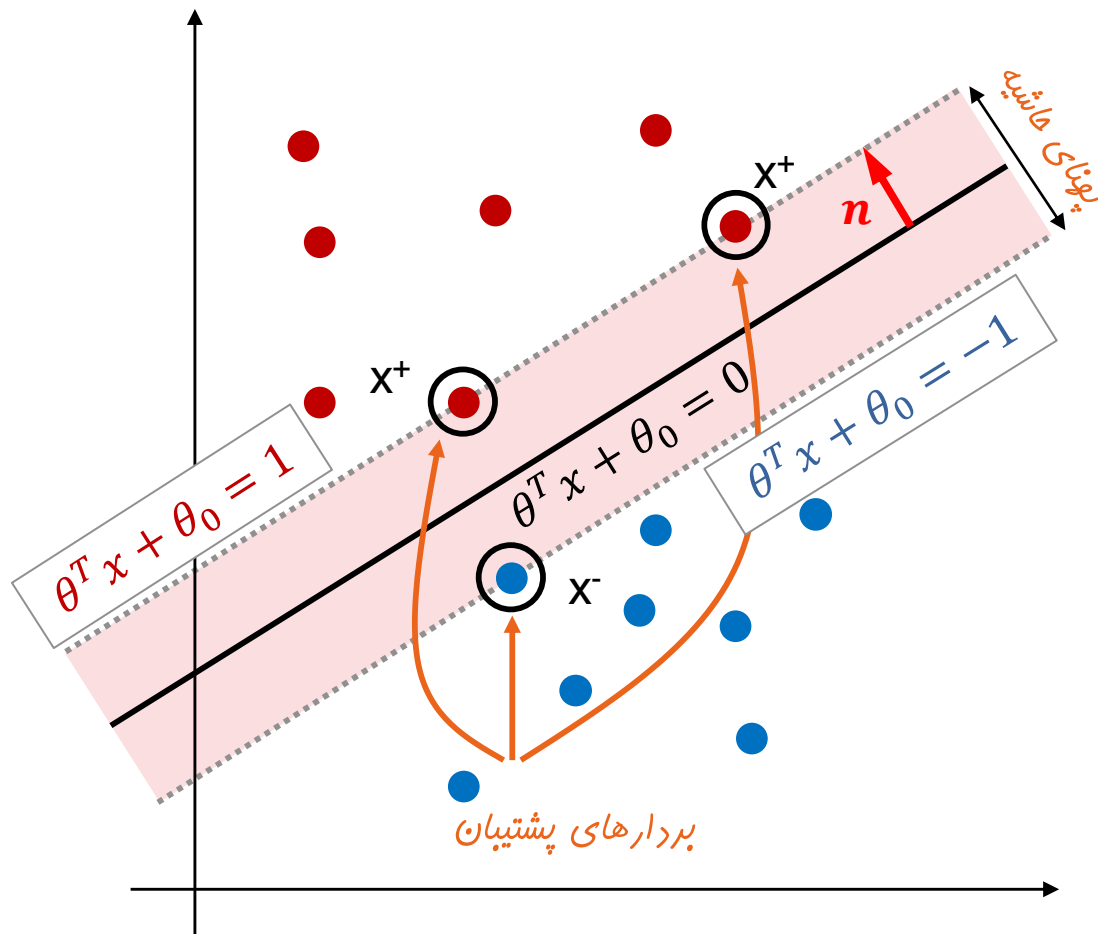
□ این معادله بی‌نهایت جواب دارد؛ اما با قرار دادن $\rho \|\boldsymbol{\theta}\| = 1$ ، خواهیم داشت:

$$\rho \|\boldsymbol{\theta}\| = 1 \Rightarrow \rho = \frac{1}{\|\boldsymbol{\theta}\|}$$

□ هدف. برای بیشینه کردن حاشیه، می‌توان اندازه بردار $\boldsymbol{\theta}$ را کمینه نمود.

□ محدودیت‌ها. مرز تصمیم‌گیری باید داده‌های دو کلاس را به درستی از یکدیگر تفکیک کند.

تابع هدف



□ می دانیم:

$$\theta^T x^+ + \theta_0 = +1$$

$$\theta^T x^- + \theta_0 = -1$$

□ بنابراین:

$$M = (x^+ - x^-) \cdot n$$

$$= (x^+ - x^-) \cdot \frac{\theta}{\|\theta\|} = \frac{2}{\|\theta\|}$$

تابع هدف: بیان رسمی

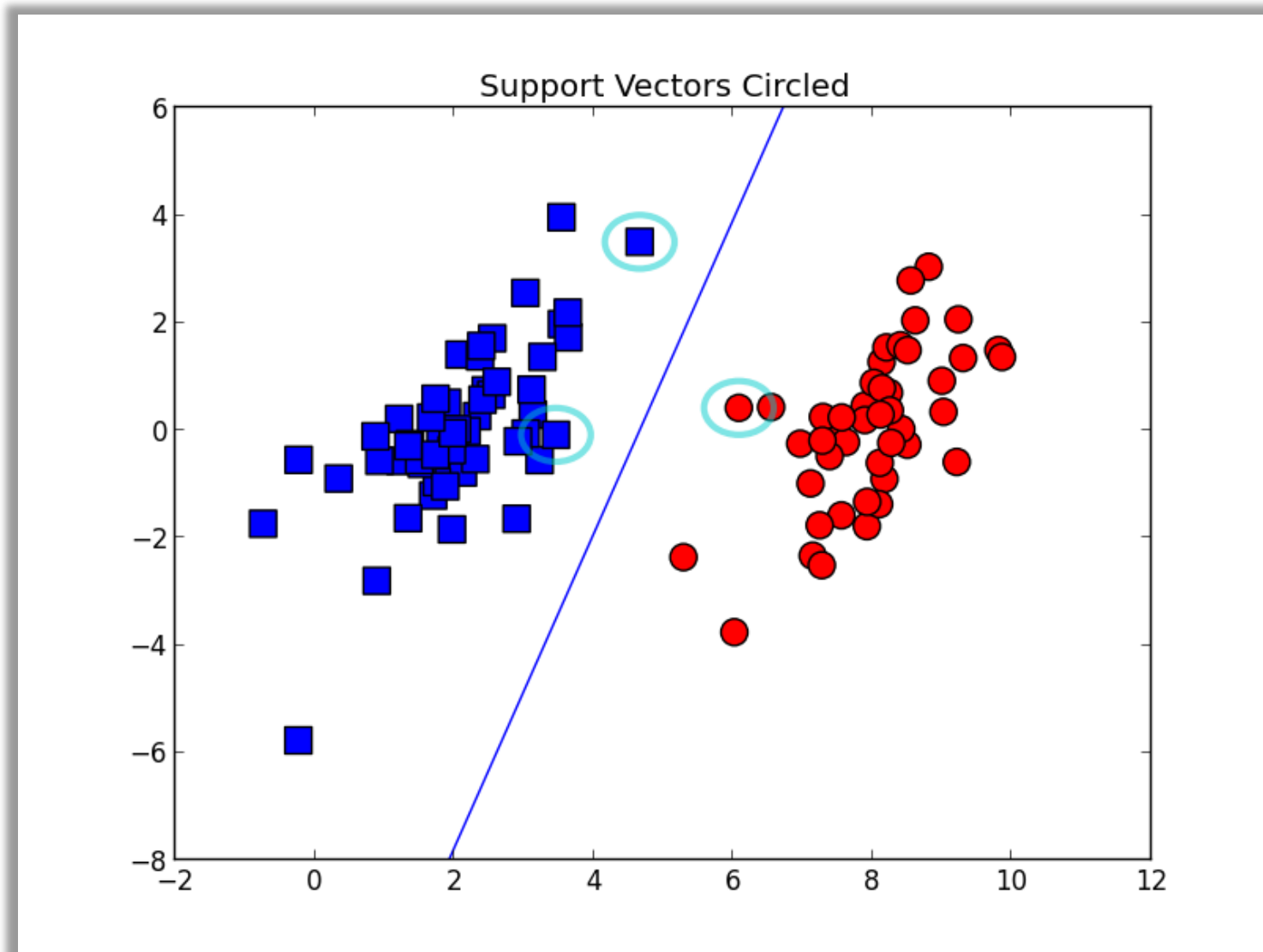
۱۲

□ تابع هدف.

$$\begin{aligned} & \min \frac{1}{2} \|\boldsymbol{\theta}\|^2 \\ \text{s.t.} \quad & (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) \geq +1 \quad \text{if } y^t = +1 \\ & (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) \leq -1 \quad \text{if } y^t = -1 \end{aligned}$$

□ ساده‌سازی.

$$\begin{aligned} & \min \frac{1}{2} \|\boldsymbol{\theta}\|^2 \\ \text{s.t.} \quad & y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) \geq +1 \end{aligned}$$



تابع هدف: بیان رسمی

۱۴

□ تابع هدف.

$$\begin{aligned} & \min \frac{1}{2} \|\boldsymbol{\theta}\|^2 \\ \text{s.t.} \quad & y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) \geq +1 \end{aligned}$$

← بهینه‌سازی مقرب

□ حل مسئله با استفاده از ضرایب لاگرانژ.

$$\begin{aligned} L_p &= \frac{1}{2} \|\boldsymbol{\theta}\|^2 - \sum_{t=1}^m \alpha^t [y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) - 1] \\ &= \frac{1}{2} \|\boldsymbol{\theta}\|^2 - \sum_{t=1}^m \alpha^t y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) + \sum_{t=1}^m \alpha^t \end{aligned}$$



یوزف لویی لاگرانژ
۱۷۳۶ - ۱۸۱۳

تابع هدف: بیان رسمی

۱۵

$$L_p = \frac{1}{2} \|\boldsymbol{\theta}\|^2 - \sum_{t=1}^m \alpha^t [y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) - 1]$$
$$= \frac{1}{2} \|\boldsymbol{\theta}\|^2 - \sum_{t=1}^m \alpha^t y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) + \sum_{t=1}^m \alpha^t$$

$$\frac{\partial L_p}{\partial \boldsymbol{\theta}} = 0 \Rightarrow \boldsymbol{\theta} = \sum_{t=1}^m \alpha^t y^t \mathbf{x}^t$$

$$\frac{\partial L_p}{\partial \theta_0} = 0 \Rightarrow \sum_{t=1}^m \alpha^t y^t = 0$$

مرز تصمیم‌گیری یک ترکیب خطی از داده‌های آموزشی

تابع هدف: بیان رسمی

۱۶

$$\begin{aligned}L_d &= \frac{1}{2}(\boldsymbol{\theta}^T \boldsymbol{\theta}) - \boldsymbol{\theta}^T \sum_{t=1}^m \alpha^t y^t \mathbf{x}^t - \theta_0 \sum_{t=1}^m \alpha^t y^t + \sum_{t=1}^m \alpha^t \\ &= -\frac{1}{2}(\boldsymbol{\theta}^T \boldsymbol{\theta}) + \sum_{t=1}^m \alpha^t \\ &= -\frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m \alpha^t \alpha^s y^t y^s (\mathbf{x}^t)^T \mathbf{x}^s + \sum_{t=1}^m \alpha^t\end{aligned}$$

الگوریتم
بهینه‌سازی ترتیبی مینیمال
پلت (۱۹۹۹)

subject to $\sum_{t=1}^m \alpha^t y^t = 0$ and $\alpha^t \geq 0 \forall t$

- مقدار بسیاری از آلفاها برابر با صفر است و تنها تعداد اندکی دارای مقدار بزرگ‌تر از صفر هستند؛
- α هایی که به ازای آنها مقدار آلفا بزرگ‌تر از صفر است، همان بردارهای پشتیبان هستند.

تابع هدف: شکل ساده شده

۱۷

$$L_d = -\frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m \alpha^t \alpha^s y^t y^s (\mathbf{x}^t)^T \mathbf{x}^s + \sum_{t=1}^m \alpha^t$$
$$= -\frac{1}{2} \alpha^T Q \alpha + e^T \alpha$$

الگوریتم
بهینه‌سازی ترتیبی مینیمال
پلت (۱۹۹۹)

$$Q_{ts} = y^t y^s (\mathbf{x}^t)^T \mathbf{x}^s, \quad e = [1 \ 1 \ \dots \ 1]^T \in \mathbb{R}^m$$

subject to $\sum_{t=1}^m \alpha^t y^t = 0$ and $\alpha^t \geq 0 \ \forall t$

- مقدار بسیاری از آلفاها برابر با صفر است و تنها تعداد اندکی دارای مقدار بزرگ‌تر از صفر هستند؛
- x هایی که به ازای آنها مقدار آلفا بزرگ‌تر از صفر است، همان بردارهای پشتیبان هستند.

داده‌های تفکیک ناپذیر خطی: ماشیه نرم

□ س. اگر داده‌ها به صورت خطی تفکیک پذیر نباشند چه می‌شود؟

$$y^t(\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) \geq 1 - \varepsilon^t$$

□ حاشیه نرم. اجازه دادن اندکی خطا در جداسازی!

$$\text{soft error} = \sum_{t=1}^m \varepsilon^t$$

□ خطای نرم.

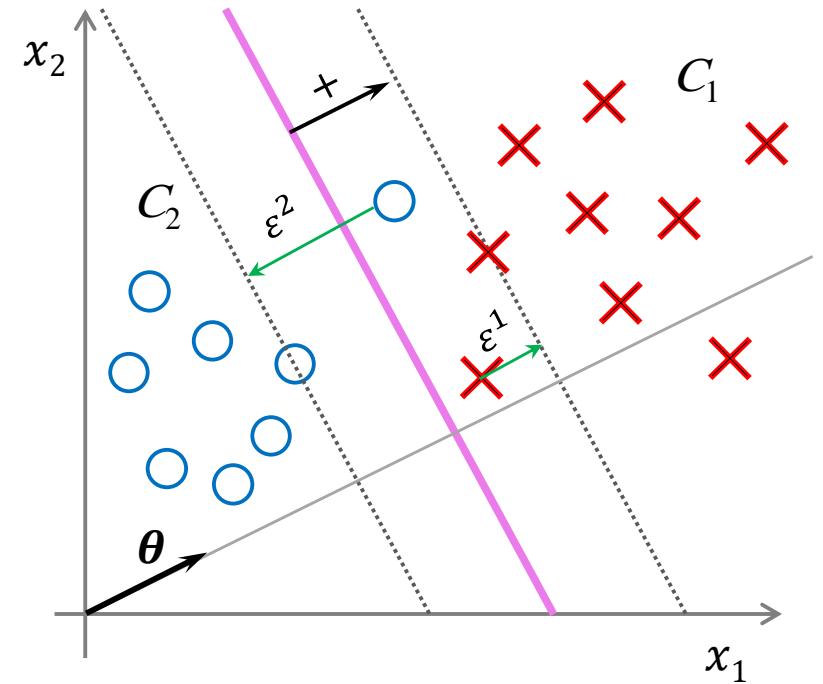
ضریب جریمه

□ تابع هدف جدید.

$$\begin{aligned} & \min \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{t=1}^m \varepsilon^t \\ \text{s.t.} \quad & y^t(\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) \geq 1 - \varepsilon^t \\ & \varepsilon^t \geq 0 \end{aligned}$$

داده‌های تفکیک‌ناپذیر خطی: ماشین بردار پشتیبان

$$\begin{aligned} & \min \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{t=1}^m \varepsilon^t \\ \text{s.t.} \quad & y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) \geq 1 - \varepsilon^t \\ & \varepsilon^t \geq 0 \end{aligned}$$



ضرایب لاگرانژ

$$L_p = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t [y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) - 1 + \varepsilon^t] - \sum_{t=1}^m \mu^t \varepsilon^t$$

داده‌های تفکیک‌ناپذیر خطی: ماشین بردار پشتیبان

۲۰

$$L_p = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t [y^t (\boldsymbol{\theta}^T \mathbf{x}^t + \theta_0) - 1 + \varepsilon^t] - \sum_{t=1}^m \mu^t \varepsilon^t$$

$$\frac{\partial L_p}{\partial \boldsymbol{\theta}} = 0 \Rightarrow \boldsymbol{\theta} = \sum_{t=1}^m \alpha^t y^t \mathbf{x}^t$$

$$\frac{\partial L_p}{\partial \theta_0} = 0 \Rightarrow \sum_{t=1}^m \alpha^t y^t = 0$$

$$\frac{\partial L_p}{\partial \varepsilon^t} = 0 \Rightarrow C - \alpha^t - \mu^t = 0 \Rightarrow 0 \leq \alpha^t \leq C$$

تابع هدف: دوگان

۲۱

$$L_d = -\frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m \alpha^t \alpha^s y^t y^s (\mathbf{x}^t)^T \mathbf{x}^s + \sum_{t=1}^m \alpha^t$$
$$= -\frac{1}{2} \alpha^T Q \alpha + e^T \alpha$$

الگوریتم
بهبودسازی ترتیبی مینیمال
پلت (۱۹۹۹)

subject to $\sum_{t=1}^m \alpha^t y^t = 0$ and $0 \leq \alpha^t \leq C \forall t$

□ تخمین خطا بر اساس تعداد بردارهای پشتیبان. [اوپنیک، ۱۹۹۵]

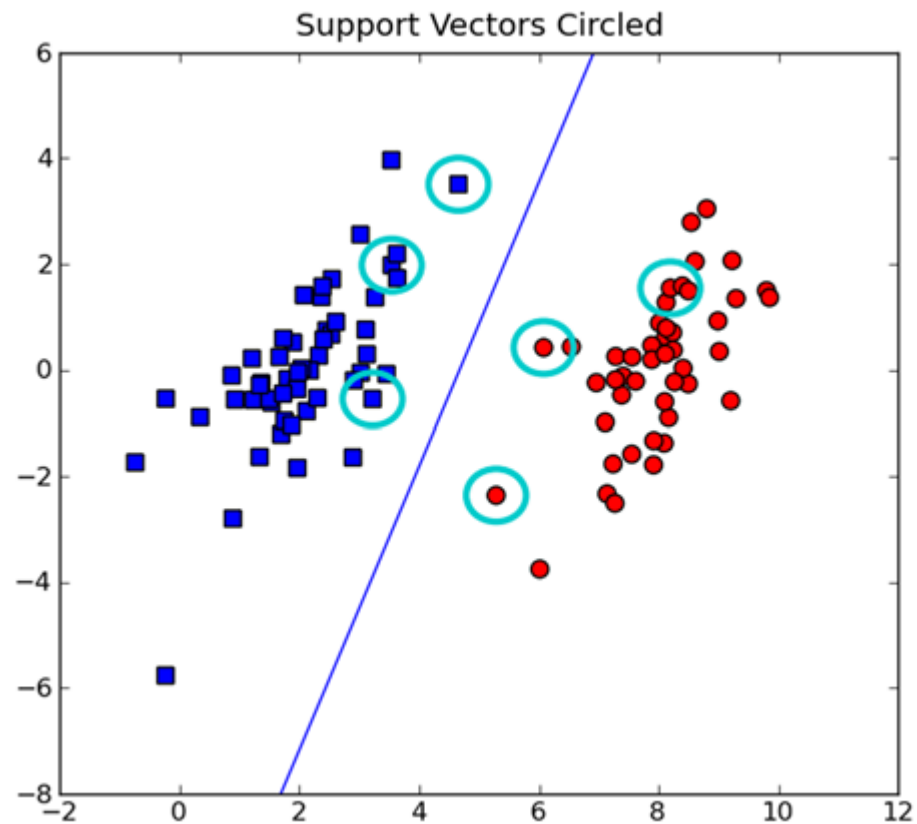
□ مقدار بسیاری از آلفاها برابر با صفر است و تنها تعداد اندکی دارای مقدار بزرگتر از صفر هستند؛

□ x هایی که به ازای آنها مقدار آلفا بزرگتر از صفر است، همان بردارهای پشتیبان هستند.

$$E_m[P(\text{error})] \leq \frac{E_m[\text{\#of support vectors}]}{m}$$

راه‌حل ماشیه نرم

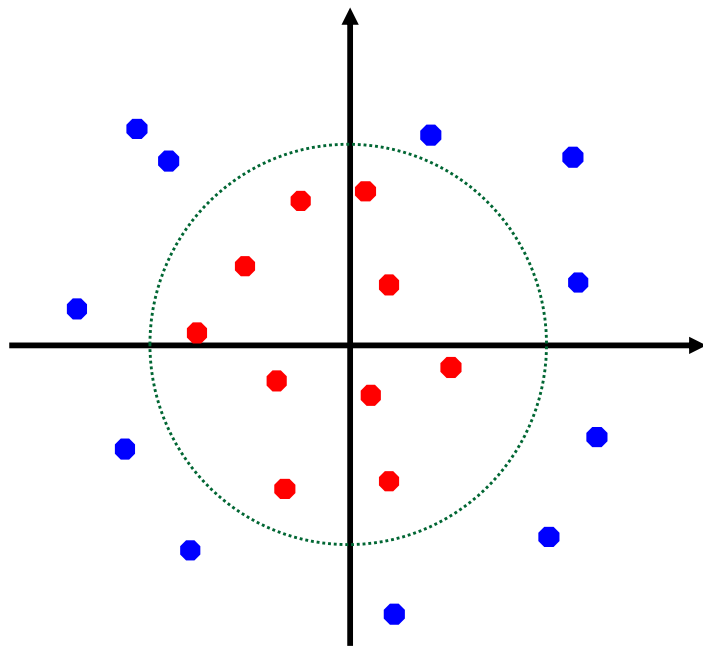
□ بردارهای پشتیبان.



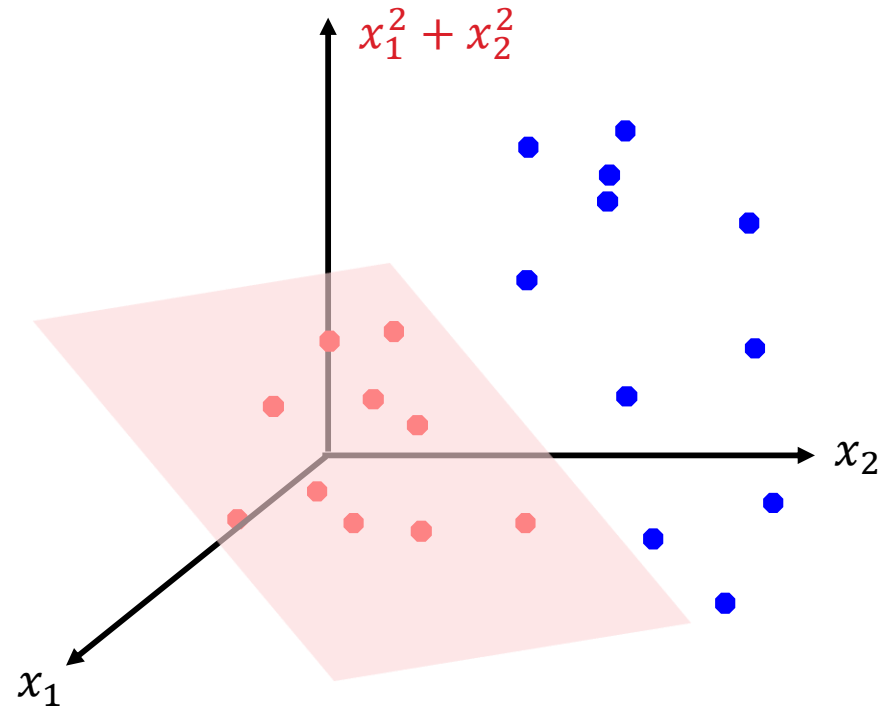
ترفند کرنل

توابع کرنل

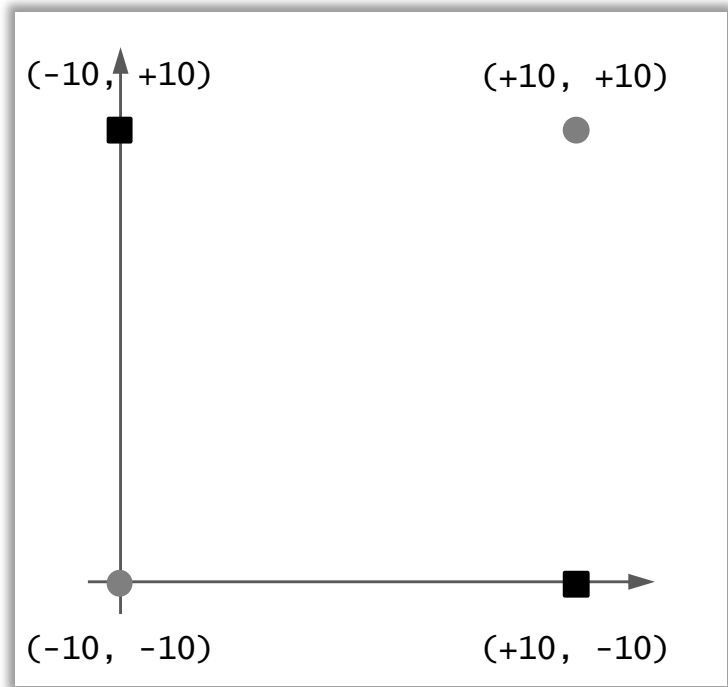
- ایده. نگاشت مسئله به یک فضای ویژگی جدید با استفاده از تبدیلات غیرخطی.
- استفاده از یک مدل خطی در فضای جدید به منظور کلاس‌بندی داده‌ها.
- مدل خطی در فضای جدید متناظر با یک مدل غیرخطی در فضای اصلی است.



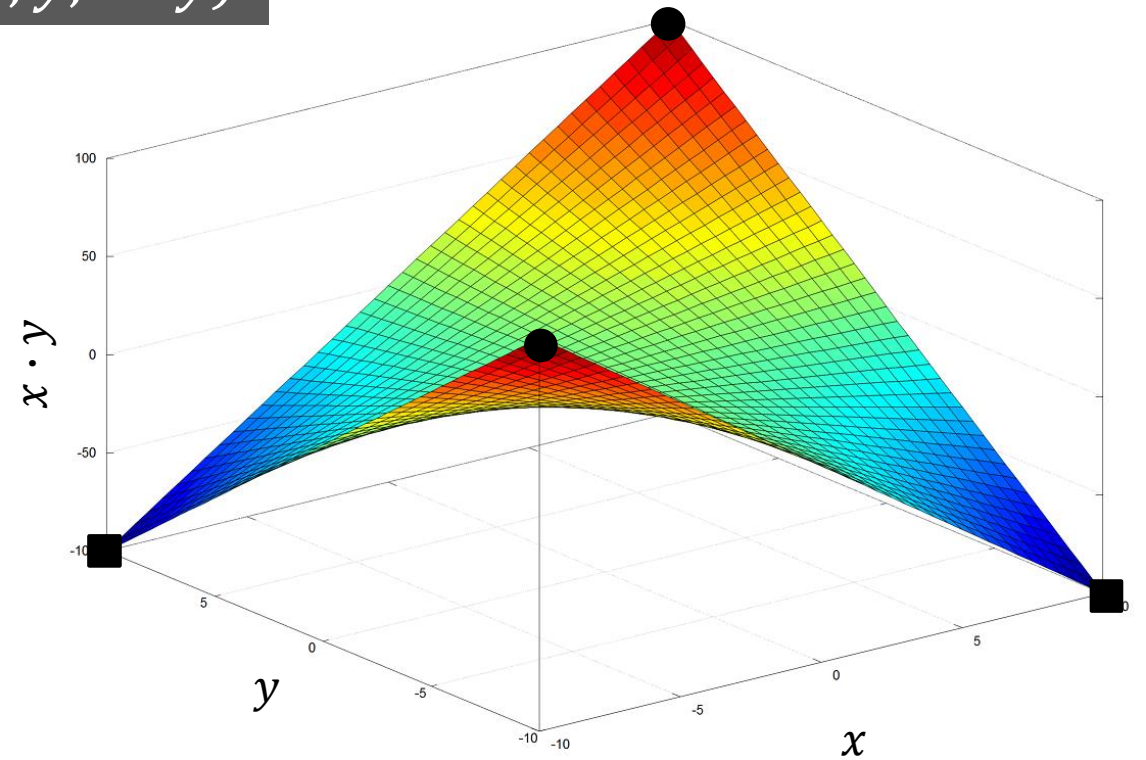
$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$



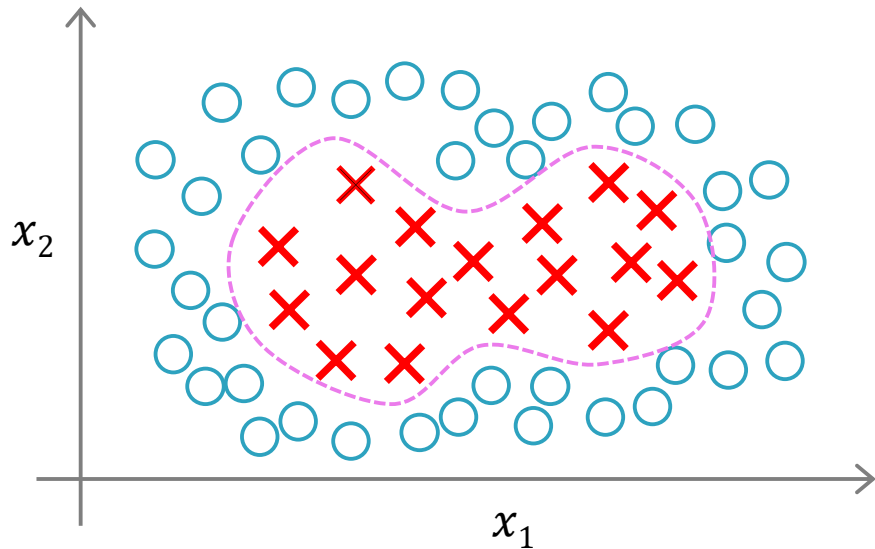
مثال: مسئله XOR



$$(x, y) \rightarrow (x, y, x \cdot y)$$



مرزهای تصمیم‌گیری غیرخطی



پیش‌بینی. $y = 1$ اگر:

$$\begin{aligned} h_{\theta}(x) &= \theta_0 \\ &+ \theta_1 x_1 + \theta_2 x_2 \\ &+ \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 \\ &+ \dots \\ &\geq 0 \end{aligned}$$

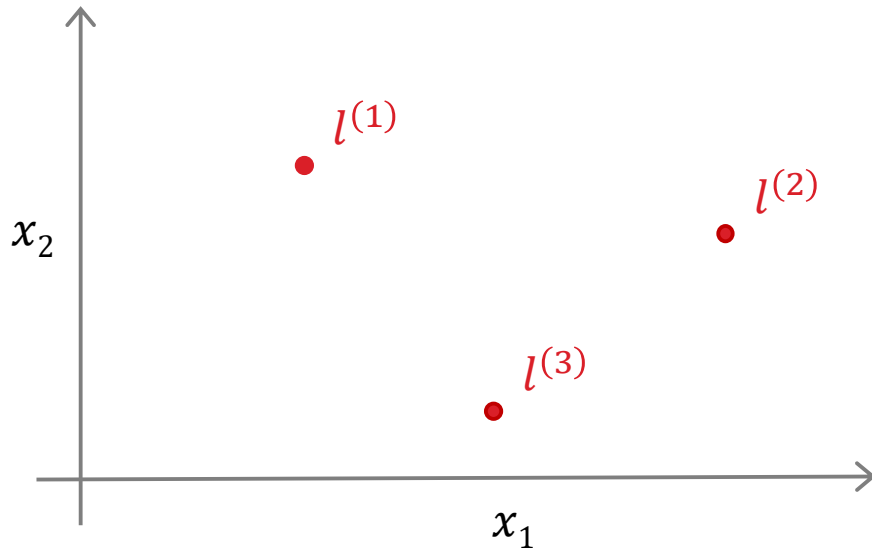
ویژگی‌ها.

$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1^2, \quad f_4 = x_2^2, \quad f_5 = x_1 x_2, \quad \dots$$

$$h_{\theta}(f) = \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 + \theta_5 f_5 + \dots \leftarrow \text{مرز تصمیم‌گیری فطی}$$

س. آیا روش بهتری برای انتخاب ویژگی‌های جدید f_1 ، f_2 و ... وجود دارد؟

□ ایده. با داشتن x ، مجموعه جدید ویژگی‌ها را بر اساس **شباهت** آن با نقاط راهنمای $l^{(1)}$ ، $l^{(2)}$ و $l^{(3)}$ انتخاب کن.



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

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

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

کرنل (کرنل گوسی)

□ تابع کرنل. معیاری به منظور محاسبه شباهت میان داده‌های x و y

کرنل‌ها به عنوان معیار شباهت

□ تابع کرنل.

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

□ حالت اول. $x \approx l^{(i)}$

$$f_i \approx \exp\left(-\frac{0}{2\sigma^2}\right) = \exp(0) = 1$$

□ حالت دوم. x بسیار دور از $l^{(i)}$

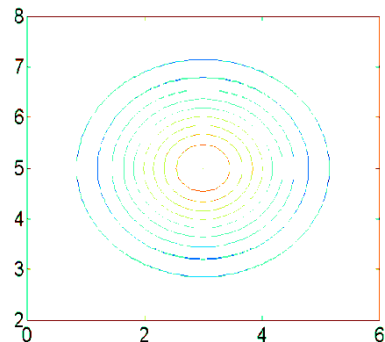
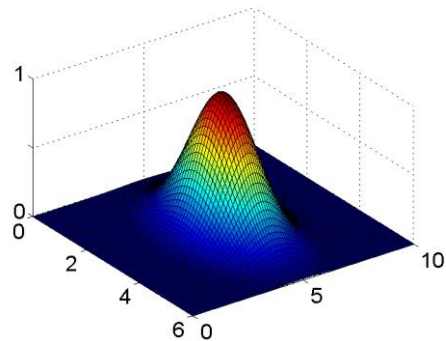
$$f_i \approx \exp\left(-\frac{\infty}{2\sigma^2}\right) = \exp(-\infty) = 0$$

کرنل‌ها به عنوان معیار شباهت

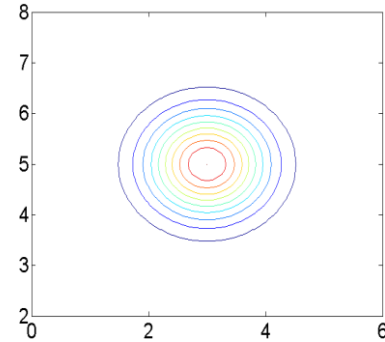
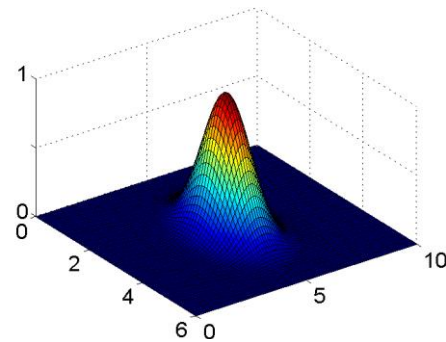
□ مثال.

$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

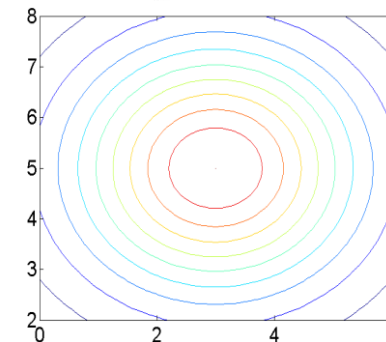
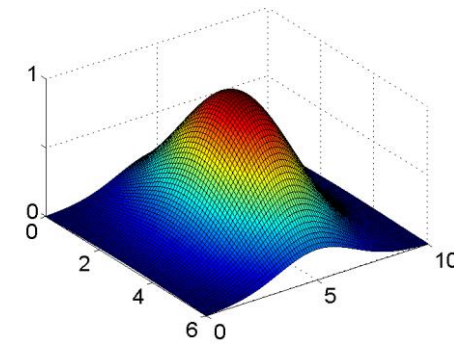
$$\sigma^2 = 1$$



$$\sigma^2 = 0.5$$



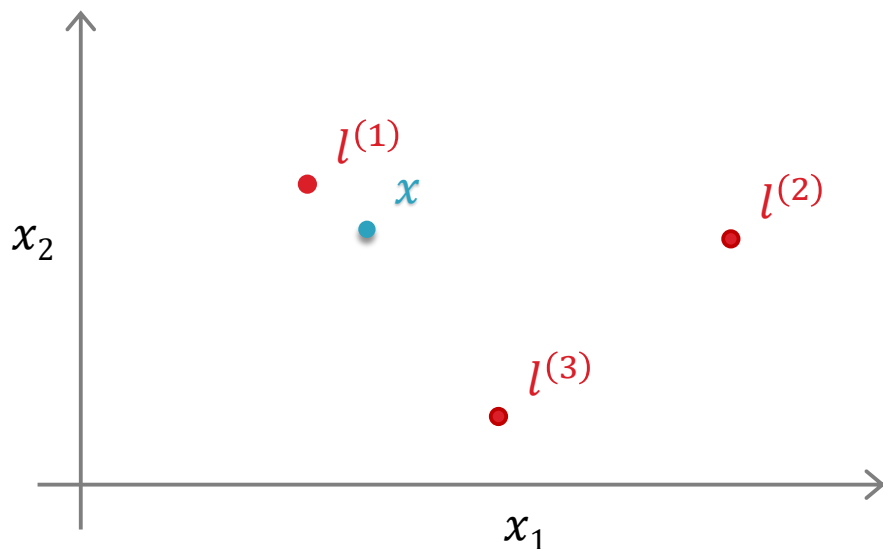
$$\sigma^2 = 3$$



کرنل‌ها به عنوان معیار شباهت

۳۰

پیش‌بینی. $y = 1$ اگر:



$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

$$\begin{array}{cccc} \uparrow & \uparrow & \uparrow & \uparrow \\ -0.5 & 1.0 & 1.0 & 0.0 \end{array}$$

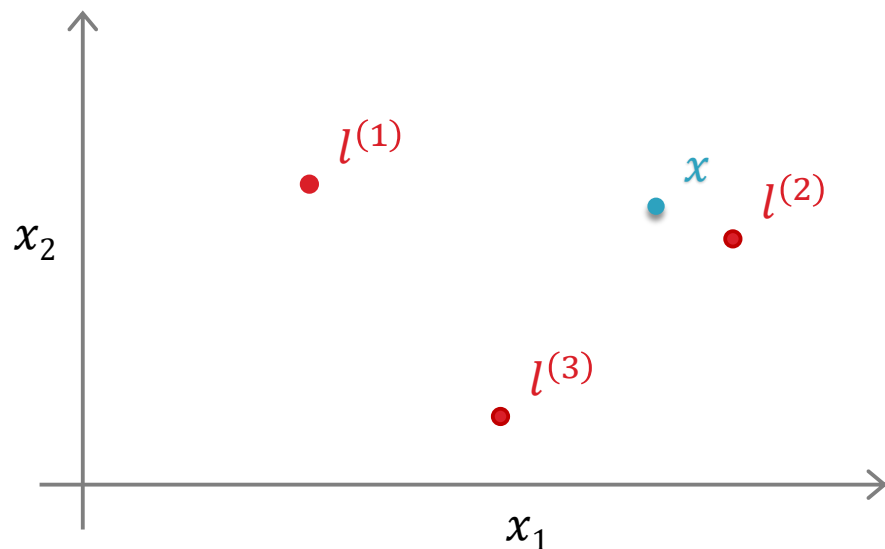
$$f_1 \approx 1, f_2 \approx f_3 \approx 0$$

$$h_{\theta}(f) \approx -0.5 + (1.0)(1.0) + (1.0)(0.0) + (0.0)(0.0) = 0.5 \geq 0 \Rightarrow y = 1$$

کرنل‌ها به عنوان معیار شباهت

۳۱

پیش‌بینی. $y = 1$ اگر:



$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

$$\begin{array}{cccc} \uparrow & \uparrow & \uparrow & \uparrow \\ -0.5 & 1.0 & 1.0 & 0.0 \end{array}$$

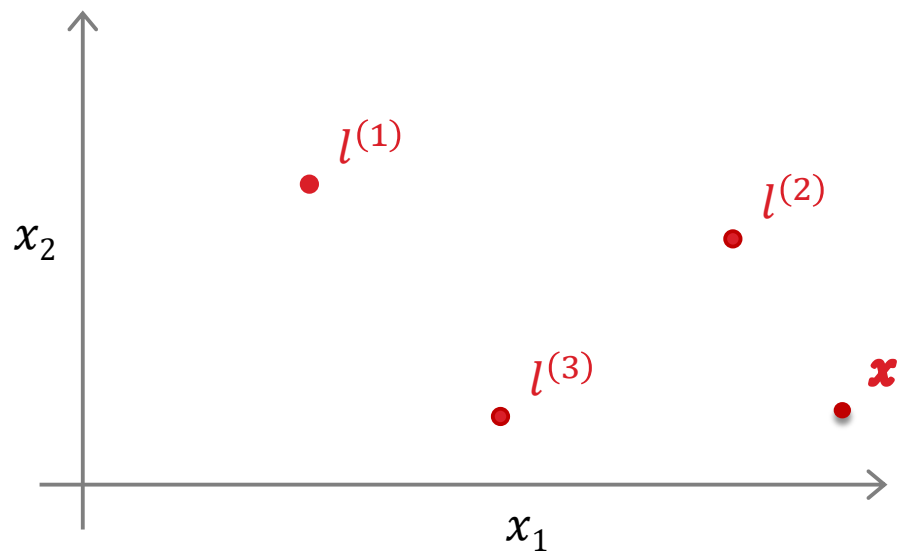
$$f_1 \approx f_3 \approx 0, f_2 \approx 1$$

$$h_{\theta}(f) \approx -0.5 + (1.0)(0.0) + (1.0)(1.0) + (0.0)(0.0) = 0.5 \geq 0 \Rightarrow y = 1$$

کرنل‌ها به عنوان معیار شباهت

۳۲

پیش‌بینی. $y = 1$ اگر:



$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

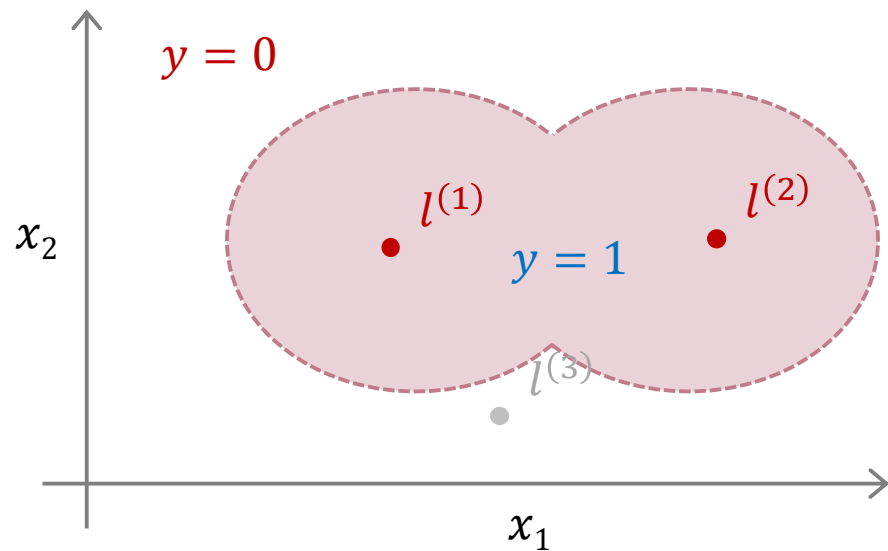
$$\begin{array}{cccc} \uparrow & \uparrow & \uparrow & \uparrow \\ -0.5 & 1.0 & 1.0 & 0.0 \end{array}$$

$$f_1 \approx f_3 \approx f_2 \approx 0$$

$$h_{\theta}(f) \approx -0.5 + (1.0)(0.0) + (1.0)(0.0) + (0.0)(0.0) = -0.5 \leq 0 \Rightarrow y = 0$$

کرنل ها به عنوان معیار شباهت

پیش بینی. $y = 1$ اگر:



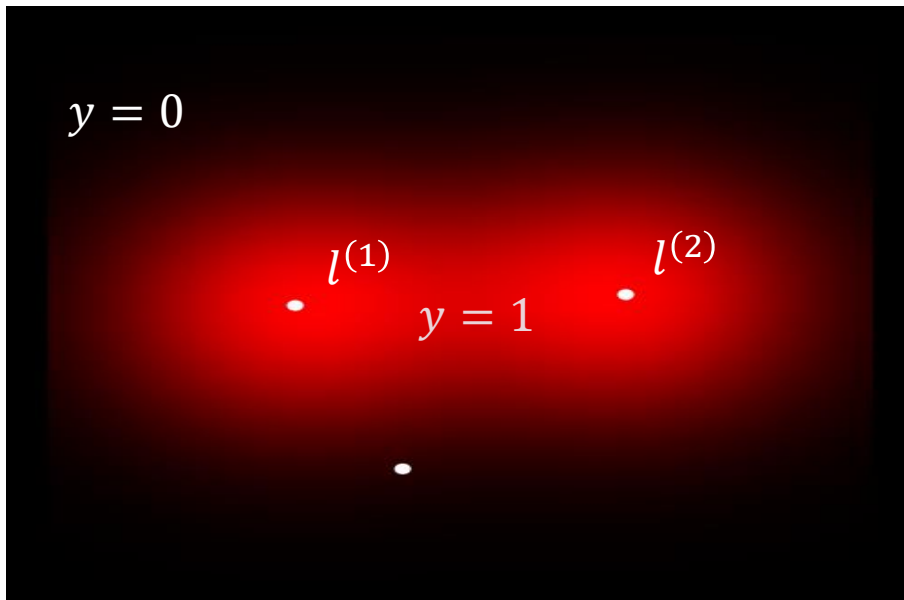
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

\uparrow \uparrow \uparrow \uparrow
-0.5 1.0 1.0 0.0

□ مرز تصمیم گیری. نقاط نزدیک به $l^{(1)}$ و $l^{(2)}$ را در کلاس ۱ و سایر نقاط را به کلاس صفر کلاس بندی می کند.

کرنل ها به عنوان معیار شباهت

۳۴



پیش بینی. $y = 1$ اگر:

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

↑ ↑ ↑ ↑
-0.5 1.0 1.0 0.0

□ مرز تصمیم گیری. نقاط نزدیک به $l^{(1)}$ و $l^{(2)}$ را در کلاس ۱ و سایر نقاط را به کلاس صفر کلاس بندی می کند.

جزئیات باقیمانده

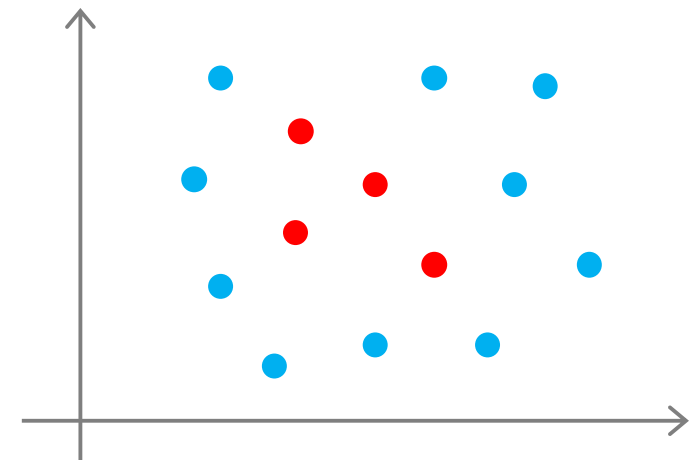
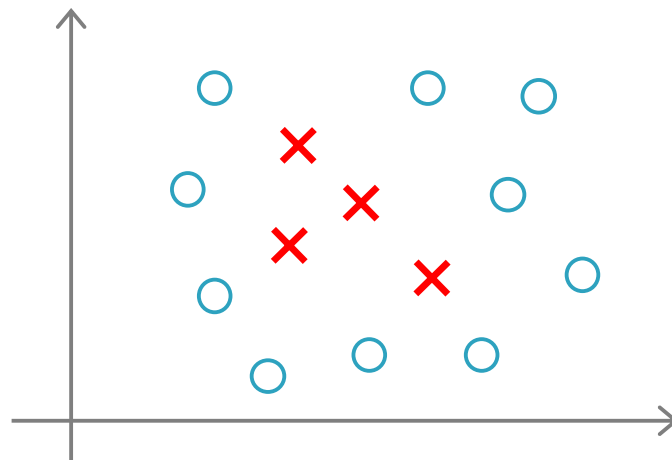
□ س. الگوریتم یادگیری نقاط راهنما را چگونه به صورت خودکار انتخاب می کند؟

□ س. مقدار مناسب برای پارامترهای تابع کرنل چگونه تعیین می شوند؟

□ س. آیا انواع دیگری از کرنل ها وجود دارد؟

انتخاب نقاط راهنما

- س. الگوریتم یادگیری نقاط راهنما را چگونه به صورت خودکار انتخاب می کند؟
- به ازای هر نمونه در مجموعه آموزشی، یک نقطه راهنما مساوی با آن نمونه انتخاب می شود.



نگاشت ویژگی‌ها

□ مجموعه آموزشی.

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

□ نقاط راهنما.

$$l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$$

□ نگاشت فضای ویژگی.

$$x = \begin{bmatrix} x_0 = 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$



$$f = \begin{bmatrix} f_0 = 1 \\ f_1 = K(x, x^{(1)}) \\ f_2 = K(x, x^{(2)}) \\ \vdots \\ f_m = K(x, x^{(m)}) \end{bmatrix}$$

ترفند کرنل

□ تابع کرنل. پیش پردازش داده x با استفاده از توابع کرنل:

$$\begin{aligned} \mathbf{z} &= \varphi(\mathbf{x}) \\ &= (\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x}), \dots, \varphi_k(\mathbf{x})) \end{aligned}$$

ممکن است بینهایت باشد!!!

$$g(\mathbf{z}) = \boldsymbol{\theta}^T \mathbf{z} + \theta_0$$

$$g(\mathbf{x}) = \boldsymbol{\theta}^T \varphi(\mathbf{x}) + \theta_0$$

□ راه حل SVM.

$$\boldsymbol{\theta} = \sum_{t=1}^m \alpha^t y^t \mathbf{z}^t = \sum_{t=1}^m \alpha^t y^t \varphi(\mathbf{x}^t)$$

$$g(\mathbf{x}) = \boldsymbol{\theta}^T \varphi(\mathbf{x}) + \theta_0 = \left(\sum_{t=1}^m \alpha^t y^t \varphi(\mathbf{x}^t)^T \right) \varphi(\mathbf{x}) + \theta_0 = \left(\sum_{t=1}^m \alpha^t y^t \varphi(\mathbf{x}^t)^T \varphi(\mathbf{x}) \right) + \theta_0$$

□ کلاس بندی داده جدید.

$$g(\mathbf{x}) = \left(\sum_{t=1}^m \alpha^t y^t k(\mathbf{x}^t, \mathbf{x}) \right) + \theta_0$$

داده آموزشی داده جدید

← مرز تصمیم گیری

$$L_p = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{t=1}^m \varepsilon^t$$

s.t. $y^t \boldsymbol{\theta}^T \boldsymbol{\varphi}(\mathbf{x}^t) \geq 1 - \varepsilon^t$

$$\varepsilon^t \geq 0$$

$$L_p = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t [y^t \boldsymbol{\theta}^T \boldsymbol{\varphi}(\mathbf{x}^t) - 1 + \varepsilon^t] - \sum_{t=1}^m \mu^t \varepsilon^t$$

ضرایب لاگرانژ

ضرایب لاگرانژ

توابع کرنل: مسئله اصلی

۴۰

$$L_p = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t [y^t \boldsymbol{\theta}^T \varphi(\mathbf{x}^t) - 1 + \varepsilon^t] - \sum_{t=1}^m \mu^t \varepsilon^t$$

$$\frac{\partial L_p}{\partial \boldsymbol{\theta}} = 0 \Rightarrow \boldsymbol{\theta} = \sum_{t=1}^m \alpha^t y^t \varphi(\mathbf{x}^t)$$

$$\frac{\partial L_p}{\partial \varepsilon^t} = 0 \Rightarrow C - \alpha^t - \mu^t = 0 \Rightarrow 0 \leq \alpha^t \leq C$$

توابع کرنل: مسئله دوگان

۴۱

$$L_d = -\frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m \alpha^t \alpha^s y^t y^s \varphi(\mathbf{x}^t)^T \varphi(\mathbf{x}^s) + \sum_{t=1}^m \alpha^t$$

subject to $\sum_{t=1}^m \alpha^t y^t = 0$ and $0 \leq \alpha^t \leq C \forall t$

□ ایده ماشین‌های کرنل. [ترفند کرنل]

□ جایگزینی حاصل ضرب داخلی توابع پایه با یک تابع کرنل به صورت $K(\mathbf{x}^t, \mathbf{x}^s)$

$$L_d = -\frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m \alpha^t \alpha^s y^t y^s K(\mathbf{x}^t, \mathbf{x}^s) + \sum_{t=1}^m \alpha^t$$

ماتریس K : یک ماتریس متقارن و مثبت معین (برای تفکیک‌پذیری فطی)

توابع کرنل: چند جمله‌ای

□ کرنل چند جمله‌ای. یک چند جمله‌ای از درجه‌ی q .

$$K(x^t, x) = (x^T x^t + 1)^q$$

□ مثال. $[q = 2, d = 2]$

$$\begin{aligned} K(x, y) &= (x^T y + 1)^2 \\ &= (x_1 y_1 + x_2 y_2 + 1)^2 \\ &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2 \end{aligned}$$

۳ ضرب، ۲ جمع

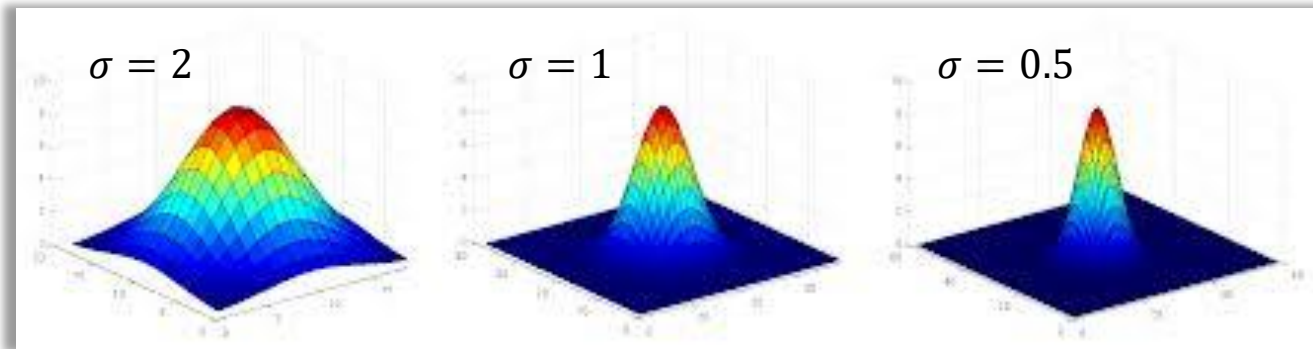
$$\varphi(x) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2]^T$$

۶ ضرب، ۵ جمع

$$\varphi(y) = [1, \sqrt{2}y_1, \sqrt{2}y_2, \sqrt{2}y_1 y_2, y_1^2, y_2^2]^T$$

تابع کرنل گوسی

۴۳



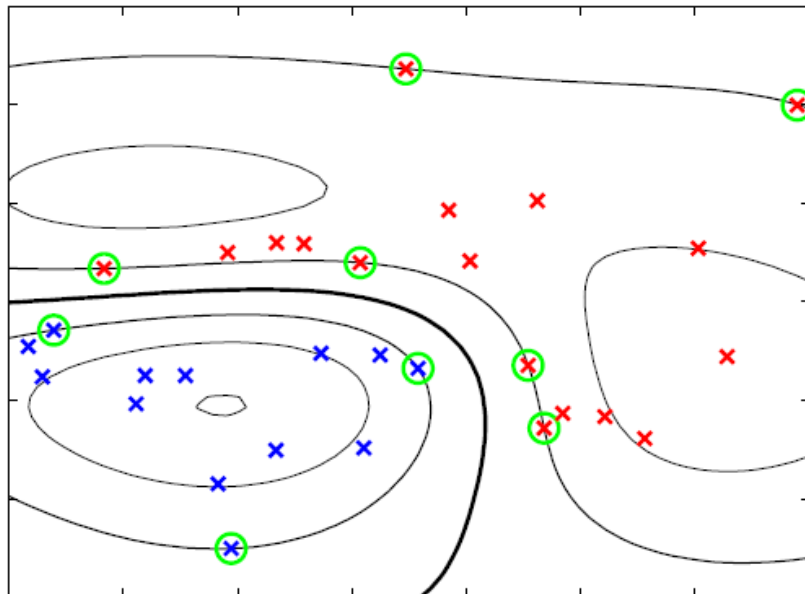
□ تابع کرنل گوسی.

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

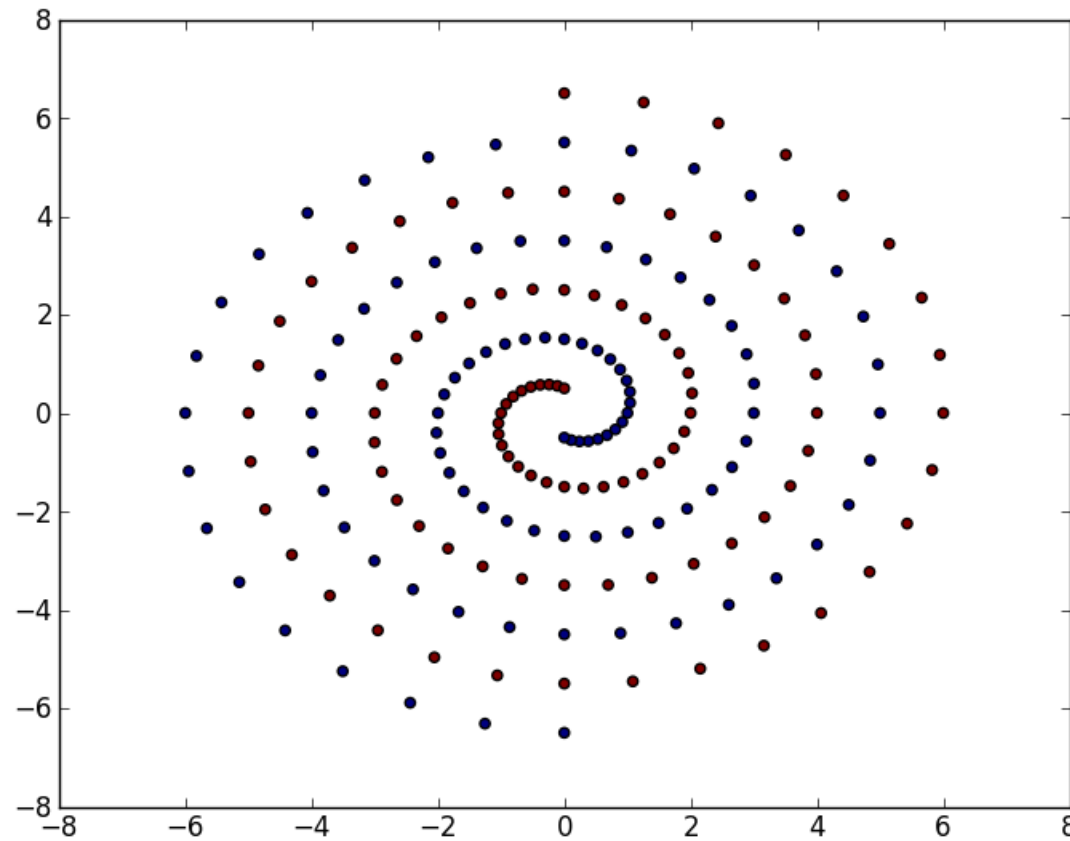
□ یافتن یک مقدار مناسب برای σ .

□ با استفاده از مجموعه اعتبارسنجی [انتخاب مدل]

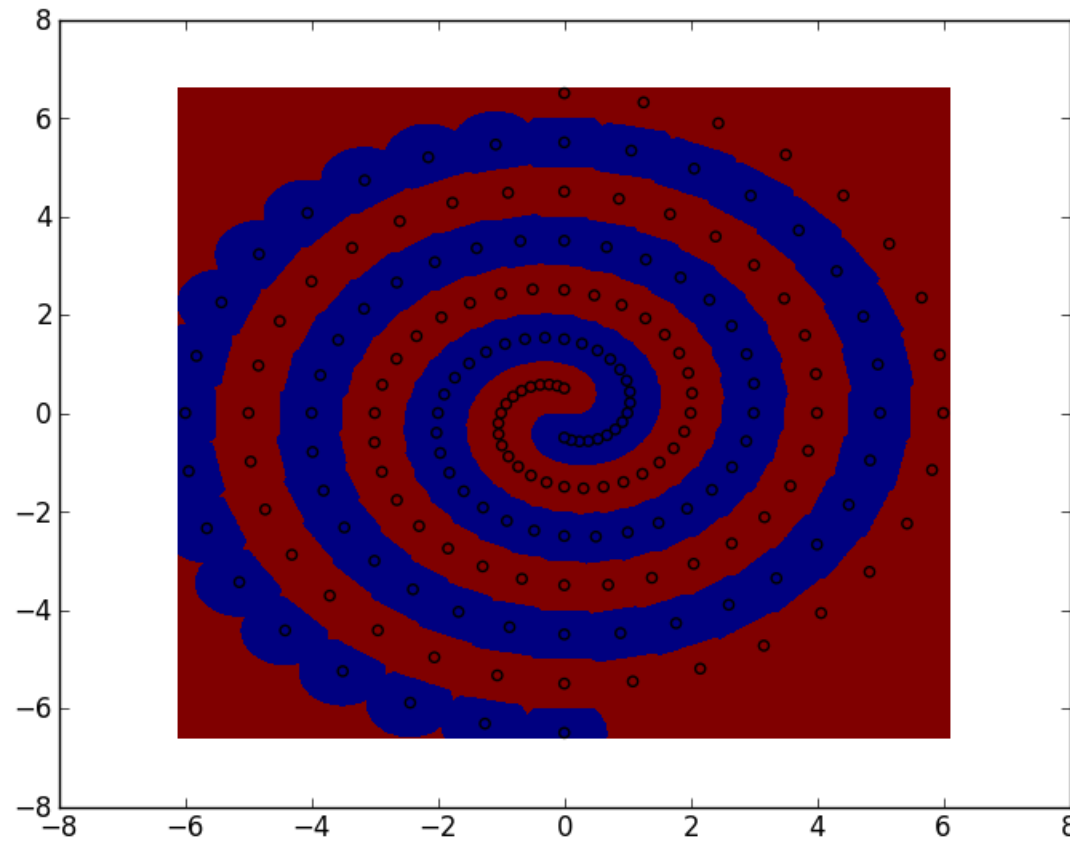
□ مقادیر بزرگ تر: مرز تصمیم گیری هموارتر



مثال: تابع کرنل گوسی

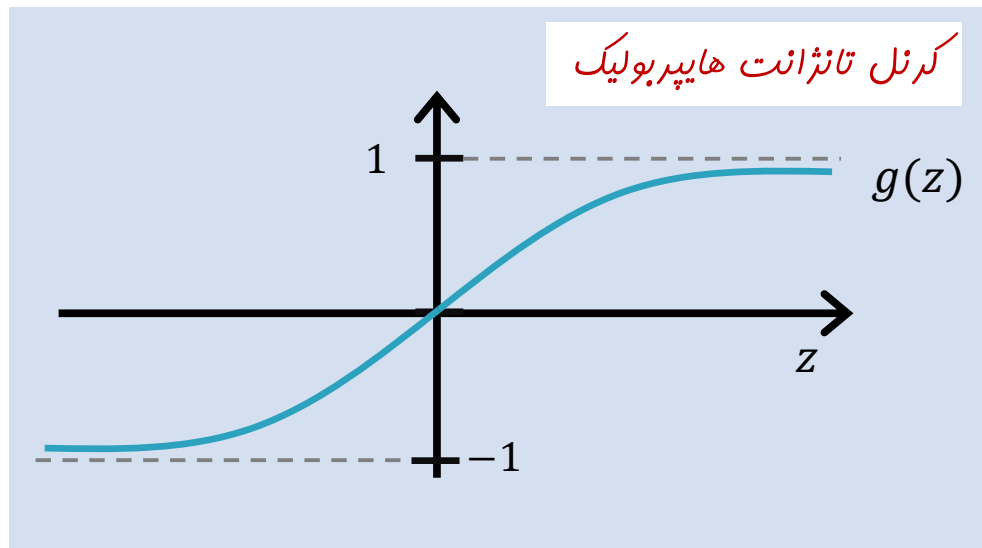


مثال: تابع کرنل گوسی



توابع کرنل: انواع دیگر

□ تابع کرنل. معیاری به منظور محاسبه شباهت میان داده‌های x و y



$$K(x^t, x) = \tanh(2x^T x^t + 1)$$

□ انواع دیگر.

□ کرنل تانژانت هایپربولیک

□ کرنل رشته‌ای

□ کرنل درختی

□ کرنل گرافی

□ ...

پارامترهای SVM

□ س. مقدار مناسب برای پارامترهای تابع کرنل چگونه تعیین می‌شوند؟

□ پارامتر C .

- مقادیر کوچک‌تر: بایاس بیشتر، واریانس کمتر
- مقادیر بزرگ‌تر: بایاس کمتر، واریانس بیشتر

□ پارامتر σ .

- مقادیر کوچک‌تر: بایاس کمتر، واریانس بیشتر
- مقادیر بزرگ‌تر: بایاس بیشتر، واریانس کمتر

تعیین مقادیر هر دو پارامتر:
بسته‌بندی شبکه‌بندی

کلاس بندی چندکلاسی

$$\arg \max_i g_i(x)$$

□ روش اول. یکی در برابر همه [روش ترجیحی]

□ آموزش: آموزش k ماشین بردار پشتیبان یکی به ازای هر کلاس

□ آزمایش: محاسبه $g_i(x)$ به ازای $0 \leq i \leq k$ و انتخاب بزرگترین مقدار

□ روش دوم. جداسازی دو به دو

□ آموزش $k(k-1)/2$ ماشین بردار پشتیبان به طوری که $g_{ij}(x)$ نمونه های دو کلاس C_i و C_j را از هم جدا می کند.

□ ساده تر و سریع تر

□ روش سوم. حل یک مسئله بهینه سازی چند کلاسی

$$\min \frac{1}{2} \sum_{i=1}^k \|\theta_i\|^2 + C \sum_i \sum_t \varepsilon_i^t$$

$$\text{s.t.} \quad \theta_{z^t} x^t + \theta_{z^t 0} \geq \theta_i x^t + \theta_{i0} + 2 - \varepsilon_i^t, \forall i \neq z^t$$

راهنمای استفاده از SVM

□ پیاده‌سازی. استفاده از بسته‌های نرم‌افزاری موجود مانند LIBSVM و SVM^{light}

□ تعیین تابع کرنل.

□ کرنل خطی (عدم استفاده از کرنل): وقتی که $n \gg m$

□ گوسی، چند جمله‌ای، رشته‌ای و ...

□ تعیین مقدار پارامترها. جستجوی شبکه‌بندی

□ انتخاب مقدار برای پارامتر C

□ انتخاب مقدار برای پارامترهای تابع کرنل (مانند σ)

		σ					
		err	.01	.1	1	10	100
C	.01						
	.1						
	1						
	10						
	100						
	err						

SVM، رگرسیون لجستیک یا شبکه عصبی؟

۵۰

□ حالت ۱. $[n \gg m]$

□ مثال: تشخیص هرزنامه (۱۰۰۰ نمونه آموزشی، ۵۰۰۰۰ ویژگی)

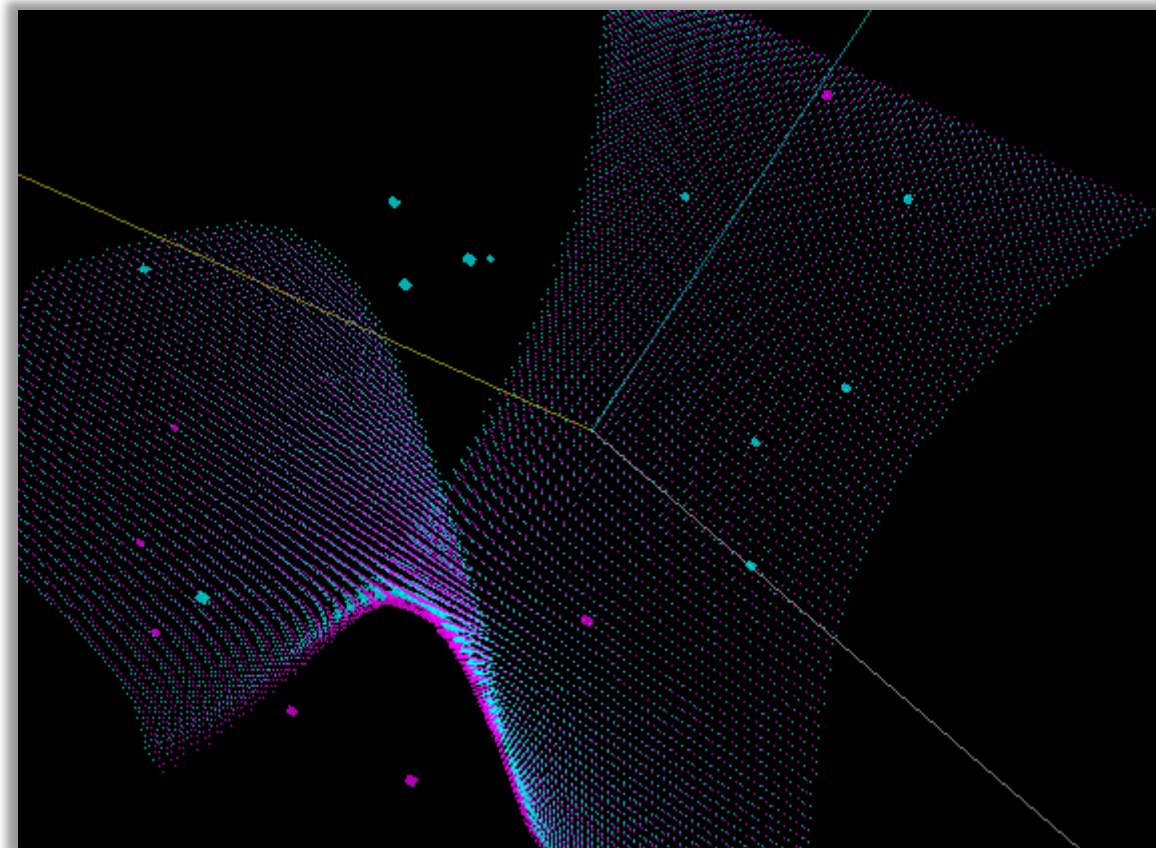
□ رگرسیون لجستیک یا SVM خطی

□ حالت ۲. [تعداد ویژگی‌ها کم، تعداد نمونه‌های آموزشی زیاد]

□ SVM با کرنل گوسی

□ توجه. شبکه‌های عصبی در تمامی حالت‌های فوق قابل استفاده هستند، اما ممکن است به زمان بیشتری برای آموزش نیاز داشته باشند.

- SVM toy 3D. [<http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/svmtoy3d/>]



مطالب بیشتر در مورد کرنل‌ها

- س. چگونه می‌دانیم استفاده از کرنل‌ها در جداسازی داده‌ها به ما کمک می‌کند؟
- در فضای n -بعدی، هر مجموعه از n بردار مستقل، به صورت خطی تفکیک‌پذیر هستند.
- اگر ماتریس K یک ماتریس مثبت معین باشد، آنگاه داده‌ها به صورت خطی تفکیک‌پذیر هستند.

□ قضیه. ماتریس K یک ماتریس مثبت معین است، زیرا $K = L^T L$

□ ستون i در ماتریس L برابر است با بردار $\phi(x^{(i)})$

□ اثبات. بردار غیر صفر v را در نظر بگیرید. در این صورت:

$$v^T K v = v^T L^T L v = (L v)^T (L v) = w^T w = \|w\|^2 \geq 0$$

□ و چون L و v هر دو مخالف صفر هستند، بردار w نیز مخالف صفر است. یعنی:

$$\|w\|^2 > 0 \Rightarrow v^T K v > 0 \Rightarrow K \text{ is positive definite}$$