

ماشین بردار پشتیبان (SVM)

(Support Vector Machine)

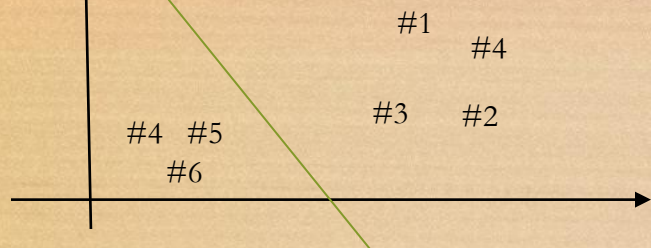
فهرست مطالب

- ❖ مقدمه
- ❖ مسئله جداسازی خطی (Linear Discrimination)
- ❖ مسئله جداسازی غیر خطی
 - بهینه سازی با حاشیه های انعطاف پذیر (نرم)
 - بکارگیری kernel trick به منظور خطی سازی مسائل غیر خطی مورد نظر
- ❖ ماشین های بردار پشتیبانی و شبکه های عصبی
- ❖ نقاط ضعف ماشین های بردار پشتیبانی

ماشین بردار پشتیبان (SVM) یک الگوریتم نظارت شده یادگیری ماشین است که هم برای مسائل طبقه بندی و هم رگرسیون قابل استفاده است، با این حال از آن بیشتر در مسائل طبقه بندی استفاده می شود. در الگوریتم SVM هر نمونه داده را به عنوان یک نقطه در فضای n بعدی روی نمودار پراکندگی داده ها ترسیم کرده و مقدار هر ویژگی مربوط به داده ها، یکی از مولفه های مختصات نقطه روی نمودار را مشخص می کند. سپس با ترسیم یک خط راست، داده های مختلف و متمایز از یکدیگر را دسته بندی می کند.

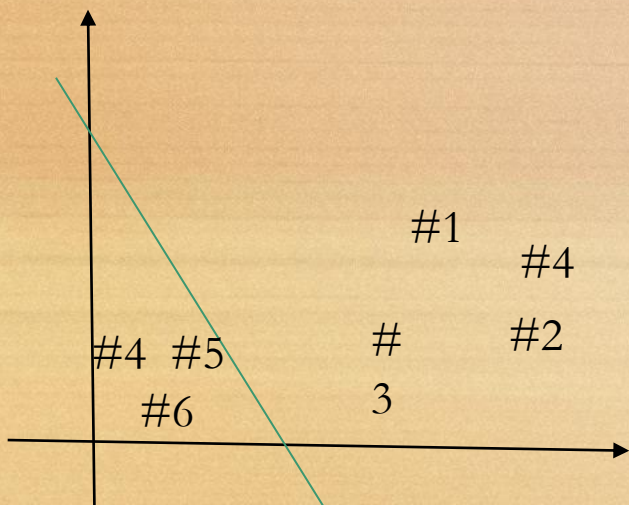
مقدمه

- پرسپترون در واقع یک طبقه‌بند خطی (linear classifier) است که می‌تواند تفاوت بین دو طبقه را تشخیص دهد. یک مثال پرسپترون به این صورت بود که می‌خواستیم تفاوت میان پراید و اتوبوس را با توجه به ویژگی‌های آن‌ها که همان طول و ارتفاع ماشین بود، به دست بیاوریم.

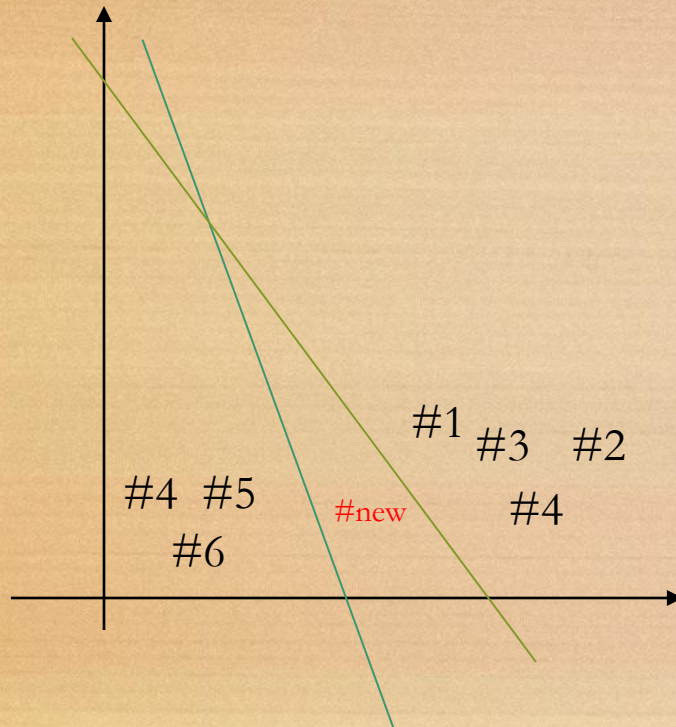


- خطِ آبی که در تصویر بالا مشاهده می‌کنید، در واقع خطی است که پرسپترون (با توجه به پارامترهای W - شیب خط و b - انحراف خط) رسم کرده است. پرسپترون با استفاده از این خط می‌تواند بفهمد نقاطی که پایین‌تر از این خط هستند، نقاط پراید و نقاطی که بالاتر از این خط هستند نقاط اتوبوس هستند.

- حال فرض کنید که خطِ آبی بالا به صورتی دیگر کشیده می‌شد. مثلاً مانند شکل زیر:
در واقع این خطی است که باز هم به صورت کامل پراید و اتوبوس را از هم دیگر جدا کرده است.



- حال فرض کنید یک نمونه جدید ماشین می‌آید و الگوریتم نمی‌داند که این ماشین چیست. پس باید با توجه به خط جدا کننده‌ای که توسط داده‌های آموزشی یادگرفته است، این ماشین جدید را طبقه بندی کند. مانند #new در شکل زیر:



اگر خط آبی یعنی همان خط جداکننده باشد، این اتومبیل جدید به دسته پرایدها می‌رود، این در حالی است که اگر خط نارنجی خط جداکننده باشد، این اتومبیل به دسته‌ی اتوبوس‌ها می‌رود. البته ما به عنوان یک انسان می‌توانیم بفهمیم که این اتومبیل بیشتر به پرایدها شبیه است چون نزدیک نمونه‌های پرایدهاست. ولی الگوریتم با توجه به خطی که از داده‌های آموزشی یاد گرفته است، این تمایز و طبقه‌بندی را انجام می‌دهد. در واقع **مشکل اصلی پرسپترون** همین جاست که خطی که رسم می‌کند، خط بهینه‌ای نیست و ممکن است خطا ایجاد کند. بنابراین به روش به مراتب پیشرفته‌تری به اسم **ماشین بردار پشتیبان** می‌رسیم که این مشکل را برطرف خواهد کرد.

ماشین بردار پشتیبان (Support Vector Machine) چیست؟

✓ الگوریتم SVM اولیه در ۱۹۶۳ توسط وپنیک ابداع شد و در سال ۱۹۹۵ توسط وپنیک و کورتس برای حالت **غیر خطی** تعمیم داده شد.

✓ ماشین بردار پشتیبان که به اختصار به آن **SVM** گفته می‌شود یک الگوریتم **یادگیری ماشین با ناظر** است که نمونه‌ی داده‌هایی را به صورتی نقاطی در فضا نشان داده شده است، با استفاده از یک خط یا هایپرپلین (Hyperplane)، از هم جدا می‌کند. این جداسازی به‌گونه‌ای است که نقاط داده‌ای که در یک طرف خط هستند مشابه به هم و در یک گروه قرار می‌گیرند. نمونه داده‌های جدید هم بعد از اضافه شدن به همان فضا در یکی از دسته‌های موجود قرار خواهند گرفت.

✓ اصول نظری مربوط به این الگوریتم ها بر پایه **تئوری یادگیری آماری** (statistical learning theory) است. مبنای نظریه مزبور، بر قائل شدن تفاوت بین نمونه های مختلف در حین یادگیری است.

✓ ماشین های بردار پشتیبانی از **دقیق ترین و نیرومندترین** الگوریتم های داده کاوی است.

✓ در اغلب موارد هیچ حساسیتی نسبت به ابعاد داده ها ندارند.

✓ شامل دسته بندی کننده بردار پشتیبانی (SVC) و رگرسیون (رگرسیون) بردار پشتیبانی (SVR) است.

✓ این روش از جمله روش‌های نسبتاً جدیدی است که در سال‌های اخیر کارایی خوبی نسبت به روش‌های قدیمی‌تر برای طبقه‌بندی از جمله شبکه‌های عصبی پرسپترون نشان داده است.

✓ مبنای کاری دسته‌بندی‌کننده SVM دسته‌بندی خطی داده‌ها است و در تقسیم خطی داده‌ها سعی می‌کنیم خطی را انتخاب کنیم که حاشیه اطمینان بیشتری داشته باشد. حل معادله پیدا کردن خط بهینه برای داده‌ها به وسیله روش‌های QP (Quadratic Programming) که روش‌های شناخته شده‌ای در حل مسائل محدودیت‌دار هستند صورت می‌گیرد.

✓ SVM‌ها به خانواده‌ای از مدل‌های خطی تعمیم یافته تعلق دارد.

✓ SVM دسته‌بندی‌کننده‌ای است که جزو شاخه روش‌های کرنل (Kernel Methods) در یادگیری ماشین محسوب می‌شود.

✓ یکی از روش‌هایی که در حال حاضر به صورت گسترده برای مسئله دسته‌بندی (Classification) مورد استفاده قرار می‌گیرد، روش ماشین بردار پشتیبان است. شاید به گونه‌ای بتوان محبوبیت کنونی روش ماشین بردار پشتیبان را با محبوبیت شبکه‌های عصبی در دهه گذشته مقایسه کرد. علت این قضیه نیز قابلیت استفاده این روش در حل مسائل گوناگون می‌باشد، در حالیکه روش‌هایی مانند درخت تصمیم‌گیری را نمی‌توان به راحتی در مسائل مختلف به کار برد.

- در صورت استفاده مناسب از SVM این الگوریتم قدرت تعمیم خوبی خواهد داشت.
- علیرغم داشتن ابعاد زیاد (high dimensionality) از overfitting پرهیز میکند. این خاصیت ناشی از optimization این الگوریتم است
- **فشرده سازی اطلاعات:**
- بجای داده های آموزشی از بردارهای پشتیبان استفاده میکند.

➤ **هدف** این دسته الگوریتم ها تشخیص و متمایز کردن الگوهای پیچیده در داده هاست

➤ **مسایل مطرح:**

- الگوهای پیچیده را چگونه نمایش دهیم
- چگونه از مسئله overfitting پرهیز کنیم

کاربردهای SVM:

- تشخیص های پزشکی
- بیوانفورماتیک
- پردازش تصویر
- متن کاوی

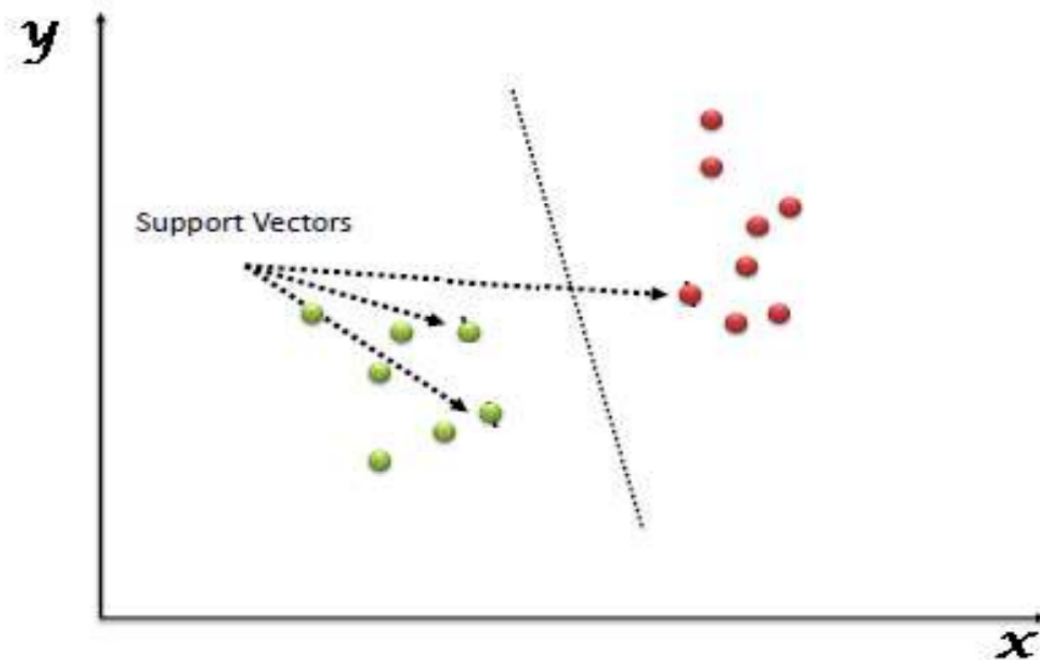
ماشین های بردار پشتیبانی **مشابه** با شبکه های عصبی قادرند تا برای هر تابع چند متغیره تقریب هایی را با درجه دقت دلخواه به دست بیاورند. بنابراین به منظور مدل کردن سیستم ها و فرآیند های غیر خطی و بسیار پیچیده می توان از **SVM** ها استفاده کرد.

- با فرض اینکه دسته ها بصورت خطی جداپذیر باشند، ابرصفحه هائی با حداکثر حاشیه (maximum margin) را بدست می آورد که دسته ها را جدا کنند.
- در مسایلی که داده ها بصورت خطی جداپذیر نباشند داده ها به فضای با ابعاد بیشتر نگاشت پیدا میکنند تا بتوان آنها را در این فضای جدید بصورت خطی جدا نمود.

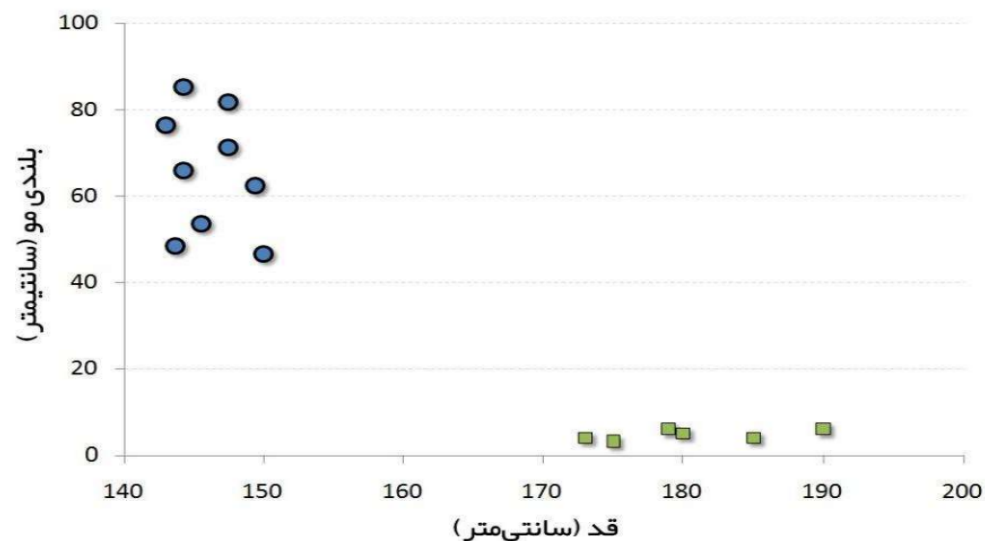
✓ **هدف** عبارت است از یافتن ابر صفحه جداکننده نقاط داده ای متعلق به دو کلاس با **حاشیه ماکسیمال** و **بهترین توانایی تعمیم**.

✓ **حاشیه** از دیدگاه هندسی عبارت است از فاصله موجود بین ابرصفحه و نزدیک ترین نمونه آموزشی. از یک زاویه دیگر حاشیه اینگونه تعریف می شود : مقدار فضا یا جدایی موجود میان دو کلاس که توسط ابرصفحه تعریف می شود.

✓ منظور از **توانایی تعمیم** آن است که دسته بندی کننده علاوه بر داشتن عملکرد خوب در دسته بندی داده های آموزشی (دقت)، دارای یک دقت **پیش بینی** بالا در قبال داده های مشاهده نشده (که توزیع آن ها با توزیع داده های آموزشی یکسان است) نیز باشد.



به بیان ساده، بردارهای پشتیبان در واقع مختصات یک مشاهده منفرد هستند. ماشین بردار پشتیبان مرزی است که به بهترین شکل دسته‌های داده‌ها را از یکدیگر جدا می‌کند.



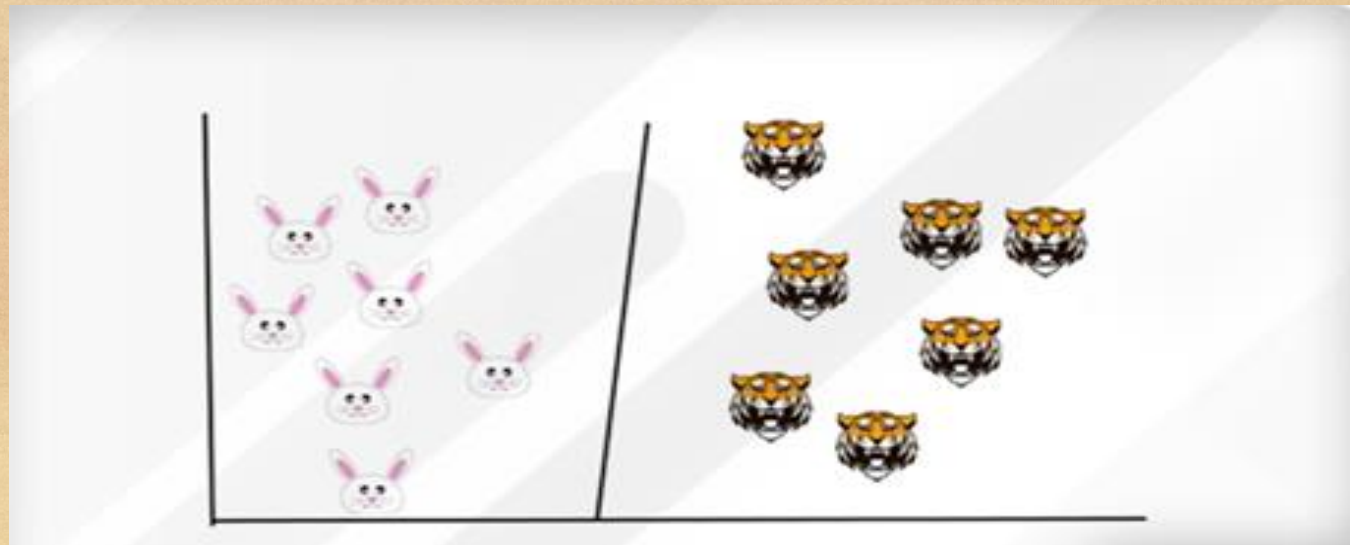
دایره‌های آبی موجود در نمودار نماینده زنان و مربع‌های سبز نماینده مردان هستند. برخی از برداشت‌هایی که می‌توان از این نمودار داشت عبارتند از:

- (1) مردان در جمعیت مثال ما، میانگین قد بلندتری دارند.
 - (2) زنان در جمعیت مثال ارائه شده، بلندی موی بیشتری دارند.
- بر این اساس اگر فردی با قد ۱۸۰ سانتی‌متر و طول موی ۴ سانتی‌متر در جمعیت وجود داشته باشد، بهترین حدسی که می‌توان زد آن است که فرد در دسته مردان قرار می‌گیرد.

ماشین بردار پشتیبان چطور کار می‌کند؟

برای درک نحوه‌ی عملکرد ماشین بردار پشتیبان فرض کنید صاحب مزرعه‌ای هستید و بنا به دلایلی می‌خواهید حصاری برای محافظت از خرگوش‌های خود در برابر ببرها ایجاد کنید. **حصار خود را کجا بسازیم؟**

- یکی از راه‌حل‌های این مشکل این است که یک طبقه‌بندی براساس موقعیت خرگوش‌ها و ببرها ایجاد کنیم. می‌توانیم گروه خرگوش‌ها را به‌عنوان یک گروه و ببرها را به‌عنوان گروه دیگر طبقه‌بندی کنیم.
- در حال حاضر، اگر سعی کنیم یک مرز میان خرگوش‌ها و ببرها بکشیم، یک خط مستقیم خواهد شد. ماشین بردار پشتیبان نیز دقیقاً به این شکل عمل می‌کند؛ یک مرز تصمیم‌گیری ترسیم می‌کند که درواقع یک هایپرپلین میان دو کلاس است تا آن‌ها را از هم جدا و طبقه‌بندی کند.



چگونه بدانیم هایپرپلین را کجا باید بکشیم؟

اصل اساسی ماشین بردار پشتیبان این است که یک هایپرپلین ترسیم کنیم که دو کلاس را به بهترین شکل از یکدیگر جدا کند. حال در مورد مثال ما دو کلاس خرگوش و ببر هستند؛ بنابراین ما با ترسیم یک هایپرپلین رندوم شروع می‌کنیم و سپس فاصله‌ی میان هایپرپلین و نزدیک‌ترین نقاط داده‌ی هر کلاس را بررسی می‌کنیم.

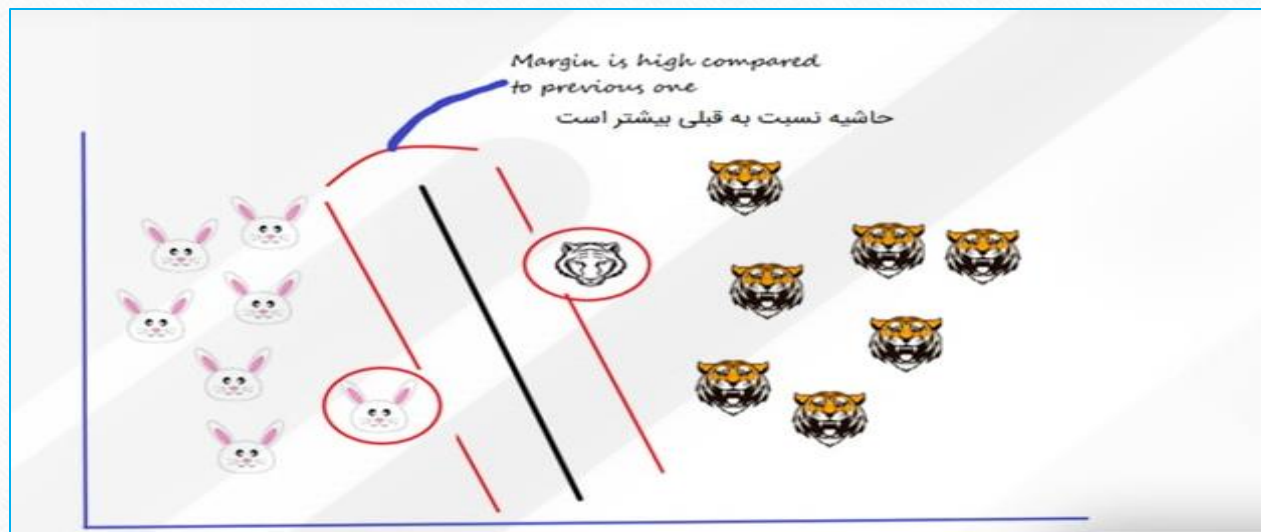


این نقاط داده نزدیک به هایپرپلین بردارهای پشتیبان (Support Vectors) نامیده می‌شوند؛ به همین دلیل است که به این الگوریتم ماشین بردار پشتیبان گفته می‌شود. اساساً هایپرپلین بر اساس این بردارهای پشتیبان ترسیم می‌شود. به‌طور معمول هایپرپلینی که بیشترین فاصله از بردارهای پشتیبان را داشته باشد بهینه‌ترین هایپرپلین است. این فاصله‌ی میان هایپرپلین و بردارهای پشتیبان حاشیه (Margin) نامیده می‌شود.

حال بیابید فرض کنیم یک نقطه‌ی داده‌ی جدید اضافه کنیم (در این مثال ببر دیگری اضافه شده است). اکنون می‌خواهیم یک هایپرپلین بکشیم تا این دو کلاس را به‌بهترین شکل از هم جدا کند؛ بنابراین، با ترسیم هایپرپلین، همان‌طور که در تصویر نشان داده شده است، شروع می‌کنیم؛ سپس فاصله‌ی میان این هایپرپلین و بردارهای پشتیبان را بررسی می‌کنیم که آیا حاشیه‌ی این هایپرپلین حداکثر است یا خیر. در این شکل حاشیه‌ی خیلی هم زیاد نیست.



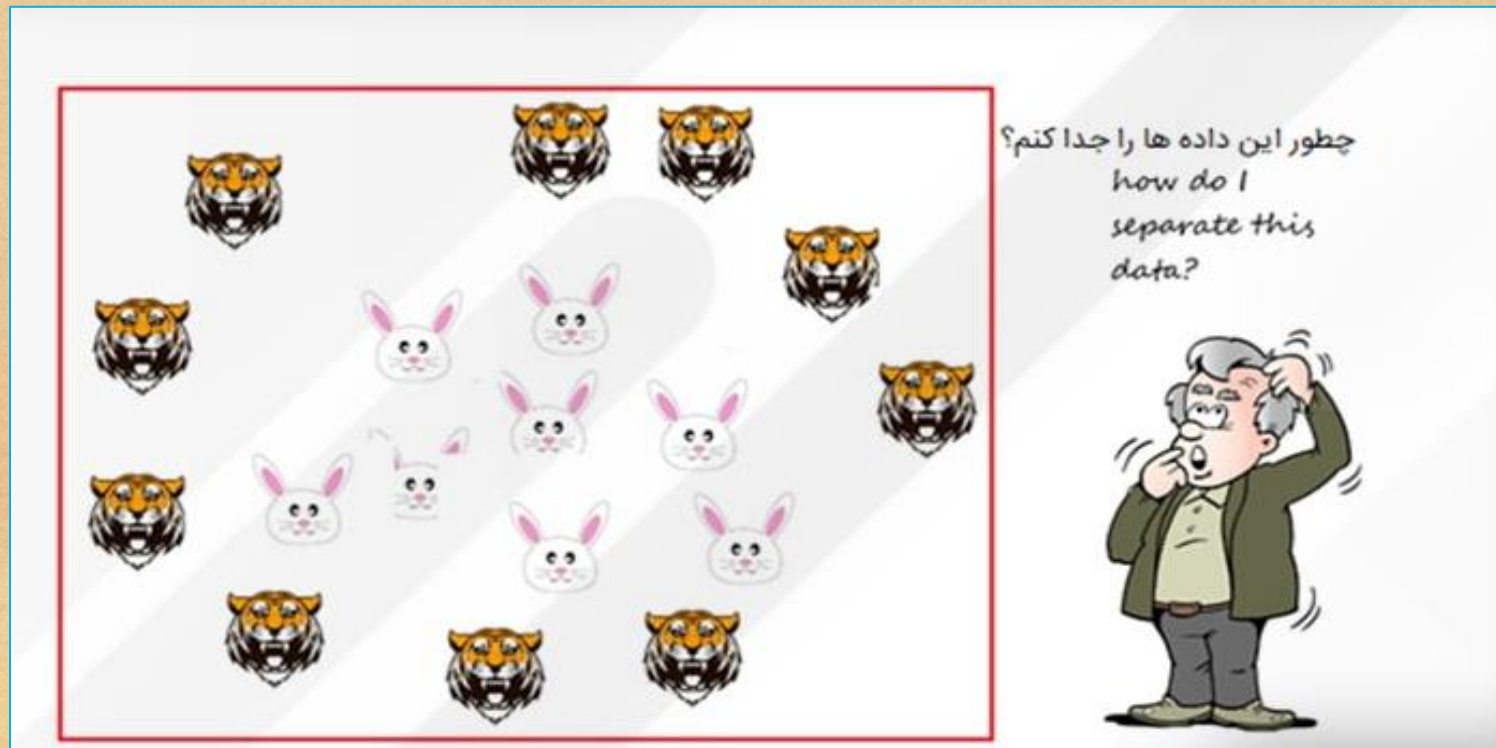
در سناریوی دوم، یک هایپرپلین متفاوت، مانند تصویر زیر، ترسیم می‌کنیم و سپس فاصله‌ی میان هایپرپلین و بردارهای پشتیبان را بررسی می‌کنیم که آیا حاشیه‌ی این هایپرپلین حداکثر است یا خیر؟

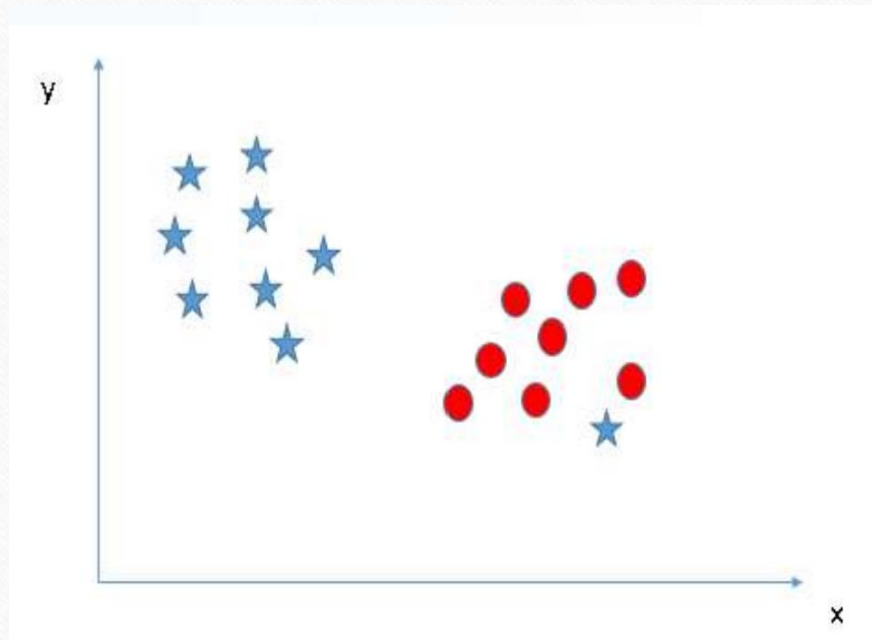
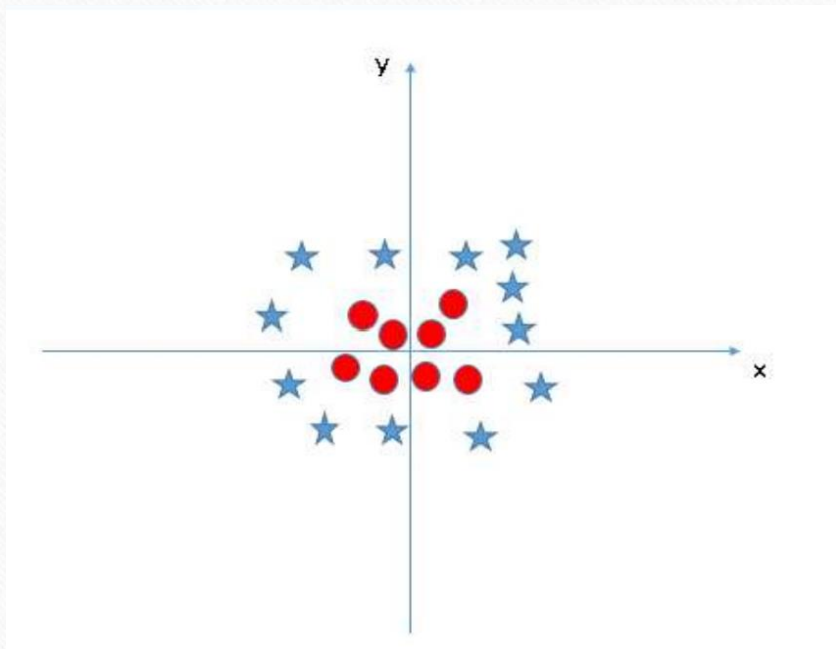


می‌بینیم که حاشیه در مقایسه با هایپرپلین قبلی بسیار زیاد است؛ بنابراین ما این هایپرپلین را انتخاب می‌کنیم.

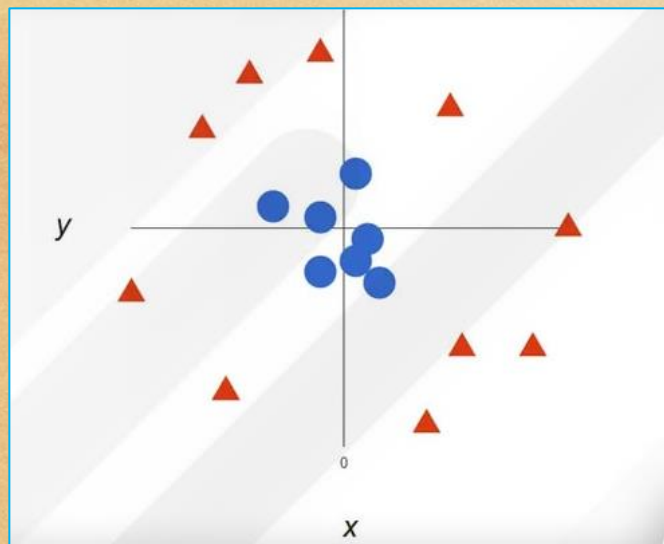
ماشین بردار پشتیبان و داده‌های غیرخطی

- در مثال قبل دیدیم که داده‌های ما تا کنون به صورت خطی تفکیک‌پذیر بودند، یعنی می‌توانستیم یک خط مستقیم برای جدا کردن دو کلاس بکشیم، اما اگر نقاط داده‌ی ما به شکل زیر باشد، چه کنیم؟ این دو کلاس با یک خط مستقیم از هم جدا نمی‌شوند.
- در الگوریتم ماشین بردار پشتیبان، هنگام ساختن یک هایپرپلین برای نقاط داده‌ی تفکیک‌پذیر خطی کار آسانی است، اما وقتی داده‌ها **غیرخطی** تفکیک‌پذیر باشند، کار بسیار چالش‌برانگیز خواهد بود.





هنگامی که نقاط داده را نمی‌توان با یک خط مستقیم یا یک هایپرپلین مستقیم جدا کرد، مسئله غیرخطی نامیده می‌شود. در چنین شرایطی **کرانل‌های** ماشین بردار پشتیبان وارد عمل می‌شوند و ابعاد فضا را افزایش می‌دهند تا نقاط داده به صورت خطی تفکیک پذیر شوند.



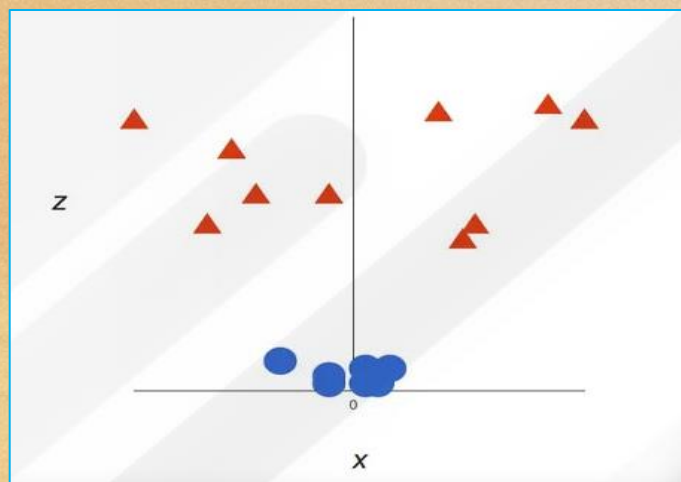
پس اگر داده‌های اولیه ما به شکل روبرو باشند، لازم است که ما بعد سوم را اضافه کنیم. تابع حال ما دو بعد داشتیم: x و y حال یک بعد z جدید ایجاد می‌کنیم که به این شکل محاسبه می‌شود:

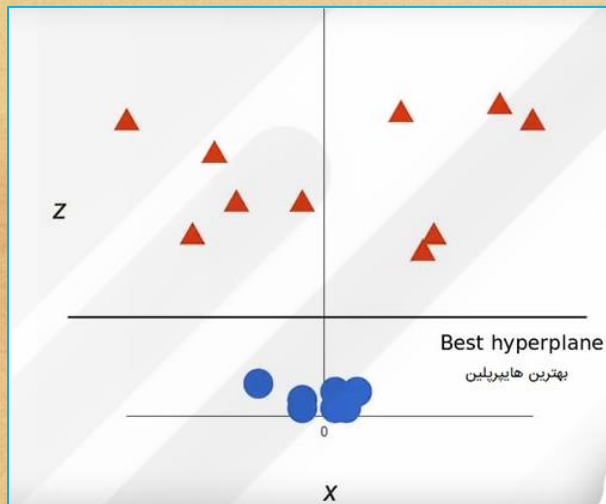
$$z = x^2 + y^2$$

(اگر دقت کنید این معادله یک دایره است).



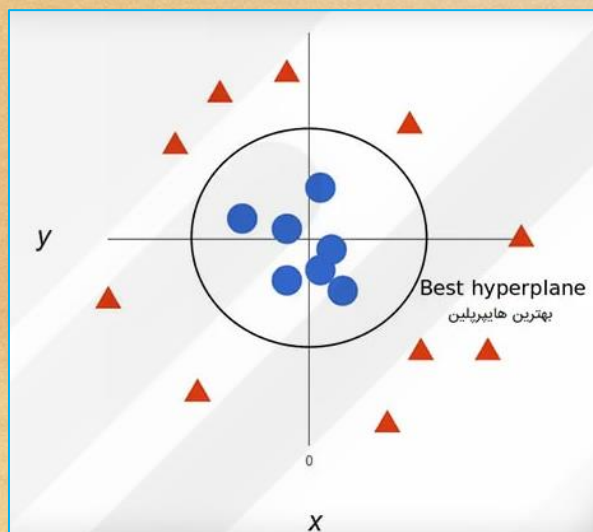
این کار به ما فضایی سه بعدی می‌دهد که به این شکل می‌توان آن را در اینجا نشان داد:



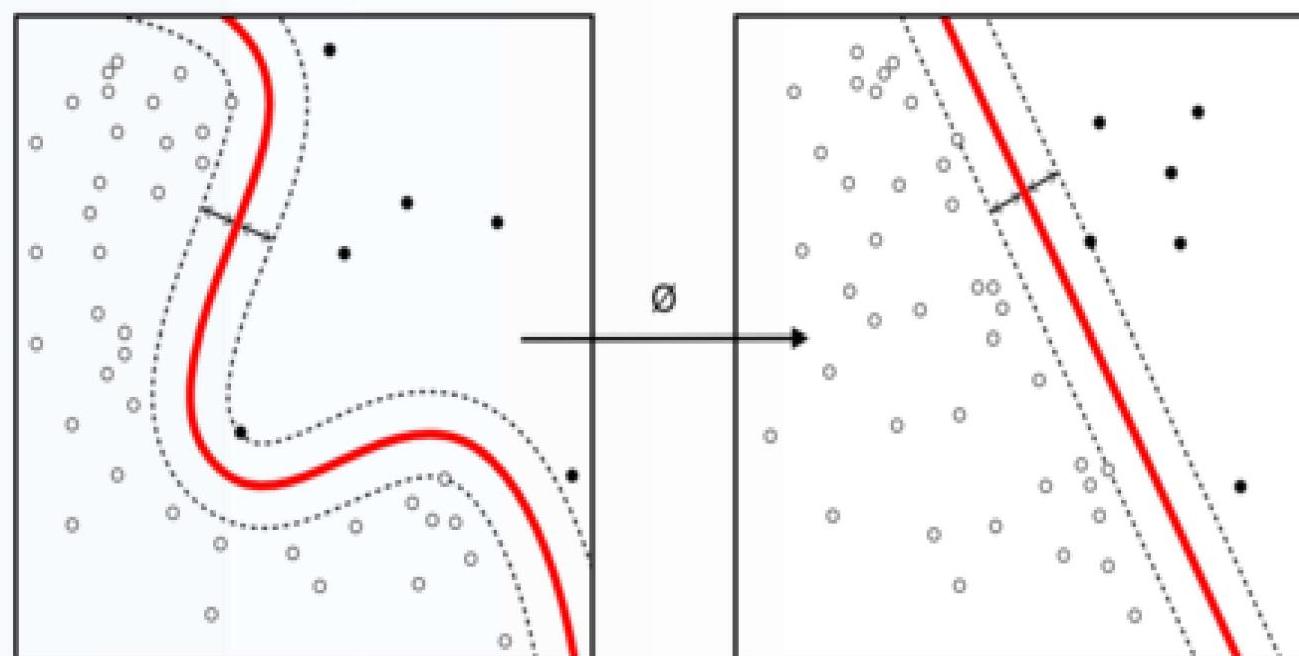


اما ماشین بردار پشتیبان الان چطور تفکیک را انجام می‌دهد؟ اجازه بدهید با هم ببینیم:

پس نقاط داده‌ی ما در حال حاضر به راحتی با یک خط تفکیک شدند. آنچه باقی مانده ترسیم مجدد آن به شکل دوبعدی است:

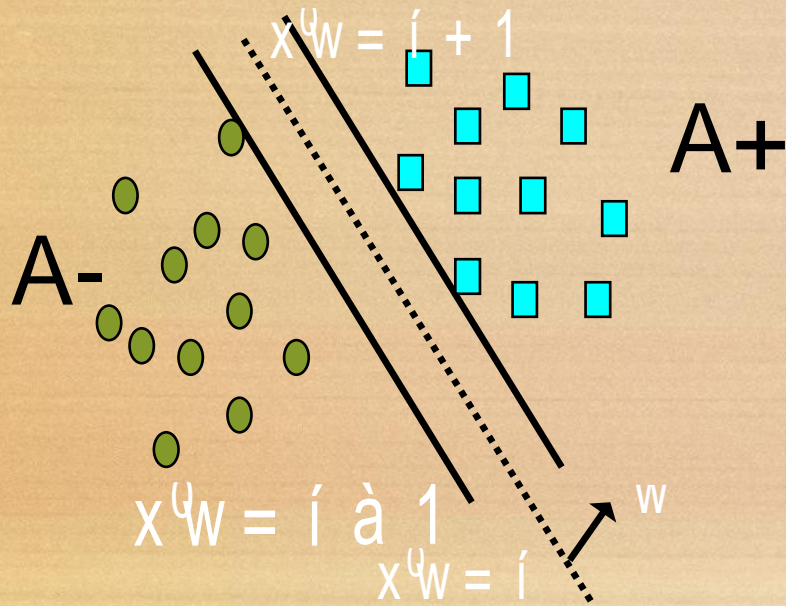


در حال حاضر کار طبقه‌بندی داده‌ها به پایان رسیده است. در واقع به این فرایند که توضیح داده شد حقه‌ی کرنل (Kernel Trick) گفته می‌شود.



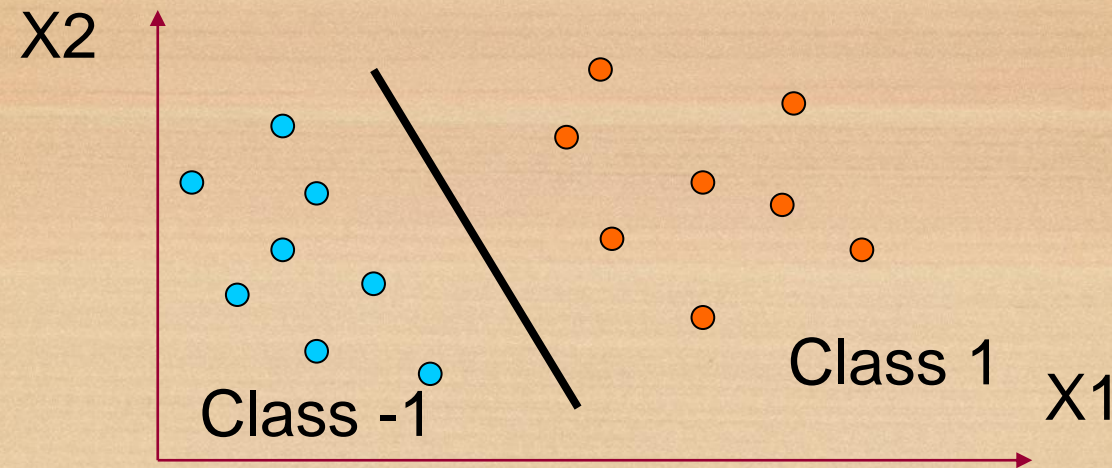
مسئله جداسازی خطی: Linear Discrimination

- اگر دو دسته وجود داشته باشند که بصورت خطی از هم جداپذیر باشند، بهترین جدا کننده این دو دسته چیست؟
- الگوریتم های مختلفی از جمله پرسپترون میتوانند این جداسازی را انجام دهند.
- آیا همه این الگوریتم ها به خوبی از عهده اینکار بر می آیند؟



Separating Surface:

خط یا ابر صفحه جدا کننده



✓ **هدف:** پیدا کردن بهترین خط (ابر صفحه) که دو دسته را از هم جدا کند. در حالت دو بعدی معادله این خط بصورت زیر است:

$$w_1 X_1 + w_2 X_2 + b = 0$$

در ابر صفحه بهینه w و b به ترتیب عبارتند از بردار وزن و یک اسکالر.

✓ در حالت n بعدی خواهیم داشت:

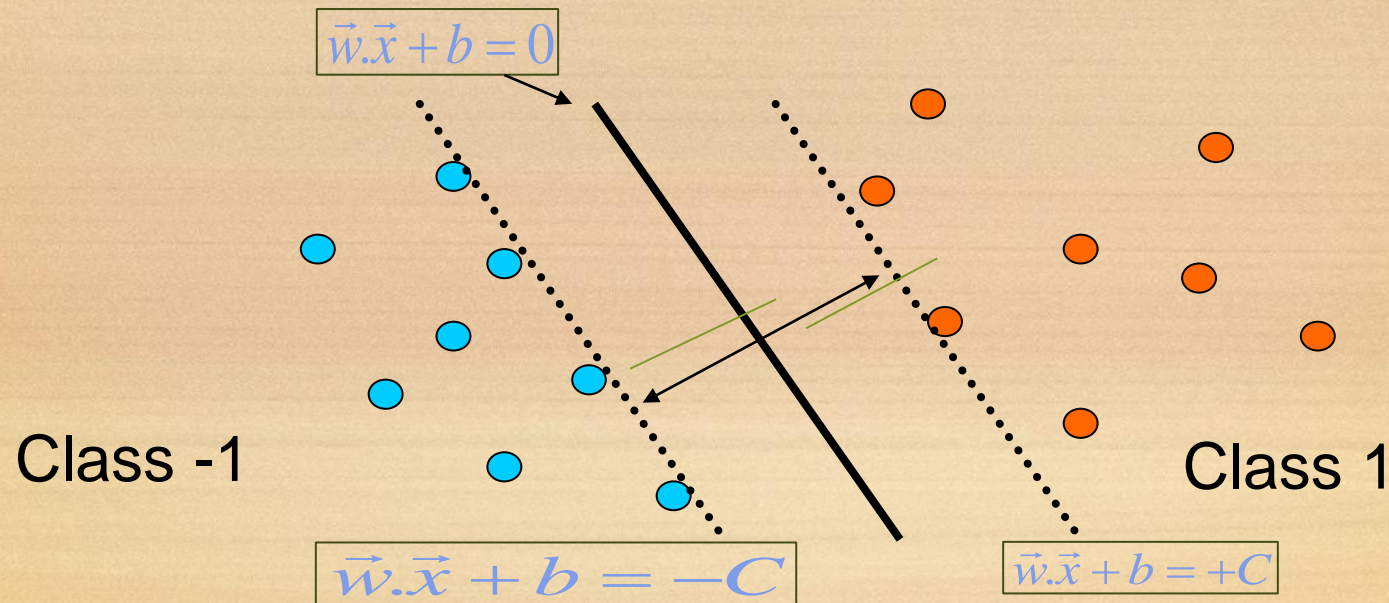
$$\sum_{i=0}^n w_i . x_i + b = 0$$

$$\vec{w}^T . \vec{x} + b = 0$$

ایده SVM برای جدا سازی دسته ها:

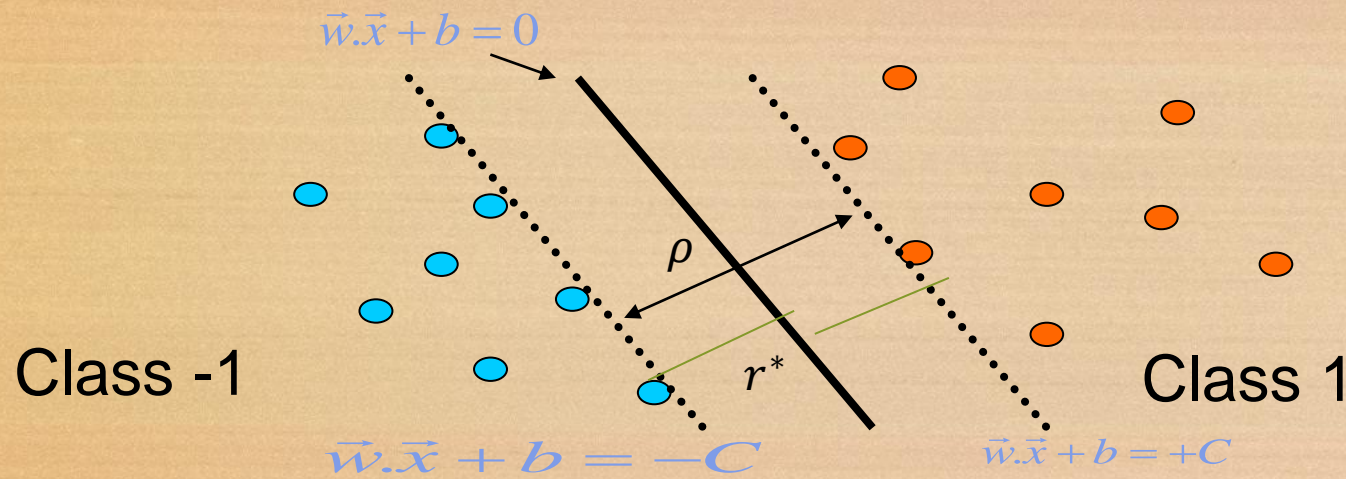
دو صفحه مرزی بسازید :

- دو صفحه مرزی موازی با صفحه دسته بندی رسم کرده و آن دو را آنقدر از هم دور می‌کنیم که به داده ها برخورد کنند.
- صفحه دسته بندی که بیشترین فاصله را از صفحات مرزی داشته باشد، بهترین جدا کننده خواهد بود.



حداکثر حاشیه

- بر طبق قضیه ای در تئوری یادگیری اگر داده های آموزشی بدرستی دسته بندی شده باشند، از بین جداسازهای خطی، آن جداسازی که حاشیه داده های آموزشی را حداکثر میکند خطای تعمیم را حداقل خواهد کرد.



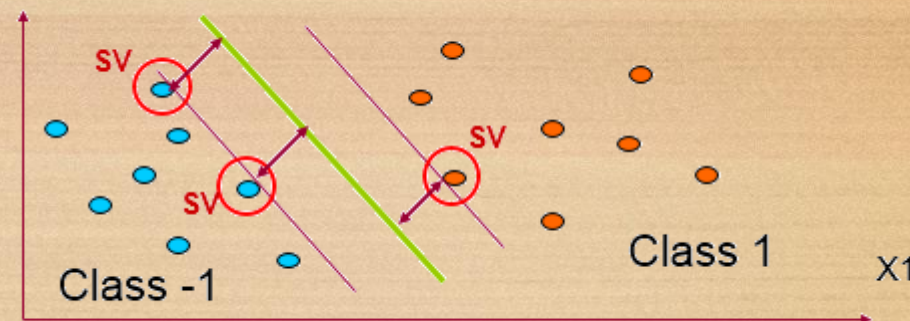
بنابراین در ارتباط با ابرصفحه بهینه پارامترهای w و b را می بایست به نحوی تعیین کرد که حاشیه جداسازی (ρ) بیشینه گردد. این حاشیه با توجه به کوتاه ترین فواصل هندسی r^* از دو کلاس تعیین می شود. با توجه به مطالب بیان شده **SVC** را دسته بندی کننده با **حاشیه ماکسیمال** می نامند.

بردار پشتیبان

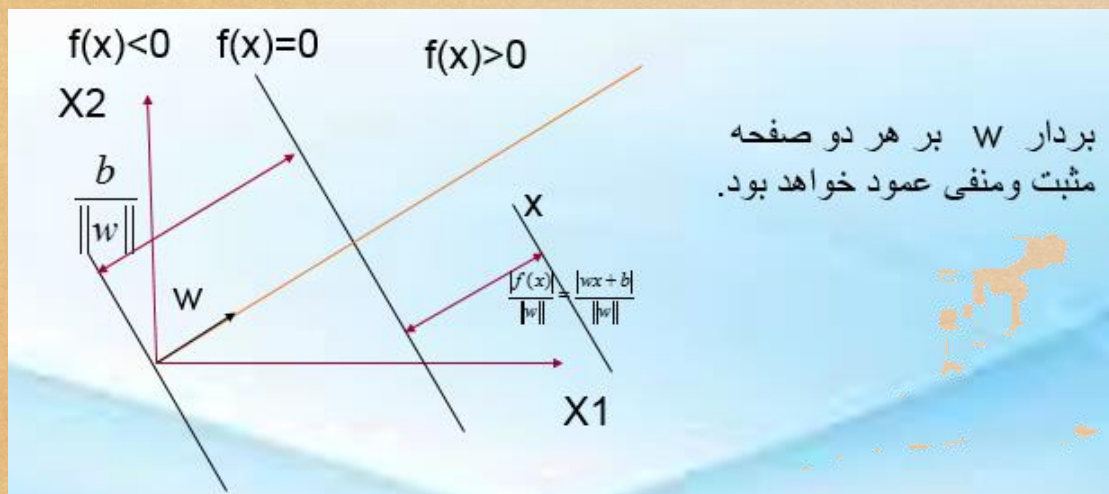
- نزدیکترین داده های آموزشی به ابر صفحه های جدا کننده بهینه **بردار پشتیبان** نامیده میشوند. نقاط داده ای ویژه ای که نامساوی های زیر را به تساوی تبدیل می کنند.

$$\begin{aligned} w^T x_i + b &\geq 1 & \text{for } y_i = +1 \\ w^T x_i + b &\leq -1 & \text{for } y_i = -1 \end{aligned}$$

X2



حل مسئله برای حالت دو بعدی



- فاصله خط جداکننده از مبدا برابر است با
- فاصله نمونه ای مثل x از خط جدا کننده برابر است با

$$\frac{|f(x)|}{\|w\|} = \frac{|wx + b|}{\|w\|}$$

تعیین حاشیه بین خطوط جدا کننده

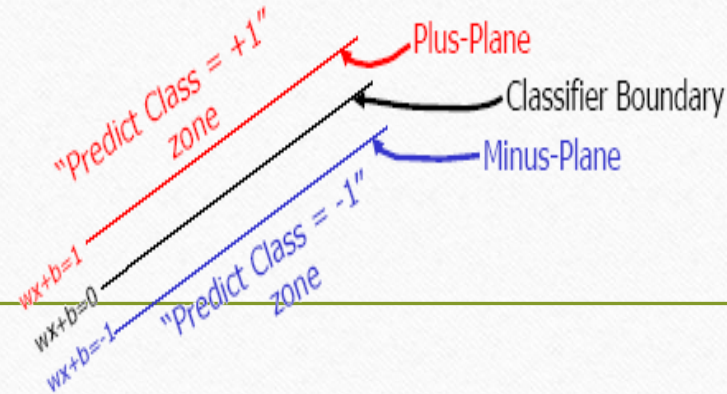
$$\text{Plus-plane} = \{ \mathbf{x} : \mathbf{w}^T \cdot \mathbf{x} + b = +1 \}$$

$$\text{Minus-plane} = \{ \mathbf{x} : \mathbf{w}^T \cdot \mathbf{x} + b = -1 \}$$

Classify as..

$$-1 \quad \text{if } \mathbf{w}^T \cdot \mathbf{x} + b \leq -1$$

$$+1 \quad \text{if } \mathbf{w}^T \cdot \mathbf{x} + b \geq 1$$



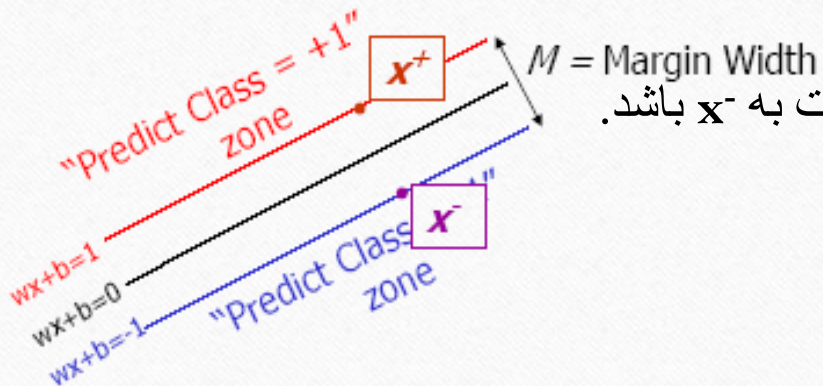
محاسبه پهنای حاشیه

• صفحه مثبت و منفی را بصورت زیر در نظر میگیریم:

- Plus-plane = $\{ \mathbf{x} : \mathbf{w}^T \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w}^T \cdot \mathbf{x} + b = -1 \}$

• بردار w بر صفحه مثبت و منفی عمود خواهد بود.

• فرض کنید x^- نقطه ای در صفحه منفی بوده و x^+ نزدیکترین نقطه در صفحه مثبت به x^- باشد.



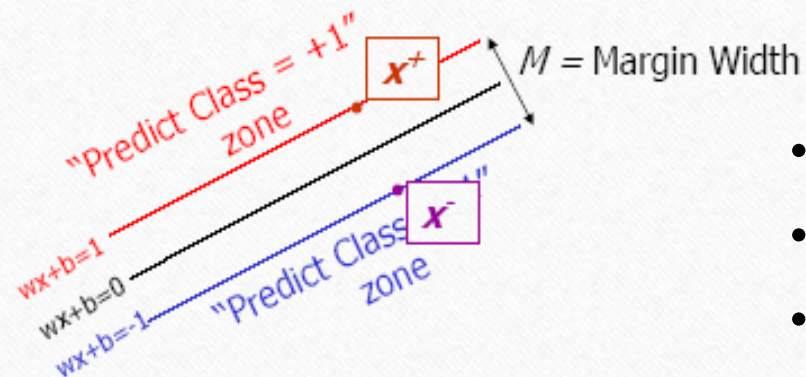
محاسبه پهنای حاشیه

- خطی که x^- را به x^+ وصل میکند بر هر دو صفحه عمود خواهد بود. لذا فاصله بین دو صفحه ضربی از W خواهد بود.

در اینصورت خواهیم داشت:

$$x^+ = x^- + \lambda w \quad \text{for some value of } \lambda.$$

- میدانیم که:



- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$

- لذا میتوان M را برحسب W و b محاسبه کرد.

- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$

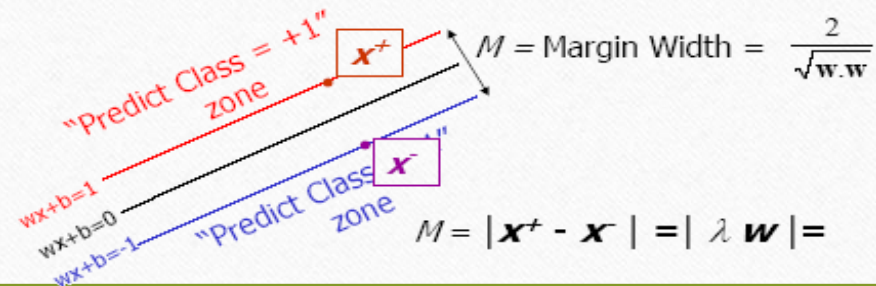

$$w \cdot (x^- + \lambda w) + b = +1$$

$$w \cdot x^- + \lambda w \cdot w + b = +1$$

$$-1 + \lambda w \cdot w = +1$$

$$\lambda = 2 / w \cdot w$$

محاسبه پهنای حاشیه



What we know:

- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$
- $\lambda = \frac{2}{w \cdot w}$

$$= \lambda |w| = \lambda \sqrt{w \cdot w}$$

$$= \frac{2\sqrt{w \cdot w}}{w \cdot w} = \frac{2}{\sqrt{w \cdot w}}$$

محدودیت

• اگر برای مثال دو بعدی فوق مقدار دسته ها را با 1 و -1 مشخص کنیم داریم:

$$\langle w, x_i \rangle + b \geq 1 \text{ for } y_i = 1$$

$$\langle w, x_i \rangle + b \leq -1 \text{ for } y_i = -1$$

• که میتوان آنها را بصورت زیر نوشت:

$$y_i (\langle w, x_i \rangle + b) \geq 1 \text{ for all } i$$

- در SVM بدنبال حل همزمان معادلات زیر هستیم:

- با داشتن داده های آموزشی (x_i, y_i) که $y_i \in \{+1, -1\}$ $i=1,2,\dots,N$

- Minimise $\|w\|^2$

- Subject to : $y_i (<w, x_i> + b) \geq 1$ for all i

-Note that $\|w\|^2 = w^T w$

- این یک مسئله **quadratic programming** با محدودیت هائی بصورت نامعادلات خطی است. روشهای شناخته شده ای برای چنین مسئله هائی بوجود آمده اند.

$$\min_{w,b} \left(\frac{1}{2} \|w\|^2 \right)$$

Subject to :

$$y_i (w^T x_i + b) \geq 1 \quad i = 1, 2, \dots, n$$

معمولا به منظور حل مساله ای که در بالا مطرح شده است مساله اولیه از روش **ضرایب لاگرانژ** استفاده می کنیم .

معمولا به منظور حل مساله ای که در بالا مطرح شده است مساله اولیه از روش **ضرایب لاگرانژ** استفاده می کنیم . تابع لاگرانژ زیر را تشکیل می دهیم :

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

α_i عبارت است از **ضریب لاگرانژ** با توجه به α نامساوی .

از $L(w, b, \alpha)$ نسبت به w و b مشتق گرفته و نتایج حاصل را برابر با صفر قرار می دهیم . در نهایت دو شرط زیر را برای بهینگی به دست می آوریم :

$$\begin{cases} \frac{\partial L(w, b, \alpha)}{\partial w} = 0 \\ \frac{\partial L(w, b, \alpha)}{\partial b} = 0 \end{cases}$$

$$\begin{cases} w = \sum_{i=1}^n \alpha_i y_i x_i \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

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

Subject to :

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

$$\alpha_i \geq 0 \quad i = 1, 2, \dots, n$$

همزمان:

$$\alpha_i [y_i (w^T x_i + b) - 1] = 0 \quad i = 1, 2, \dots, n$$

بنابر این تنها آن بردارهای پشتیبانی (x_i, y_i) که کمترین فاصله را از ابر صفحه بهینه داشته و حاشیه ماکسیمال را تعیین می کنند متناظر با α_i های غیر صفر می باشند. بقیه α_i ها برابر با صفر هستند.

پس از تعیین ضرایب لاگرانژ بهینه، می توانیم بردار وزن بهینه w^* را محاسبه نماییم :

$$w^* = \sum_{i=1}^n \alpha_i y_i x_i$$

سپس به کمک بردار پشتیبانی مثبت، مقدار بهینه اسکالر b را به صورت زیر تعیین می کنیم :

$$b^* = 1 - w^{*T} x_s \quad \text{for } y_s = +1$$

فرمول لاگرانژ..

• توجه شود که برچسب های مختلف لاگرانژ تاکید بر این نکته دارد که $L(w,b,a)$ برای اصلی و D برای دوگان) از تابع هدف یکسان ولی با شرایط متفاوت؛ دو راه-حل، از دو مسیر زیر بدست می آید:

• مینیمم کردن $L(w,b,a)$

• ماکزیمم کردن L_D

• محاسبه a_i ، به حل معادله بالا، موسوم به quadratic problem منجر شود.

➤ مسئله جداسازی غیر خطی:

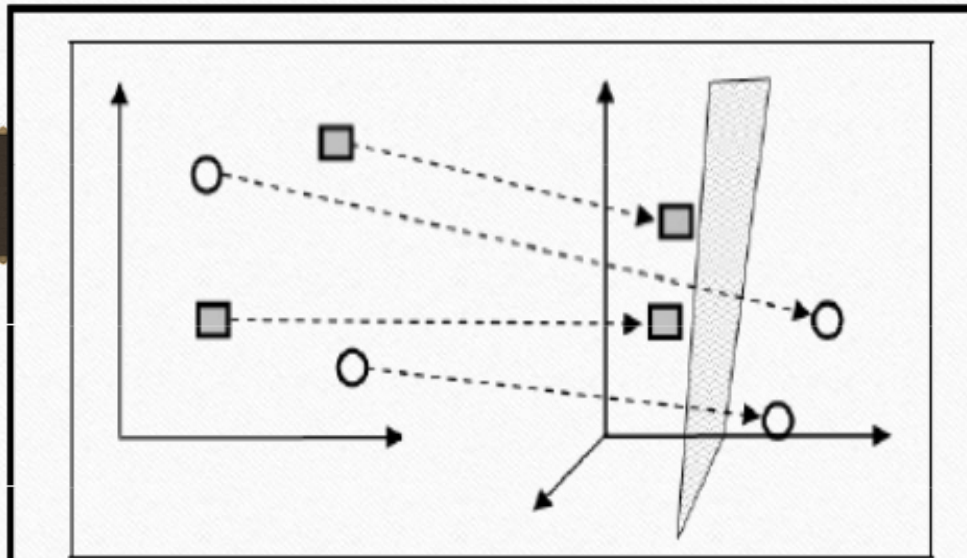
- ✓ SVC با حاشیه ماکسیمال، نقطه آغازین الگوریتم های SVM قلمداد می شود.
- ✓ زمانی که تفکیک خطی نمونه ها بطور کامل شدنی نیست این امکان وجود دارد که حاشیه ها منفی باشند. در چنین مواردی فضای شدنی مساله اولیه تهی و حل مساله بهینه سازی غیر ممکن می باشد.

➤ دو رویکرد برای برای حل مسائل تفکیک ناپذیر (غیر خطی):

- ✓ بهینه سازی با حاشیه های انعطاف پذیر (نرم)
- ✓ بکارگیری kernel trick به منظور خطی سازی مسائل غیر خطی مورد نظر

جدائی پذیری خطی و غیر خطی، مفهوم کرنل

نمایش دستیابی به جدایی پذیری خطی، برای مسئله ای که دارای این خاصیت نیست، به کمک نگاشت



Kernel Mapping from input to feature space.

- در حالتی که جدا پذیری بصورت خطی نباشد (جدا پذیری غیر خطی) ایده اصلی این است که نمونه ها را به یک فضای با بعد بالا (feature space) فضای مشخصه نگاشت دهیم که در فضای جدید مشخصه ها، نمونه ها می توانند به صورت خطی از هم جدا شوند.
- این امر نیاز به اعمال یک تابع هسته (کرنل) را به همراه خواهد آورد.

SVC با حاشیه انعطاف پذیر (نرم) و بهینه سازی:

✓ مواردی را در نظر بگیرید که در آن ها تعدادی از نقاط متعلق به کلاس های مختلف با یکدیگر مخلوط شده اند . این نقاط بیانگر خطای آموزش می باشند

✓ هدف از بکارگیری حاشیه انعطاف پذیر توسعه دادن الگوریتم SVC است به نحوی که در تعیین ابرصفحه امکان وجود داده های حاوی اختلال نیز فراهم شود .

به منظور دستیابی به این هدف یک متغیر کمبود ξ_i معرفی می شود تا از این طریق بتوان مقدار خطای دسته بندی را در نظر گرفت :

$$\min_{w,b} \left(\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \right)$$

Subject to :

$$y_i (w^T x_i + b) \geq 1 - \xi_i \quad \xi_i \geq 0 \quad i = 1, 2, \dots, n$$

نقش C در رابطه فوق ایجاد تعادل میان پیچیدگی ماشین و تعداد نقاط تفکیک ناپذیر است . این پارامتر توسط کاربر و بر اساس تجربه یا تحلیل تعیین می شود .

متغیر کمبود ξ_i بیانگر فاصله موجود میان ابرصفحه و داده ای است که دسته بندی آن نادرست می باشد . در واقع این فاصله میزان انحراف یک نمونه را از وضعیت ایده ال تفکیک پذیری می سنجد .

با استفاده از روش ضرایب لاگرانژ می توان مساله دوگان متناظر با حاشیه انعطاف پذیر را به صورت زیر فرموله کرد :

$$\begin{aligned} \max_{\alpha} W(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{Subject to:} \\ \sum_{i=1}^n \alpha_i y_i &= 0 \\ 0 \leq \alpha_i &\leq C \quad i = 1, 2, \dots, n \end{aligned}$$

در حالت تفکیک ناپذیر شرط مکمل کרוش - کان - تاکر به صورت زیر است

$$\alpha_i [y_i (w^T x_i + b) - 1 + \xi_i] = 0 \quad i = 1, 2, \dots, n$$

$$\gamma_i \xi_i = 0 \quad i = 1, 2, \dots, n$$

و

$$\xi_i = 0 \quad \text{if} \quad \alpha_i < C$$

در نقطه زینی $\alpha_i + \gamma_i = C$ است. بنابراین

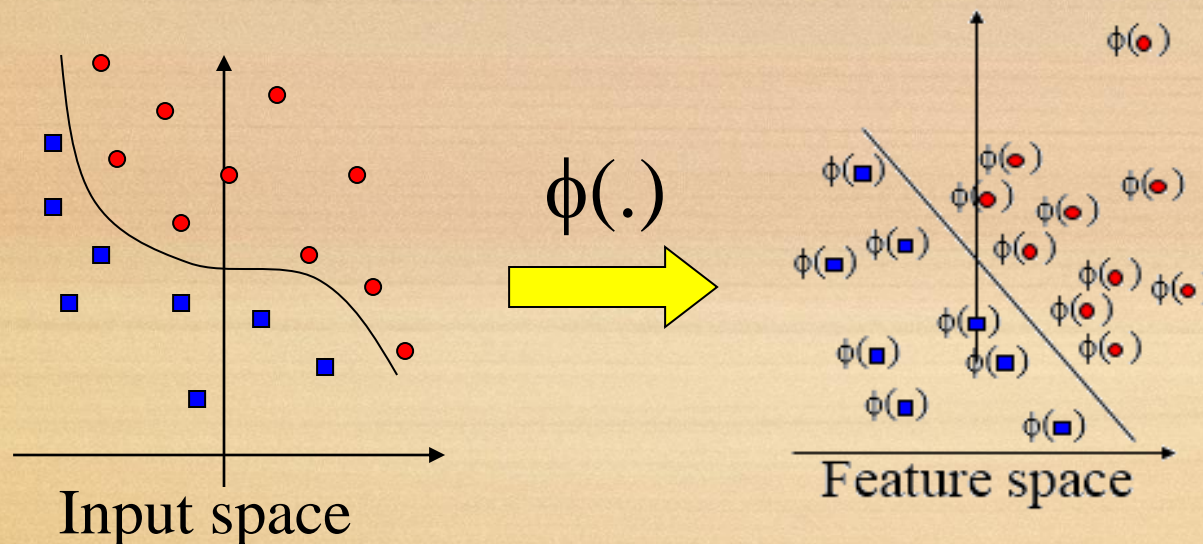
$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

در نتیجه وزن بهینه برابر است با :

مقدار بهینه برای اسکالر b با توجه به رابطه اول و هر نقطه داده ای دلخواه در مجموعه آموزشی که برای آن داشته باشیم $0 \leq \alpha_i \leq C$ به دست می آید

➤ Kernel trick:

- ✓ بر اساس ضرب داخلی داده های مفروض یک تابع کرنل مناسب تعریف می شود
 - ✓ با یک تبدیل غیر خطی از فضای ورودی به فضای خصیصه با ابعاد بیش تر (حتی نامتناهی) مواجه هستیم تا از این طریق بتوان مسائل را به صورت خطی تفکیک پذیر ساخت
 - ✓ این رویکرد را می توان با توجه به قضیه پوشش در مورد تفکیک پذیری الگوها توجیه کرد : احتمال آنکه یک مساله پیچیده دسته بندی الگو در فضایی با ابعاد بیش تر ، به صورت خطی تفکیک پذیر باشد بیش تر از احتمال تفکیک پذیری خطی این مساله در فضایی با ابعاد کم تر است .
- فرض کنید که $\Phi: X \rightarrow H$ تبدیلی غیر خطی از فضای ورودی X به فضای خصیصه H باشد که در آن مساله به صورت خطی تفکیک پذیر است.



- انجام محاسبات در فضای ویژگی می تواند پرهزینه باشد برای اینکه ابعاد بیشتری دارد.
- در حالت کلی ابعاد این فضا بی نهایت است.
- برای غلبه بر این مشکل از kernel trick استفاده میشود.

کرل ضرب داخلی: کرل، برای تمامی $x, x' \in X$ یک تابع $K(x, x')$ بوده و به صورت زیر تعریف می شود:

$$K(x, x') = \Phi^T(x) \Phi(x')$$

- ✓ علاوه بر توابع کرل خطی می توان از توابع کرل چند جمله ای یا سیگموئید نیز استفاده کرد. در سال های اخیر تحقیقات زیادی درباره توابع کرل مختلف انجام شده است.
- ✓ به دلیل پیچیدگی مسائل **Kernel trick** همیشه تضمین نمی کند که آن ها بطور کامل تفکیک پذیر خطی باشند، بنابراین در اغلب اوقات ایده های مربوط به **حاشیه انعطاف پذیر** و **Kernel trick** را یکپارچه ساخته تا از مزایای هر دو رویکرد بهره مند شویم.
- ✓ در این حالت فرم دوگان مساله بهینه سازی به صورت زیر خواهد بود:

$$\begin{aligned} \max_{\alpha} W(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{\gamma} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{Subject to:} \\ \sum_{i=1}^n \alpha_i y_i &= 0 \\ 0 \leq \alpha_i &\leq C \quad i = 1, 2, \dots, n \end{aligned}$$

با استفاده از روش ضرایب لاگرانژ دسته بندی کننده بهینه را به دست می آوریم:

$$f(x) = \sum_{i=1}^n \alpha_i^* y_i K(x_i, x) + b^*$$

$$b^* = 1 - \sum_{i=1}^n \alpha_i^* y_i K(x_i, x_s)$$

در نتیجه

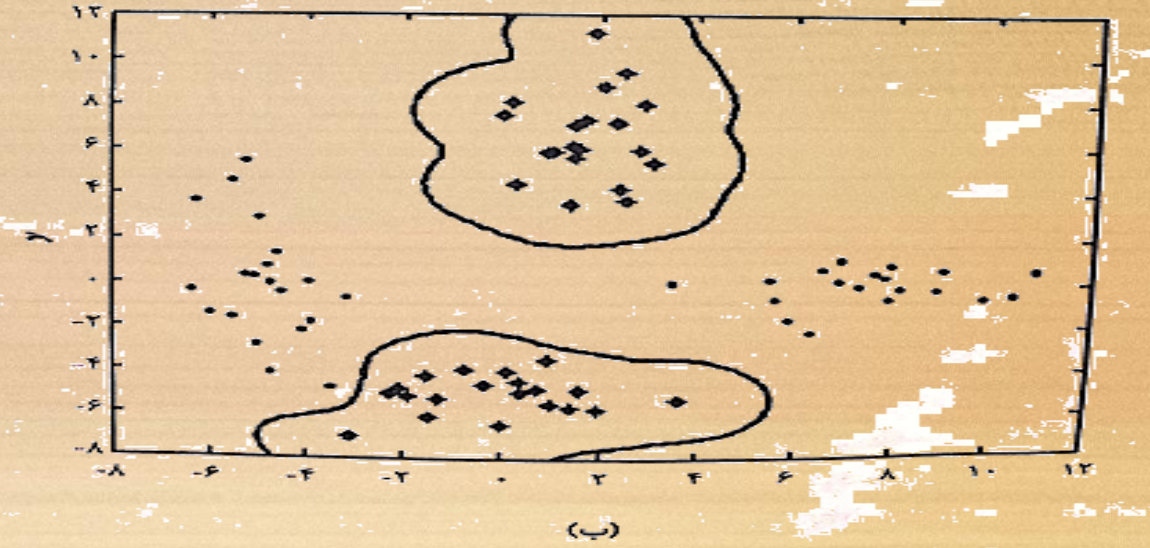
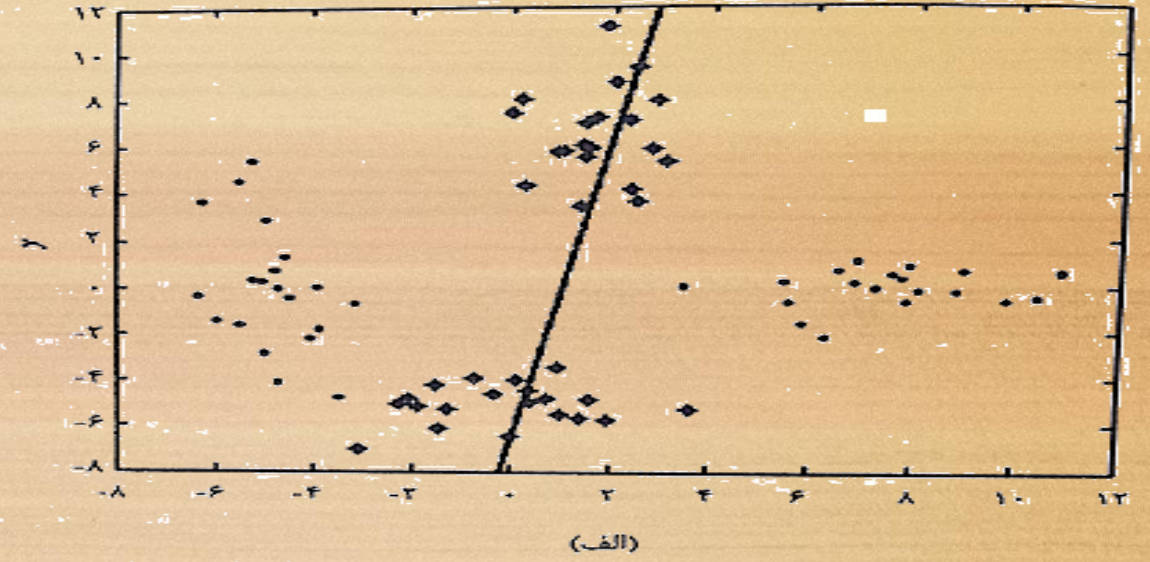
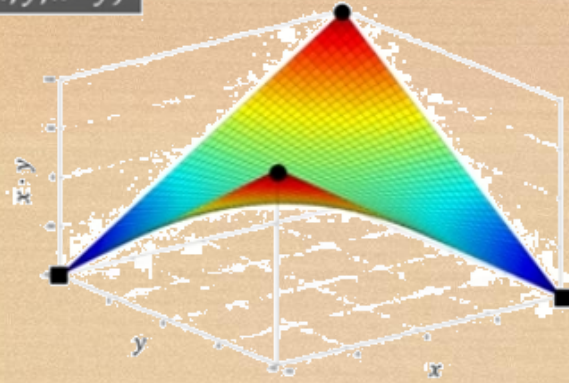
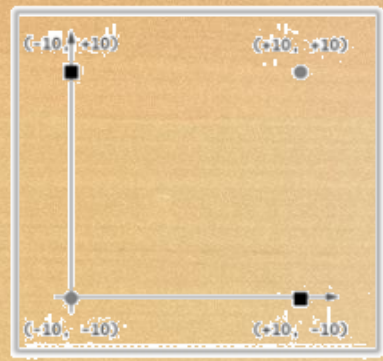
➤ مثال:

مساله XOR یکی از مسائل دسته بندی بوده که به صورت خطی تفکیک پذیر نمی باشد . به منظور نمایش اهمیت یکپارچه سازی SVC با **حاشیه انعطاف پذیر** و **Kernel trick** در حل مسائل پیچیده دسته بندی از این مساله استفاده می کنیم. یک مجموعه داده XOR دو بعدی را می توان به تصادف و توسط چهار توزیع گوسی متفاوت تولید کرد . " * " و " • " نمونه های متعلق به دو کلاس را نمایش می دهند . همانگونه که در شکل **الف** نشان داده شده است SVC مقدماتی و کرنل خطی در حل مساله XOR با شکست کامل مواجه می شوند. یک مرز خطی توانایی تفکیک دو کلاس مورد نظر را ندارد . همانگونه که ملاحظه می کنید این مرز خطی تمامی نمونه ها را به دو قسمت تقسیم کرده است . بنابراین این رویکرد نمی تواند هدف مورد نظر از دسته بندی را برآورده سازد . در نتیجه از ترکیب SVC با **حاشیه انعطاف پذیر** و **کرنل با پایه شعاعی گوسی** به منظور حل این مساله استفاده می کنیم.

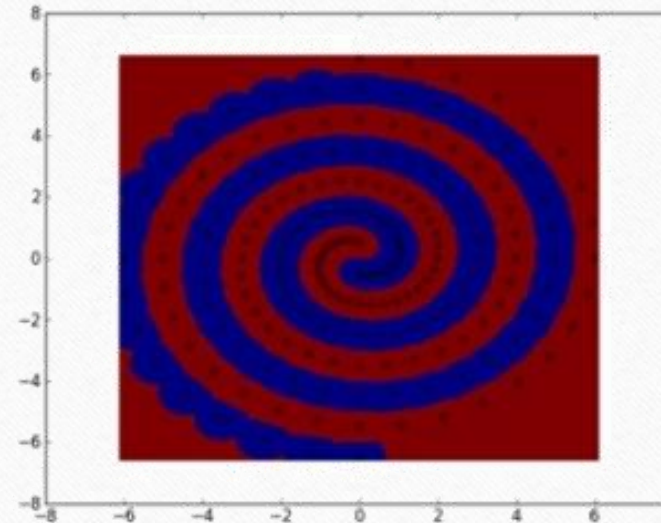
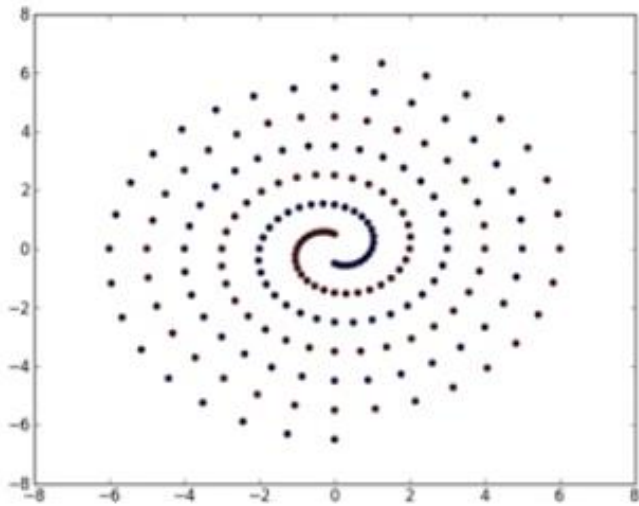
$$K(x_i, x) = \exp\left(-\frac{\|x - x_i\|^2}{2\sigma^2}\right)$$

پارامتر C را برابر با 1 و پارامتر کرنل را مساوی با 1 قرار می دهیم ($\sigma = 1$) . مرز تفکیک کننده متناظر را در شکل ب ملاحظه می فرمایید با استفاده از **Kernel trick** مرز دیگر خطی نبوده و تنها یک کلاس را شامل می شود. با بررسی نمونه هایی که در داخل و خارج از مرز قرار دارند در می یابیم که دسته بندی کننده نمونه ها را به درستی دسته بندی کرده است.

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



مثال تابع کرنل گوسی:



مقدار مناسب برای پارامترهای تابع کرنل چگونه تعیین می‌شود؟

□ 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

➤ ماشین‌های بردار پشتیبانی و شبکه‌های عصبی:

- توسعه شبکه‌های عصبی قبل از آنکه بر اساس اصول نظری انجام گرفته باشد، بر اساس روش‌های ابتکاری و تعداد زیادی از کاربردها و آزمایش‌ها صورت گرفته است ولی ماشین‌های بردار پشتیبان بر عکس.
- از مزیت‌های SVM ها آن است که جواب یک SVM سراسری و یکتا است، درحالی‌که شبکه‌های عصبی عمدتاً از چندین کمینه محلی رنج می‌برند.
- برخلاف شبکه‌های عصبی، پیچیدگی محاسباتی مربوط به ماشین‌های بردار پشتیبانی به ابعاد فضای ورودی بستگی ندارد.
- یکی از دلایلی که به موجب آن عملکرد ماشین‌های بردار پشتیبانی در اغلب موارد از عملکرد شبکه‌های عصبی بهتر است، کمتر رخ دادن مشکل برازش بیش از اندازه در SVM ها است.
- در اغلب موارد از کرنل‌های گوسی استفاده می‌شود.

➤ نقاط ضعف ماشین‌های بردار پشتیبانی:

- با وجود آنکه SVM ها از دیدگاه عملی و کاربردی دارای مزیت‌های قابل توجهی هستند، اما محدودیت‌هایی نیز دارند.
- یکی از پرسش‌هایی که هنوز پاسخ کاملی برای آن نیست، نحوه انتخاب پارامترهای تابع کرنل است.
- محدودیت دوم، در رابطه با اندازه و سرعت است (هم در آموزش و هم در آزمایش).
- پردازش داده‌های گسسته یکی دیگر از مشکلات است.

با این وجود، SVM ها دارای یک شالوده نظری منسجم بوده و جواب‌های ایجاد شده توسط آنها، سراسری و یکتا است. امروزه ماشین‌های بردار پشتیبانی به متداول‌ترین تکنیک‌های پیش‌بینی در داده‌کاوی تبدیل و در بیشتر ابزارهای تجاری داده‌کاوی بکار گرفته شده‌اند.