



توابع جداساز خطی

در این فصل توابع جداساز خطی را مورد بررسی قرار خواهیم داد. توابع جداساز خطی شامل دو دسته مجزای مبتنی بر فواصل و مبتنی بر توزیع های آماری می باشند. از آنجایی که توابع جداساز آماری در فصل جداگذاری بررسی شده اند، بحث در مورد توابع جداساز خطی آماری نیز به انجام خواهد شد.

تعریف جداساز انواع آنها، پرسپکتیو، جداساز خطی فیشر، ماشین های بردار پشتیبان خطی و روش های پراکنده سازی خطی را به گونه ای که در این فصل به طور مفصل مورد بحث و بررسی قرار گرفته اند.

یک جداساز، تابعی است که هر بردار ویژگی x را به یکی از k کلاس موجود نسبت می دهد. جداسازها از یک دید به دو دسته قطعی و غیر قطعی (آماری) تقسیم می شوند. جداسازهای قطعی بر پایه فواصل و جداسازهای آماری بر پایه توزیع های احتمالاتی عمل می کنند.

در دسته بندی بر پایه فواصل، از این واقعیت استفاده می شود که «بردارهای ویژگی نزدیک به هم در فضای ویژگی، بیان گر اشیای شبیه به هم در دنیای واقعی می باشند». به عنوان مثال در یک مساله دسته بندی دو کلاسه، اگر t_1 و t_2 نماینده های دو کلاس C_1 و C_2 باشند، بردار ویژگی x به کلاس C_1 تعلق دارد اگر فاصله x از t_1 کمتر از فاصله x از t_2 باشد و بردار ویژگی x به کلاس C_2 تعلق دارد اگر فاصله x از t_2 کمتر از فاصله x از t_1 باشد. یعنی،

$$\|x - t_1\| \leq_{C_2}^{\leq_{C_1}} \|x - t_2\| \quad (1)$$

با توجه به اینکه $\|x - t_i\|^2 = (x - t_i)^t (x - t_i)$ ، خواهیم داشت:

$$(x - t_1)^t (x - t_1) \leq_{C_2}^{\leq_{C_1}} (x - t_2)^t (x - t_2) \Rightarrow t_1^t t_1 - 2x^t t_1 \leq_{C_2}^{\leq_{C_1}} t_2^t t_2 - 2x^t t_2$$

و با در نظر گرفتن $g_i(x) = t_i^t t_i - 2x^t t_i$ ، رابطه (1) به فرم زیر در می آید:

$$g_1(x) \leq_{C_2}^{\leq_{C_1}} g_2(x)$$

که می توان آن را به صورت زیر نیز بیان نمود:

$$g(x) = (g_2(x) - g_1(x)) \geq_{C_2}^{\leq_{C_1}} 0 \quad (2)$$

¹ Linear Discriminan Functions (LDF)

این تابع برای داده‌های کلاس اول مقدار بزرگتر از صفر و برای داده‌های کلاس دوم مقدار کوچکتر از صفر را بر می‌گرداند. از آنجایی که از تابع $g(x)$ برای تفکیک دو کلاس استفاده می‌شود، به آن تابع جداساز گفته می‌شود. به توابعی جداسازی که همانند بردار ویژگی فوق از درجه اول باشند، توابع جداساز خطی گفته می‌شود.

برای $x = (x_1, x_2, \dots, x_d)$ ، یک بردار ویژگی در یک فضای d بعدی، تابع جداساز خطی به یکی از صورت‌های زیر می‌باشد:

$$g(x) = w^t x + w_0 \quad (3)$$

که در آن $w = (w_1, w_2, \dots, w_d)$ می‌باشد. فرم فوق را فرم معمولی توابع جداساز خطی می‌نامند. پس در دسته‌بندی دو کلاسه بوسیله توابع جداساز، باید تابعی ارائه کنیم که بوسیله آن بتوان، بطور قطعی کلاس هر بردار ویژگی جدید را تعیین نمود. برای ارائه تابع جداساز نیز باید پارامترهای (w, w_0) مشخص شوند.

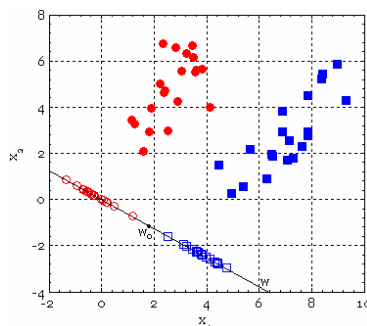
اگر بردارهای ویژگی را به صورت افزوده^۲ در نظر بگیریم، یعنی $x = (1, x_1, x_2, \dots, x_d)$ و همچنین برای بردار وزن‌ها نیز داشته باشیم $w = (w_0, w_1, w_2, \dots, w_d)$ ، رابطه (۳) به فرم زیر در می‌آید:

$$g(x) = w^t x \quad (4)$$

این فرم توابع جداساز را فرم همگن^۳ می‌نامند. در حل تحلیلی مسایل، کار کردن با فرم همگن توابع جداساز، راحت‌تر است. مطابق رابطه (۴)، اگر بردار ویژگی x متعلق به کلاس C_1 باشد، $w^t x > 0$ و اگر بردار ویژگی x متعلق به کلاس C_2 باشد، $w^t x < 0$ یا $w^t(-x) > 0$ خواهد بود.

اگر تمامی بردارهای ویژگی کلاس C_2 را قرینه کنیم، برای تمامی بردارهای ویژگی متعلق به هر دو کلاس $g(x) = w^t x > 0$ خواهد شد. با این کار فرم نرمال تابع جداساز خطی بدست می‌آید. استفاده از فرم نرمال در حل مسایل، در بسیاری از حالات، محاسبات را ساده‌تر می‌کند.

در فرم معمولی یک تابع جداساز خطی، جزء $w^t x$ بیان می‌کند که داده‌ها باید بر روی ابر صفحه با بردار مشخصه w تابانده شوند. بر روی این ابر صفحه نیز w_0 مرز جداسازی را تعیین می‌کند. این شهود در شکل ۱، برای حالت دوبعدی (که ابر صفحه برابر خط خواهد بود) نمایش داده شده است.



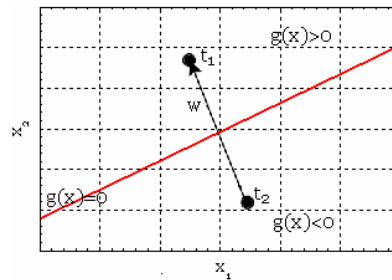
شکل ۱: شهود هندسی دسته‌بندی با توابع جداساز خطی

در یک تابع جداساز، حالت تساوی با صفر را بررسی نکردیم. اگر $g(x) = 0$ (در تمامی فرم‌های نمایش)، در این صورت بردار ویژگی x به کدام کلاس تعلق خواهد داشت؟ واقعیت این است که در این مورد نمی‌توان به درستی اظهار نظر کرد، چون این بردار ویژگی بر روی مرز بین دو کلاس قرار گرفته است. در حالت کلی به مجموعه تمامی نقاطی که به ازای آنها $g(x) = w^t x + w_0 = 0$ می‌باشد، مرز جداسازی (مرز تصمیم‌گیری) گفته می‌شود.

برای یک تابع جداساز خطی، مرز جداسازی در یک فضای دو بعدی، یک خط، در یک فضای سه بعدی یک صفحه و در حالت کلی در یک فضای d بعدی، یک ابر صفحه $(d-1)$ بعدی خواهد بود. مفاهیم مطرح شده فوق در شکل ۲ در یک فضای ویژگی دو بعدی به تصویر کشیده شده‌اند:

² Augmented Feature Vector

³ Homogeneous



شکل ۲: مرز جداسازی در فضای دو بعدی

در این شکل، t_1 و t_2 مراکز دو کلاس C_1 و C_2 بوده و خط واصل t_1 و t_2 می‌باشد. خط قرمز رنگ مرز جداسازی دو کلاس (عمود بر w) می‌باشد. یک طرف خط متعلق به کلاس C_1 بوده و در آن $g(x) > 0$ می‌باشد (بالای خط مرز در شکل). طرف دیگر خط نیز متعلق به کلاس C_2 بوده و در آن $g(x) < 0$ می‌باشد (پایین خط مرز در شکل).

مفاهیم دسته‌بندی خطی را برای حالت دو کلاسه مورد بررسی قرار دادیم. در بسیاری از مسایل، تعداد کلاس‌ها ممکن است بیش از دو کلاس باشد. در این حالت برای دسته‌بندی می‌توان به دو صورت عمل نمود. یکی اینکه، برای جداسازی بردارهای ویژگی هر کلاس از سایر کلاس‌ها، یک تابع جداساز ارائه نمود. در این حالت برای دسته‌بندی n کلاس به n تابع جداساز نیاز است (البته دسته‌بندی با $n-1$ تابع جداساز نیز قابل انجام است. چگونه؟). و دیگر اینکه برای جداسازی بردارهای ویژگی یک کلاس از هر کلاس دیگر، یک تابع جداساز داشته باشیم. در این حالت برای دسته‌بندی n کلاس به $n(n-1)/2$ تابع جداساز نیاز است. در حالت اول مجموعه بردارهای ویژگی را جدپذیر خطی کامل^۴ و در حالت دوم مجموعه بردارهای ویژگی را جدپذیر خطی دویبدو^۵ گویند.

مثال ۱: یک مساله دسته‌بندی سه کلاسه را در نظر بگیرید. در این مساله برای هر کلاس یک تابع معرف به صورت زیر ارائه شده است:

$$g_1(x) = -x_1 + x_2, \quad g_2(x) = x_1 + x_2 - 1, \quad g_3(x) = -x_2$$

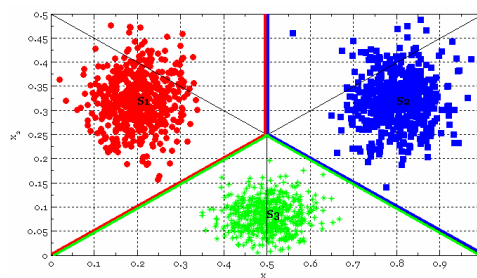
برای تشخیص تعلق یک داده به یک کلاس نیز از قانون زیر استفاده می‌شود:

$$x \in C_i \Leftrightarrow g_i(x) < g_j(x); \forall j \neq i$$

هر جداساز با فرم فوق را یک ماشین خطی می‌نامند. این ماشین خطی ویژگی را به چه صورت افراز می‌کند؟

در این مثال ناحیه مربوط به هر کلاس، مطابق زیر، با دو نامساوی خطی مشخص می‌شود:

ناحیه ۳	ناحیه ۲	ناحیه ۱
$\underline{x} \in S_3 : \begin{cases} g_3 > g_1 \Rightarrow x_2 < x_1 / 2 \\ g_3 > g_2 \Rightarrow x_2 < 1/2(-x_1 + 1) \end{cases}$	$\underline{x} \in S_2 : \begin{cases} g_2 > g_1 \Rightarrow x_1 > 1/2 \\ g_2 > g_3 \Rightarrow x_2 > 1/2(-x_1 + 1) \end{cases}$	$\underline{x} \in S_1 : \begin{cases} g_1 > g_2 \Rightarrow x_1 < 1/2 \\ g_1 > g_3 \Rightarrow x_2 > x_1 / 2 \end{cases}$



⁴ Totally Linearly Seperable
⁵ Pairwise Linearly Seperable

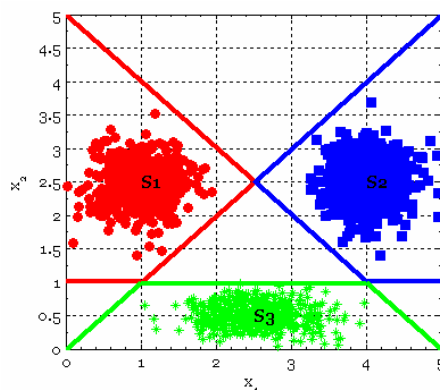
مثال ۲: یک مساله دسته‌بندی با سه کلاس را در نظر بگیرید که در آن بردارهای ویژگی جداپذیر خطی کامل هستند. توابع جداساز را نیز به صورت زیر در نظر بگیرید:

$$g_1(x) = -x_1 + x_2, \quad g_2(x) = x_1 + x_2 - 5, \quad g_3(x) = -x_2 + 1$$

توابع جداساز به قسمی ارائه شده‌اند که $g_i(x) > 0$ نمایان‌گر تعلق x به کلاس i ام و $g_i(x) < 0$ نمایان‌گر عدم تعلق آن به کلاس i ام باشد. این جداساز فضای ویژگی را به چه صورت افراز می‌کند؟ نواحی‌ای که جداساز قادر به تصمیم‌گیری در مورد آنها نیست، را نیز مشخص نمایید.

هر ناحیه با سه نامساوی خطی روی x تعریف می‌شود (یعنی هر ناحیه بین سه خط محصور است).

ناحیه ۳	ناحیه ۲	ناحیه ۱
$x \in S_3 : \begin{cases} -x_1 + x_2 < 0 \\ x_1 + x_2 - 5 < 0 \\ -x_2 + 1 > 0 \end{cases}$	$x \in S_2 : \begin{cases} -x_1 + x_2 < 0 \\ x_1 + x_2 - 5 > 0 \\ -x_2 + 1 < 0 \end{cases}$	$x \in S_1 : \begin{cases} -x_1 + x_2 > 0 \\ x_1 + x_2 - 5 < 0 \\ -x_2 + 1 < 0 \end{cases}$



نواحی‌ای نیز که مقدار بیش از یک تابع برای داده‌های آنها، مثبت باشد، نواحی غیر قابل تصمیم‌گیری می‌باشند (شامل چهار ناحیه مثلثی در شکل).

□

مثال ۳: توابع جداساز، برای یک جداساز جداپذیر خطی دوبعدی سه کلاسه، را به صورت زیر در نظر بگیرید:

$$g_{12}(x) = -x_1 - x_2 + 5, \quad g_{13}(x) = -x_1 + 3, \quad g_{23}(x) = -x_1 + x_2$$

(و $g_{ij}(x) = -g_{ji}(x)$). در این جداساز، قانون تصمیم‌گیری به صورت زیر می‌باشد:

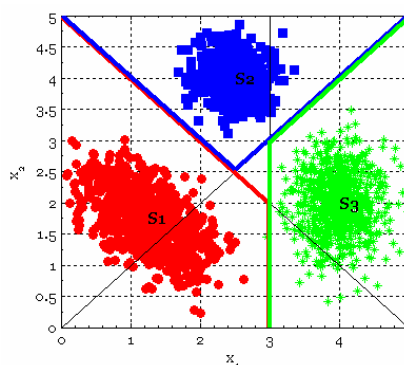
$$x \in C_i \Leftrightarrow \forall j \neq i \quad g_{ij}(x) > 0$$

این جداساز فضای ویژگی را به چه صورت افراز می‌کند؟ نواحی‌ای که جداساز قادر به تصمیم‌گیری در مورد آنها نمی‌باشد را مشخص نمایید.

هر ناحیه با دو نامساوی خطی مشخص می‌شود که هر نامساوی متناظر با جداسازی کلاس مورد نظر از کلاس‌های دیگر است. در نتیجه هر ناحیه با دو خط محصور است.

ناحیه ۳	ناحیه ۲	ناحیه ۱
$x \in S_3 : \begin{cases} x_1 - 3 > 0 \\ x_1 - x_2 > 0 \end{cases}$	$x \in S_2 : \begin{cases} x_1 + x_2 - 5 > 0 \\ -x_1 + x_2 > 0 \end{cases}$	$x \in S_1 : \begin{cases} -x_1 - x_2 + 5 > 0 \\ -x_1 + 3 > 0 \end{cases}$

همانطور که در شکل مشخص است، ناحیه مثلثی وسط، در هیچ یک از کلاس‌ها قرار نمی‌گیرد. زیرا در آن فقط $g_{13}(x)$ بزرگ‌تر از صفر می‌باشد.



□

تمرین ۱: نشان دهید یک ماشین خطی، حالت خاصی از جداساز جدپذیر خطی دودو می‌باشد.

تمرین ۲: نشان دهید در یک مساله دسته‌بندی چند کلاسه با یک ماشین خطی، نواحی بدست آمده برای هر کلاس (ناحیه حاصل از مرزهای جداساز آن کلاس) محدب می‌باشند. (راهنمایی: برای نشان دادن محدب بودن یک ناحیه کافی است نشان دهید که $\lambda x_1 + (1-\lambda)x_2 \in R_i$ if $0 \leq \lambda \leq 1$ and $x_1 \in R_i$ and $x_2 \in R_i$ چرا؟)

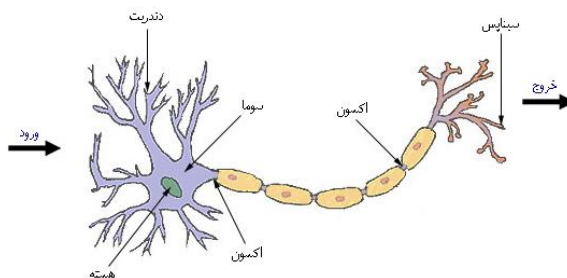
نحوه بدست آوردن توابع
جداساز برای مسائل مختلف

مساله اصلی در بحث توابع جداساز خطی، نحوه بدست آوردن توابع به ازای هر مساله خاص می‌باشد. در ادامه این مبحث به روش‌های مختلف بدست آوردن توابع جداساز خطی می‌پردازیم. ابتدا نرون‌های عصبی مصنوعی را که شکل دیگری از نمایش توابع جداساز می‌باشند، معرفی نموده و الگوریتم‌های پرسپترون و Relaxation را، که بر پایه کاهش گرادیان می‌باشند، برای بدست آوردن پارامترهای آنها توضیح می‌دهیم. سپس به جداساز خطی فیشر می‌پردازیم که بر اساس تابش بردارهای ویژگی در فضای کاهش یافته جدیدی، که جداسازی بردارهای ویژگی را آسان‌تر می‌کند، عمل می‌کند. پس از آن به ماشین‌های بردار پشتیبان خطی می‌پردازیم که جداسازهایی با بیشترین حاشیه را می‌یابند. پیدا کردن توابع جداساز با روش‌های بر پایه کمینه‌سازی خطا نیز بعد از آن مورد بررسی قرار می‌گیرند که شامل روش‌های تخمین میانگین مربعات خطا، الگوریتم Widrow-Hoff، روش ماتریس شبه معکوس و متد Ho-Kahyap می‌باشند. مقایسه روش‌ها و چند بحث تکمیلی نیز پایان بخش این مبحث خواهند بود.

۱- پرسپترون

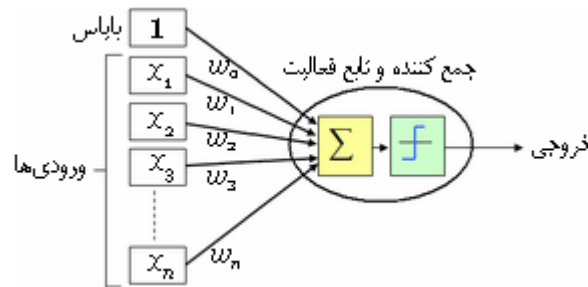
پرسپترون و نرون‌های عصبی

پرسپترون و دیگر شبکه‌های عصبی مصنوعی الهام گرفته از ساختار عصبی بدن انسان هستند. کوچکترین جز در این ساختار نرون‌ها می‌باشند. ساختار یک نرون طبیعی در شکل ۳ نشان داده شده است. در این ساختار بدنه نرون وظیفه جذب پتانسیم توسط دندریت‌ها را دارد. این پتانسیم در بدنه ذخیره شده و هر وقت مقدار آن از حد معینی تجاوز کند، بدنه از طریق سیناپس‌ها شروع به پمپاژ پتانسیم به بیرون می‌کند. پتانسیم پمپاژ شده در فضای سیناپسی توسط دندریت‌های سایر نرون‌هایی که در فضای مجاور قرار دارند، جذب می‌شود.



شکل ۳: ساختار یک نرون طبیعی در شبکه عصبی موجودات زنده

در طراحی پرسپترون نیز از همین ایده الهام گرفته شده است. یک پرسپترون برداری از ورودی‌های با مقادیر حقیقی را گرفته و یک ترکیب خطی از این ورودی‌ها را محاسبه می‌کند. اگر حاصل از یک مقدار آستانه بیشتر بود، مقدار ۱ و در غیر این صورت مقدار ۰- (یا صفر) را به خروجی خود می‌فرستد (شکل ۴).



شکل ۴: ساختار یک نرون مصنوعی (پرسپترون)

پرسپترون به هر کدام از ورودی‌های خود، وزن مشخصی می‌دهد. اگر دارای n ورودی x_1, x_2, \dots, x_n باشیم و وزن‌های معادل این ورودی‌ها w_1, w_2, \dots, w_n باشند، در این صورت ورودی پرسپترون برابر خواهد بود با:

$$w_1x_1 + w_2x_2 + \dots + w_nx_n = \sum_{i=1}^n w_i x_i$$

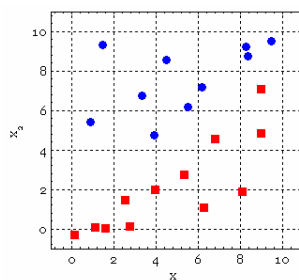
پرسپترون بر اساس بیشتر شدن یا نشدن ورودی‌اش از حد مشخصی (w_0)، خروجی‌های خود را تولید خواهد کرد. این مقایسه را می‌توان بصورت نشان داد. و اگر ورودی دیگری با مقدار ۱ و وزن w_0 برای پرسپترون قائل شویم، مقایسه بصورت $\sum_{i=0}^n w_i x_i \geq 0$ (فرم همگن) در خواهد آمد. در نتیجه تابع یک پرسپترون بصورت زیر مشخص می‌شود:

$$f(x_1, x_2, \dots, x_n) = \begin{cases} 1 & \text{if } \sum_{i=0}^n w_i x_i \geq 0 \\ -1 & \text{otherwise} \end{cases} \quad (5)$$

حال که ساختار پرسپترون را بررسی کردیم، نحوه استفاده از آن را برای حل مسایل دسته‌بندی بررسی می‌کنیم. یک مساله دسته‌بندی دو کلاسه در یک فضای d بعدی را در نظر بگیرید. و همچنین فرض کنید برای طراحی یک جداساز مناسب برای این مساله، n بردار ویژگی (داده‌های آموزشی) در اختیار داریم. این مساله چگونه با یک پرسپترون حل می‌شود؟

جداساز مورد نیاز برای حل این مساله، یک پرسپترون با $d+1$ ورودی (d ورودی برای دریافت ابعاد مختلف هر کدام از بردارهای ویژگی و یک ورودی نیز برای مشخص کردن حد آستانه) و یک خروجی با مقادیر ۱ یا -۱ (برای نشان دادن تعلق به کلاس اول یا دوم) می‌باشد. وزن‌های نامشخص پرسپترون (w_0, w_1, \dots, w_d) باید با توجه به n داده آموزشی و ساختار تعیین شده، مشخص شوند. با مشخص شدن این وزن‌ها، تمامی قسمت‌های پرسپترون شناخته شده و تابع پرسپترون قادر خواهد بود که برای هر بردار ویژگی، کلاس آن را تعیین نماید.

مثال ۴: یک مساله دسته‌بندی دو کلاسه در یک فضای دو بعدی را با بردارهای ویژگی آموزشی نشان داده در شکل زیر در نظر بگیرید. یک پرسپترون برای حل این مساله طراحی نمائید.



برای حل این مساله نیاز به یک پرسپترون با سه ورودی و یک خروجی داریم. اگر خروجی پرسپترون صفر باشد، نشان دهنده

یکی از کلاس‌ها (مثلاً مربع‌ها) و اگر ۱ باشد، نشان دهنده کلاس دیگر (دایره‌ها) می‌باشد. برای طراحی کامل پرسپترون مربوطه، تنها مجهول‌های باقیمانده، وزن‌ها (w_2, w_1, w_0) می‌باشند. تابع پرسپترون برای این مثال نیز بصورت زیر می‌باشد:

$$f(x_1, x_2) = \begin{cases} 1 & \text{if } w_2x_2 + w_1x_1 - w_0 \geq 0 \\ -1 & \text{otherwise} \end{cases} \quad (۶)$$

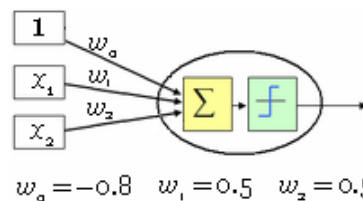
از طرفی، با توجه به شکل، می‌دانیم که معادله خط جداساز دو کلاس $x_2 - x_1 = 0$ بوده و بالای خط $(x_2 - x_1 > 0)$ ، دایره‌ها (کلاس ۱) و پایین خط $(x_2 - x_1 < 0)$ ، مربع‌ها (کلاس -۱) قرار گرفته‌اند.

با معادل‌سازی $w_2x_2 + w_1x_1 - w_0 \geq 0$ و $x_2 - x_1 > 0$ و همچنین معادل‌سازی $w_2x_2 + w_1x_1 - w_0 < 0$ و $x_2 - x_1 < 0$ به این نتیجه می‌رسیم که $(w_2, w_1, w_0) = (1, -1, 0)$. این نتیجه طراحی پرسپترون ما را کامل می‌کند.

□

پرسپترون و شبیه‌سازی
عملگرهای منطقی

یک پرسپترون می‌تواند بسیاری از توابع منطقی نظیر and, or, nand و nor را شبیه‌سازی. به عنوان مثال پرسپترون زیر عملکرد and را شبیه‌سازی می‌کند (خروجی پرسپترون برابر x_1 and x_2 می‌باشد).



شکل ۵: شبیه‌سازی عملکرد and توسط پرسپترون

تمرین ۳: توابع and, or, nand و nor را با پرسپترون شبیه‌سازی کنید.

یک پرسپترون فقط قادر به دسته‌بندی خطی است. مسائلی که توسط یک ابر صفحه قابل جداسازی باشند، در نتیجه قادر به شبیه‌سازی توابعی نظیر XOR نیز نمی‌باشد. در فصل‌های آتی، شبکه‌های عصبی چند لایه را بررسی خواهیم کرد که قادر به حل مسائل غیر خطی نیز می‌باشند.

تعیین سیستماتیک وزن‌های
پرسپترون

همانطور که مشاهده کردید، در مثال‌های فوق برای تعیین وزن‌ها در طراحی پرسپترون از دانش شهودی خود استفاده کردیم. این دانش شهودی همیشه وجود ندارد یا کسب آن به راحتی امکان‌پذیر نیست. پس باید بدنبال روش مناسبی (یک الگوریتم) برای تعیین وزن‌ها، بصورت سیستماتیک باشیم.

اگر مساله را به شکل نرمال در نظر بگیریم، هدف از این الگوریتم این است که $w = (w_0, w_1, \dots, w_d)$ را طوری بیاید که به ازای همه بردارهای ویژگی ورودی (x) ، داشته باشیم

$$g(x) = w^T x > 0 \quad (۷)$$

روش‌های تکراری برای حل
مسایل بهینه‌سازی

یکی از روش‌های پیدا کردن w استفاده از الگوریتم‌های بهینه‌سازی تکراری می‌باشد. برای این کار، ابتدا w را بطور تصادفی انتخاب کرده و میزان خوبی آن را می‌سنجیم. اگر w بردار خوبی بود، جواب مساله پیدا شده است، ولی اگر به اندازه کافی خوب نبود، باید آن را طوری تغییر دهیم که عملکرد تابع جداساز معادل آن بهتر شود. برای بردار جدید بدست آمده نیز، باید همین عملیات آزمون و تغییر را به صورت متوالی انجام دهیم، تا جایی که بردار خوبی بدست آید.

روش پیشنهاد شده، از لحاظ تئوری روش مناسبی می‌باشد. ولی جهت اجرا، باید بر دو مشکل غلبه کنیم. یکی اینکه باید روشی برای تغییر بردار w در هر مرحله، به قسمی بیابیم که بردار جدید، بهتر از بردار قبل باشد و دیگر اینکه باید معیاری برای سنجش خوبی w ارائه دهیم.

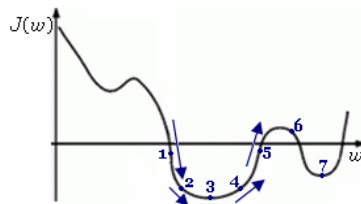
روش کاهش گرادیان

برای تولید بردار جدید از بردار قبلی از روش کاهش گرادیان^۶ استفاده می‌کنیم. در این روش، هدف، پیدا کردن بردار بهینه $w = (w_0, w_1, \dots, w_d)$ است که تابع هزینه‌ای نظیر $J(w)$ را کمینه کند. در نتیجه برای بردار بهینه w باید داشته باشیم:

^۶ Gradient Descent

$$\nabla J = \frac{\partial J}{\partial w} = \left[\frac{\partial J}{\partial w_0}, \frac{\partial J}{\partial w_1}, \dots, \frac{\partial J}{\partial w_d} \right] = 0$$

در شکل ۶ جهت و مقدار گرادیان $\nabla J(w)$ برای یک w یک بعدی در نقاط مختلف نشان داده شده است. در نقطه کمینه مقدار گرادیان برابر صفر می‌باشد (نقطه ۳). در نقاطی که قبل از نقطه کمینه قرار دارند، جهت گرادیان (جهت شیب خط مماس بر آن نقاط) منفی بوده (نقاط ۱ و ۲) و در نقاطی که بعد از نقطه کمینه قرار دارند، جهت گرادیان مثبت می‌باشد (نقاط ۴ و ۵). مقدار گرادیان نیز در نقاطی که شیب تندتری داریم، بیشتر است (صرفنظر از علامت مقادیر، اندازه گرادیان در نقاط نمایش داده شده در شکل، به ترتیب نزولی در نقاط ۱، ۴، ۵، ۲ و ۳ می‌باشد).



شکل ۶: مقدار و جهت گرادیان در نقاط دور و بر یک نقطه کمینه

با توجه به مطالب ذکر شده، اگر در مساله اصلی، در یک مرحله، w مقدار $J(w)$ را کمینه نکند، باید مقدار آن به سمت بهبود (به سمت w بهینه) تغییر پیدا کند. با توجه به تعبیر هندسی گرادیان (و با توجه به مفاهیم نشان داده شده در شکل ۶) میزان این تغییرات باید در راستای $-\nabla J$ و ضریب مثبتی از مقدار گرادیان باشد. یعنی،

$$w_{new} = w_{old} - \eta \nabla J(w) \quad (۸)$$

η را نرخ یادگیری^۷ گوئیم که می‌تواند در مراحل مختلف مقادیر متفاوتی را به خود بگیرد (هر چند که در بسیاری از موارد نرخ یادگیری، ثابت در نظر گرفته می‌شود). اگر η مقدار کوچکی داشته باشد، تعداد قدم‌ها تا رسیدن به هدف بسیار زیاد خواهد شد (زیرا در هر مرحله مقدار خیلی کمی به هدف نزدیک می‌شویم) و اگر η مقدار بزرگی داشته باشد (قدم بزرگی به سمت هدف برداشته شود)، الگوریتم ممکن است قادر به پیدا کردن کمینه عمومی نبوده و یک کمینه محلی را پیدا کند. و یا اینکه اصلاً نتواند نقطه بهینه‌ای را بیابد. مثلاً در شکل ۶، فرض کنید در نقطه ۲ قرار داریم. می‌دانیم که برای نزدیک شدن به هدف باید به سمت راست حرکت کنیم. اگر نرخ یادگیری مقدار بزرگی باشد، قدم بزرگی به سمت راست برداشته می‌شود و مثلاً به نقطه ۶ می‌رسیم. در مراحل بعدی، نقطه کمینه عمومی ۳ توسط الگوریتم قابل کشف نیست و نقطه کمینه محلی ۷ پسته به مقدار نرخ یادگیری در آن مراحل، ممکن است یافت بشود یا نشود (الگوریتم با یک قدم بزرگ دیگر خیلی راحت می‌تواند از این نقطه نیز رد شود).

با توجه به مطالب گفته شده، می‌توان الگوریتم کاهش گرادیان را به صورت زیر بیان نمود:

الگوریتم کاهش گرادیان برای بهینه‌سازی تابع $J(w)$

الگوریتم کاهش گرادیان برای

مسائل بهینه‌سازی

الف) یک مقدار اولیه برای w در نظر بگیر $(v_{(0)})$ و متغیر k را نیز با مقدار اولیه صفر در نظر بگیر.

ب) تا وقتی که $\nabla J(w) \neq 0$ (یا در بسیاری از مسائل تا وقتی که $\nabla J(w) > \varepsilon$) مراحل زیر را تکرار کن:

ب(۱) مقداری مثبت برای $\eta_{(k)}$ در نظر بگیر.

ب(۲) مقدار w را مطابق قانون $w_{(k+1)} = w_{(k)} - \eta_{(k)} \nabla J(w)$ بروز کن.

ب(۳) مقدار k را یکی اضافه کن.

برای همگرا شدن الگوریتم (متوقف شدن الگوریتم و یافتن w بهینه)، نرخ یادگیری، علاوه بر مثبت بودن، باید دو شرط زیر را

⁷ Learning Rate

نیز دارا باشد:

$$\sum_{k=1}^{\infty} \eta_{(k)} \rightarrow \infty, \sum_{k=1}^{\infty} (\eta_{(k)})^2 < \infty$$

تمرین ۴: نرخ یادگیری با دارا بودن شروط همگرایی، مقادیر در چه بازه‌هایی را می‌تواند بپذیرد؟

تمرین ۵: اگر در یک مساله دسته‌بندی ندانیم که بردارهای ویژگی به صورت خطی جداپذیر هستند یا نه، بهتر است از ضریب یادگیری متغیر کوچک شونده در طی پیشروی الگوریتم استفاده نمائیم تا اثرات داده‌های هم‌پوشان، در صورت وجود، کاهش یابد (نظیر الگوریتم شبیه‌سازی حرارت). چرا این عمل باعث کاهش اثرات داده‌های هم‌پوشان می‌شود؟ آیا $\eta_{(k)} = 1/k$ می‌تواند نرخ یادگیری مناسبی برای این منظور باشد (شروط همگرایی را داراست)؟

توجه کنید که الگوریتم کاهش گرادیان تضمینی نمی‌دهد که کمینه عمومی را پیدا کند. ولی به دلیل سادگی آن و اینکه برای هر تابع بهینه‌سازی‌ای قابل اعمال هست، استفاده زیادی پیدا کرده است.

به مساله تعیین سیستماتیک وزن‌های پرسپترون برگردیم. در این مساله می‌بایستی بر دو مشکل غلبه می‌کردیم. تا اینجا بر مشکل اول غلبه کرده و روشی برای بهبود بردار w در هر مرحله، تا رسیدن به یک جواب بهینه (محلی یا عمومی)، ارائه کردیم. حال باید به دنبال غلبه بر مشکل دوم، یعنی ارائه معیاری برای سنجش خوبی بردار w (یا همان ارائه تابع $J(w)$) باشیم.

باید $J(w)$ را به قسمی پیدا کنیم که به ازای w بهینه کمینه شود. بدیهی‌ترین تعریفی که برای $J(w)$ می‌توان در نظر گرفت به صورت زیر می‌باشد:

$$J(w) = |X_m(w)|; \quad X_m(w) = \{\text{samples } x_{(i)} \mid w^t x_{(i)} < 0\}$$

$J(w)$ فوق بیان‌گر تعدادی از بردارهای ویژگی ورودی است که $g(w)$ در دسته‌بندی آنها دچار اشتباه می‌شود. این تابع، یک تابع چند تکه‌ای، با مقادیر ثابت در هر تکه بوده (چرا؟) و برای بدست آوردن ∇J مناسب نمی‌باشد (زیرا گرادیان برای مقادیر ثابت برابر صفر می‌شود). بنابر این تابع دیگری به فرم زیر برای $J(w)$ پیشنهاد می‌کنیم:

$$J(w) = \sum_{x \in X_m(w)} -w^t x; \quad X_m(w) = \{\text{samples } x_{(i)} \mid w^t x_{(i)} < 0\}$$

می‌دانیم برای هر بردار w و هر بردار ویژگی x رابطه $w^t x / \|w\|$ نشان دهنده فاصله بردار ویژگی x از بردار w می‌باشد. در نتیجه تابع ارائه شده فوق، مجموع فواصل بردارهای ویژگی اشتباه دسته‌بندی شده، را از بردار w نشان می‌دهد. این تابع، یک تابع چند تکه‌ای، با ضابطه‌ای خطی در هر بازه بوده (چرا؟) و در هر بازه می‌توان مقدار گرادیان را برای آن بازه محاسبه نمود.

فاصله یک نقطه از یک ابر صفحه

اگر x یک نقطه در فضا و P ابر صفحه‌ای با معادله $g(x) = w^t x + w_0 = 0$ بوده و x_P تصویر نقطه بر روی ابر صفحه و r فاصله نقطه تا ابر صفحه باشد، می‌توان نوشت: $x = x_P + r \cdot w / \|w\|$ که در آن w بردار مشخصه صفحه و $w / \|w\|$ جهت آن (فارغ از اندازه) می‌باشد. در این صورت داریم:

$$g(x) = w^t (x_P + r \frac{w}{\|w\|}) + w_0 = w^t x_P + r \frac{w^t w}{\|w\|} + w_0 = w^t x_P + w_0 + r \frac{\|w\|^2}{\|w\|} = 0 + r \|w\|$$

در نتیجه $r = g(x) / \|w\|$.

می‌توان با استفاده از این تابع و الگوریتم کاهش گرادیان به حل مساله اقدام نمود. بخش بروز رسانی مقادیر w در الگوریتم نیز به صورت $w_{(k+1)} = w_{(k)} + \eta_{(k)} \sum_{x \in X_m} x$ در می‌آید. با توجه به این قانون در هر بار بروز رسانی، مقدار w با توجه به تمامی بردارهای ویژگی اشتباه دسته‌بندی شده (یا در واقع، متناسب با مجموع فواصل بردارهای ویژگی اشتباه دسته‌بندی شده، با مرز جداسازی)، بروز می‌شود. توجه داشته باشید که بدلیل وابستگی $J(w)$ به X_m ، پیدا کردن w بهینه برای $J(w)$ ، بصورت

تابع بهینه‌سازی مناسب برای پیدا کردن وزن‌های بهینه در پرسپترون

تحلیلی امکان پذیر نمی باشد.

مثال ۵: یک مساله دسته بندی دو کلاسه، در یک فضای چهار بعدی را با بردارهای ویژگی زیر در نظر بگیرید:

$$C_1 = \{(1, 1, -1, -1), (1, -1, -1, 1)\}, C_2 = \{(1, 1, 1, 1), (-1, -1, -1, 1)\}$$

می خواهیم دسته بندی این داده ها را با یک پرسپترون انجام دهیم. وزن های این پرسپترون $(w_0, w_1, w_2, w_3, w_4)$ را بدست آورید.

ابتدا بردارهای ویژگی داده شده را به بردارهای افزوده تبدیل کرده و سپس آنها را نرمال می کنیم:

$$C = \{(1, 1, 1, -1, -1), (1, 1, -1, -1, 1), (-1, -1, -1, -1, -1), (-1, 1, 1, 1, -1)\}$$

با در نظر گرفتن نرخ یادگیری به صورت $\eta_{(k)} = \eta = 1$ ، قانون بروز رسانی به فرم زیر در می آید:

$$w_{(k+1)} = w_{(k)} + \sum_{x \in X_m} x$$

با در نظر گرفتن بردار اولیه بصورت $w_{(0)} = \{0.25, 0.25, 0.25, 0.25, 0.25\}$ ، الگوریتم را اجرا می کنیم. داریم:

$$g(x_{(1)}) = 0.25 > 0, g(x_{(2)}) = 0.25 > 0, g(x_{(3)}) = -1.25 < 0, g(x_{(4)}) = 0.25 > 0$$

$x_{(3)}$ اشتباه دسته بندی شده است، پس $\{x_{(3)}\} = Z_m(w_{(0)})$ در نتیجه $w_{(0)}$ باید بروز شود:

$$w_{(1)} = w_{(0)} + \{x_{(3)}\} = (-0.75, -0.75, -0.75, -0.75, -0.75)$$

مجدداً به محاسبه X_m پرداخته و می بینیم که $\{x_{(1)}, x_{(2)}, x_{(4)}\} = X_m(w_{(1)})$ در نتیجه $w_{(1)}$ باید بروز شود:

$$w_{(2)} = w_{(1)} + (x_{(1)} + x_{(2)} + x_{(4)}) = (0.25, 2.25, 0.25, -1.75, -1.75)$$

مجدداً به محاسبه X_m پرداخته و می بینیم که تمامی بردارهای ویژگی درست دسته بندی شده اند. در نتیجه $w = w_{(2)}$.

□

دقت کنید که بی نهایت جواب دیگر نیز برای w در مثال قبل موجود می باشند. جواب نهایی بدست آمده کاملاً بستگی به مقدار اولیه انتخاب شده و نرخ یادگیری دارد.

به جای محاسبه X_m و بروز رسانی w با تمامی داده های اشتباه دسته بندی شده (بروز رسانی دسته ای)، می توان از بروز رسانی تک نمونه ای استفاده نمود. در بروز رسانی تک نمونه ای با مشاهده هر بردار ویژگی اشتباه دسته بندی شده، w بروز می شود. در پرسپترون، قانون بروز رسانی تک نمونه ای معادل قانون دسته ای فوق، به صورت $w_{(k+1)} = w_{(k)} + \eta_{(k)} x_i$ در می آید.

در بروز رسانی تک نمونه ای، بردارهای ویژگی را یکی یکی پشت سر هم با شبکه عصبی موجود دسته بندی کرده و با رسیدن به هر دسته بندی اشتباه، مقدار w را با توجه به همان بردار ویژگی بروز کنیم. پس از اینکه تمامی داده ها با شبکه، مورد آزمایش قرار گرفتند، اگر هنوز w بهینه به دست نیامده بود، برای بار دوم بردارهای ویژگی را با شبکه و با آخرین w بدست آمده، مجدداً مورد آزمون قرار می دهیم و همین روند را تا رسیدن به جواب بهینه ادامه می دهیم.

تمرین ۶: یک مساله کلاسه بندی دو کلاسه، در فضای دو بعدی، را با بردارهای ویژگی زیر در نظر بگیرید:

$$C_1 = \{(2, 1), (4, 3), (3, 5)\} \quad C_2 = \{(1, 3), (5, 6)\}$$

آیا الگوریتم کاهش گرادیان متوقف می شود؟ چرا؟ بردارهایی را که در مراحل مختلف بدست می آیند رسم کنید. آیا این بردارها دارای رفتار یا ویژگی خاص مشترکی هستند؟ با چه شرطی الگوریتم را متوقف کنیم که جواب مناسبی بدست آید؟

تمرین ۷: الگوریتم کاهش گرادیان با قانون بروز رسانی دسته ای سریع تر به جواب نزدیک می شود یا با قانون بروز رسانی تک نمونه ای؟ کدام یک در مواجهه با نویزهای تکین (تک داده های دور افتاده) موفق تر عمل می کند؟

تمرین ۸: نشان دهید که اگر $x_i \in X_m(w_{(k)})$ ؛ $|w_{(k)}^t x_i| / |x_i|^2$ ؛ $\eta_{(k)} > |w_{(k+1)}^t x_i|$ باشد قطعاً $w_{(k+1)}^t x_i$ منفی نمی شود.

در حل مسایل بهینه سازی با

الگوریتم کاهش گرادیان

بی نهایت جواب، بسته به مقدار

اولیه انتخاب شده و نرخ

یادگیری، می تواند بدست آید.

قانون بروز رسانی تک نمونه ای

در مقایسه با قانون بروز رسانی

دسته ای

۱-۱. معیار Relaxation برای آموزش نرون

معیار دیگری که برای آموزش یک نرون می‌تواند مورد استفاده قرار گیرد، معیار Relaxation به صورت زیر می‌باشد.

$$J(w) = \frac{1}{2} \sum_{x \in X_b(w)} \frac{(w^t x - b)^2}{|x|^t}; \quad X_b(w) = \{ \text{samples } x_{(i)} \mid w^t x_{(i)} \leq b \}$$

در این معیار، نه تنها بردارهای ویژگی اشتباه دسته‌بندی شده در تصحیح w شرکت خواهند داشت، بلکه بردارهای ویژگی درست دسته‌بندی شده ولی نزدیک به مرز (نزدیک‌تر از $b/\|w\|$) نیز در تصحیح شرکت خواهند داشت. در این روش تلاش بر این است که حاشیه‌ای با طول $b/\|w\|$ در اطراف مرز جداساز وجود داشته باشد که قدرت تعمیم جداساز را در مواجهه با بردارهای ویژگی جدید بالا ببرد. قانون بروز رسانی دسته‌ای برای این معیار به صورت زیر در خواهد آمد:

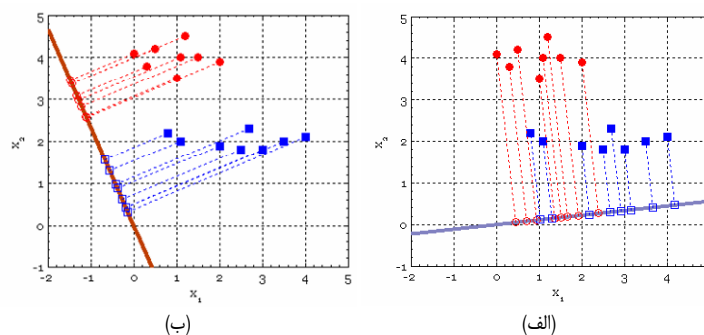
$$w_{(k+1)} = w_{(k)} + \eta_{(k)} \sum_{x \in X_b(w)} \frac{b - w_{(k)}^t x}{|x|^2} x$$

قانون بروز رسانی تک نمونه‌ای نیز برابر خواهد بود با:

$$w_{(k+1)} = w_{(k)} - \eta_{(k)} \frac{w_{(k)}^t x}{|x|^2} x$$

۲- جداساز خطی فشر

ایده اصلی در جداساز خطی فشر این است که داده‌ها در جهتی تابانده⁸ شوند که بیشترین جداسازی را داشته باشند و براحتی بتوان مرزی برای داده‌های تابیده شده مشخص نمود (کاهش ابعاد نیز انجام می‌شود). به عنوان مثال در شکل ۷ حاصل تابش در (ب) دارای جداسازی بسیار بیشتری از حاصل تابش در (الف) می‌باشد.



شکل ۷: تابش داده‌ها در جهت‌های مختلف. جداسازی در تابش (ب) بسیار بیشتر از جداسازی در تابش (الف) می‌باشد.

برای داده‌های معرفی شده در شکل ۷، جهت‌های مختلفی وجود دارد که با تابش در آن جهت‌ها، می‌توان به جداسازی داده‌ها، بدون خطا، اقدام نمود. حال سوال پیش می‌آید که کدام یک از این جهت‌ها بهترین جهت ممکن می‌باشد و چگونه می‌توان این جهت را پیدا نمود؟

در ادامه به بررسی دو سوال فوق با یک مثال دسته‌بندی دو کلاسه در فضای d بعدی می‌پردازیم. فرض کنید در این فضا دارای n بردار ویژگی آموزشی (x_1, x_2, \dots, x_n) باشیم: n_1 بردار ویژگی از کلاس اول و n_2 بردار ویژگی از کلاس دوم.

اگر بردار ویژگی x در جهت بردار w تابیده شود، از لحاظ هندسی حاصل تابش محل تقاطع خط عمود بر بردار w از نقطه x خواهد بود (به شکل ۷ توجه کنید). فاصله این نقطه از مرکز مختصات، $w^t x$ خواهد بود. با این کار بردار ویژگی از یک فضای دو بعدی به یک فضای یک بعدی نگاشت شده است.

قوانین بروز رسانی دسته‌ای و
تک نمونه‌ای برای معیار
Relaxation

در جداساز خطی فشر داده‌ها
در جهتی تابانده می‌شوند که
بیشترین جداسازی را داشته
باشند

مفهوم هندسی تابش

حال با این مقدمه به سوال اول بر می‌گردیم. چه جهتی بهترین جهت تابش بوده و داده‌ها بعد از تابش در راستای آن، به بهترین شکل قابل جداسازی خواهند بود؟ ممکن است جواب افراد مختلف به این سوال متفاوت و حتی در برخی موارد متناقض باشد. یکی از پاسخ‌ها می‌تواند این باشد که بهترین جهت، جهتی است که در آن، داده‌های دو کلاس متمایز بیشترین فاصله را از هم داشته باشند، در حالی که داده‌های هر کلاس بیشترین نزدیکی و مقاربت را به یکدیگر داشته باشند.

پاسخ ارائه شده که توسط فیشر ارائه شده، منطقی و معقول می‌باشد، ولی این جهت به چه صورت قابل شناسایی می‌باشد؟ (سوال دوم) برای پاسخ به این سوال ناچاریم مفاهیم کیفی ارائه شده در پاسخ سوال اول را کمی کنیم. فرض کنید μ_1 و μ_2 میانگین کلاس‌های C_1 و C_2 بوده و $\hat{\mu}_1$ میانگین داده‌های تابیده شده کلاس C_1 و $\hat{\mu}_2$ میانگین داده‌های تابیده شده کلاس C_2 باشند. در این صورت می‌توان فاصله داده‌های دو کلاس را با $d = |\hat{\mu}_1 - \hat{\mu}_2|$ کمی کرد. داریم:

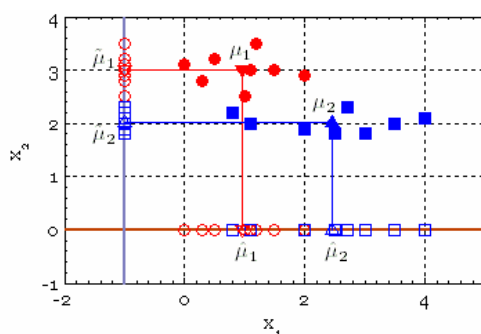
$$\mu_1 = \frac{1}{n_1} \sum_{x_i \in C_1} x_i, \quad \mu_2 = \frac{1}{n_2} \sum_{x_i \in C_2} x_i$$

و اگر راستای تابش، w باشد، داریم:

$$\hat{\mu}_1 = \frac{1}{n_1} \sum_{x_i \in C_1} w^t x_i = w^t \left(\frac{1}{n_1} \sum_{x_i \in C_1} x_i \right) = w^t \mu_1$$

و به همین ترتیب $\hat{\mu}_2 = w^t \mu_2$. در نتیجه $d = |\hat{\mu}_1 - \hat{\mu}_2| = |w^t \mu_1 - w^t \mu_2| = w^t |\mu_1 - \mu_2|$.

قبلا توضیح ندادیم که چرا این معیار به تنهایی معیار مناسبی نیست. ولی اکنون با کمی‌سازی برخی از مفاهیم می‌توانیم این سوال را راحت‌تر جواب دهیم. برای این منظور به داده‌های دو کلاس ارائه شده در شکل ۸ توجه نمایم. داده‌های نشان داده شده، دو بار در دو جهت افقی و عمودی تابانده شده‌اند. $\hat{\mu}_1, \hat{\mu}_2$ میانگین داده‌های دو کلاس تابانده شده در جهت افقی و $\tilde{\mu}_1, \tilde{\mu}_2$ میانگین داده‌های دو کلاس تابانده شده در جهت عمودی می‌باشند.



شکل ۸: تابش داده‌ها در دو جهت متفاوت برای سنجش معیار «فاصله بین کلاس‌ها»

همانطور که مشاهده می‌کنید، $|\hat{\mu}_1 - \hat{\mu}_2|$ بزرگتر از $|\tilde{\mu}_1 - \tilde{\mu}_2|$ است ولی جداسازی داده‌های تابیده شده در جهت افقی راحت‌تر از جداسازی داده‌های تابیده شده در جهت عمودی نیست (می‌توان بین داده‌های دو کلاس تابیده شده در جهت عمودی مرزی بدون خطا رسم کرد ولی برای داده‌های تابیده شده در جهت افقی این کار امکان‌پذیر نیست). دلیل این امر، همانطور که فیشر نیز بیان می‌کند، در نظر نگرفتن فواصل بین کلاسی داده‌های تابیده شده در دو جهت می‌باشد. فواصل بین کلاسی داده‌ها، یا پراکندگی آنها، را می‌توان با مجموع فواصل داده‌های هر کلاس از میانگین آن کلاس، کمی نمود:

$$S_1 = \sum_{x_i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T, \quad S_2 = \sum_{x_i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T$$

برای محاسبه پراکندگی داده‌های تابیده شده نیز داریم:

$$\begin{aligned} \hat{S}_1 &= \sum_{\hat{x}_i \in \hat{C}_1} (\hat{x}_i - \hat{\mu}_1)^2 = \sum_{x_i \in C_1} (w^t x_i - w^t \mu_1)^2 = \sum_{x_i \in C_1} (w^t (x_i - \mu_1))^2 = \sum_{x_i \in C_1} (w^t (x_i - \mu_1))^t (w^t (x_i - \mu_1)) \\ &= \sum_{x_i \in C_1} ((x_i - \mu_1)^t w) ((x_i - \mu_1)^t w) = \sum_{x_i \in C_1} w^t (x_i - \mu_1) (x_i - \mu_1)^t w = w^t S_1 w \end{aligned}$$

و به همین ترتیب $\hat{S}_2 = w^t S_2 w$. هر چه داده‌های یک کلاس متمرکزتر باشند، پراکندگی آن کلاس کمتر است و بالعکس.

تابع هدف فیشر

مطابق معیار فیشر و با توجه به تعاریف میانگین و پراکندگی، بهترین جهت از دید فیشر، جهتی است که میانگین داده‌های تائیده شده کلاس‌های مختلف، بیشترین فاصله را از هم داشته باشند و در حین حال پراکندگی داده‌های تائیده شده مربوط به هر کلاس کمترین مقدار ممکن باشد. و این یعنی اینکه $|\hat{\mu}_1 - \hat{\mu}_2|$ (یا معادل آن $(\hat{\mu}_1 - \hat{\mu}_2)^2$) بیشینه و $\hat{S}_1 + \hat{S}_2$ کمینه شود. پس باید دنبال جهتی باشیم که تابع هدف زیر را بیشینه کند:

$$J(w) = \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\hat{S}_1 + \hat{S}_2} \quad (9)$$

با بسط این رابطه و جایگزینی روابط مربوط به $\hat{\mu}_1, \hat{\mu}_2, \hat{S}_1$ و \hat{S}_2 خواهیم داشت:

$$J(w) = \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\hat{S}_1 + \hat{S}_2} = \frac{w^t (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t w}{w^t S_1 w + w^t S_2 w} = \frac{w^t (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t w}{w^t (S_1 + S_2) w}$$

که اگر $S_w = S_1 + S_2$ و $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t$ در نظر گرفته شوند، خواهیم داشت:

$$J(w) = \frac{w^t S_B w}{w^t S_w w}$$

برای اینکه $J(w)$ بیشینه شود باید داشته باشیم:

$$\frac{d}{dw} J(w) = \frac{(2S_B w)w^t S_w w - (2S_w w)w^t S_B w}{(w^t S_w w)^2} = 0$$

برای این منظور کافی است صورت کسر فوق مساوی صفر قرار داده شود، که با تقسیم صورت بر $w^t S_w w$ خواهیم داشت:

$$S_B w - \frac{w^t S_B w}{w^t S_w w} S_w w = 0 \Rightarrow S_B w = \frac{w^t S_B w}{w^t S_w w} S_w w$$

و این یعنی،

$$S_B w = J(w).S_w w \quad (10)$$

$J(w)$ یک مقدار عددی بوده و می‌تواند با λ جایگزین شود. در نتیجه $S_w^{-1} S_B w = \lambda w$. از طرفی می‌دانیم که حاصلضرب S_B در هر برداری، راستایی در جهت $(\mu_1 - \mu_2)$ خواهد داشت. زیرا:

$$S_B y = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t y = (\mu_1 - \mu_2)((\mu_1 - \mu_2)^t y) = \alpha.(\mu_1 - \mu_2)$$

پس راستای w برابر خواهد بود با:

$$\bar{w} = S_w^{-1}(\mu_1 - \mu_2) \quad (11)$$

ما در اینجا به جای حل دقیق رابطه (10) از حل شهودی مساله برای یافتن جهت بردار w استفاده کردیم. برای حل دقیق این مساله می‌توان از مساله مقادیر ویژه تعمیم یافته استفاده نمود. زیرا این رابطه، دقیقاً دارای صورت یک مساله مقادیر ویژه تعمیم یافته می‌باشد.

مساله مقادیر ویژه تعمیم یافته

برای دو ماتریس داده شده A و B مساله پیدا کردن جفت مقدار (λ, x) بطوری که $Ax = \lambda Bx$ ، را مساله مقادیر ویژه تعمیم یافته می‌گویند.

تمرین ۹: با استفاده از نحوه حل مقادیر ویژه تعمیم یافته، رابطه (10) را بطور دقیق حل نمایید.

مثال ۶: یک مساله دسته‌بندی در فضای دو بعدی با دو کلاس و بردارهای ویژگی زیر در نظر بگیرید:

$$c_1 = \{(1,2), (2,3), (3,3), (4,5), (5,5)\}, \quad c_2 = \{(1,0), (2,1), (3,1), (3,2), (5,3), (6,5)\}$$

جهت فیشر را برای این داده‌ها پیدا کنید.

ابتدا میانگین داده‌ها را بدست می‌آوریم:

$$\mu_1 = \sum_{x_i \in C_1} x_i = (3, 3.6), \quad \mu_2 = \sum_{x_i \in C_2} x_i = (3.3, 2)$$

سپس پراکندگی هر کدام از کلاس‌ها را محاسبه می‌کنیم:

$$S_1 = \sum_{x_i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T = \begin{bmatrix} 10 & 8.0 \\ 8.0 & 7.2 \end{bmatrix}, \quad S_2 = \sum_{x_i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T = \begin{bmatrix} 17.3 & 16 \\ 16 & 16 \end{bmatrix}$$

در نتیجه:

$$S_w = S_1 + S_2 = \begin{bmatrix} 27.3 & 24 \\ 24 & 23.2 \end{bmatrix}, \quad S_w^{-1} = \begin{bmatrix} 0.39 & -0.41 \\ -0.41 & 0.47 \end{bmatrix}$$

و نهایتاً اینکه:

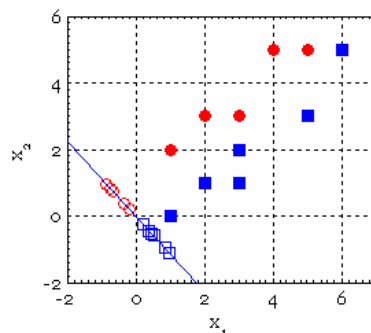
$$w = S_w^{-1}(\mu_1 - \mu_2) = \begin{bmatrix} -0.79 \\ 0.89 \end{bmatrix}$$

مختصات داده‌های جدید (تاییده شده) بر روی محور تابش عبارتند از:

$$y_{1i} = w^T x_{1i} \Rightarrow Y = \{y_{11}, y_{12}, \dots, y_{15}\} = \{0.99, 1.08, 0.29, 1.28, 0.48\}$$

$$y_{2i} = w^T x_{2i} \Rightarrow Y = \{y_{21}, y_{22}, \dots, y_{26}\} = \{-0.79, -0.7, -1.49, -0.6, -1.3, -0.31\}$$

داده‌های دو کلاس، خط تابش و داده‌های تاییده شده، در شکل زیر نشان داده شده‌اند.



□

جداساز خطی فیشر (FLD⁹) را برای حالت دو بعدی بررسی کردیم. حال می‌خواهیم ببینیم که جداساز فیشر برای حالت چند بعدی به چه صورت خواهد بود؟ اگر به مرور محاسبات انجام شده تا بدست آوردن رابطه (۱۱) پردازید، می‌بینید که تمامی معادلات و روابط نوشته شده، مستقل از تعداد ابعاد داده‌ها، برقرار بوده و صحیح می‌باشند. در نتیجه رابطه (۱۱) برای ابعاد بالاتر نیز برقرار بوده و برای هر فضای d بعدی، برداری d بعدی بدست می‌آید.

جداساز خطی فیشر را برای حالت دو کلاسه بررسی کردیم. حال می‌خواهیم جداساز فیشر را برای حالت چند کلاسه بررسی کنیم. برای این منظور در ساده‌ترین حالت می‌توان از ماشین خطی کامل (با $k-1$ جهت فیشر مستقل از هم) و یا ماشین خطی دوبعدی جدپذیر (با $k(k-1)/2$ جهت فیشر مستقل از هم) استفاده نمود.

فیشر و فضاهای چند بعدی

فیشر و مسایل چند کلاسه

⁹ Fisher Linear Discriminant

می‌توان به جای تبدیل مساله چند کلاسه به مسایل دو کلاسه کوچکتر مستقل از هم و حل این مسایل کوچکتر، روشی را برای حل یک جای مساله نیز ارائه نمود. این روش‌ها با عنوان تحلیل جداسازهای چندتایی (MDA¹⁰) مشهورند.

برای حالت تعمیم یافته فیشر، یک مساله دسته‌بندی را با k کلاس، در یک فضای d بعدی در نظر بگیرید. می‌خواهیم ماتریس فیشر (w) را پیدا کنیم به طوری که هر بردار ویژگی x_i را تحت رابطه $y_i = w^t x_i$ ، به بردار ویژگی y_i تبدیل کرده و این تبدیل به صورتی باشد که در فضای تابش بدست آمده، پراکندگی بین کلاسی داده‌ها، بیشینه و مجموع پراکندگی داده‌های داخل کلاسی، کمینه باشد.

اگر تعداد بردارهای ویژگی کلاس i ام را با n_i نشان دهیم و مجموع تمامی بردارهای ویژگی برابر n باشد، پراکندگی کلی داده‌ها (S_T)، مجموع پراکندگی داده‌های داخل کلاسی (S_W) و پراکندگی بین کلاسی داده‌ها (S_B) در فضای تابش برابر خواهد بود با:

$$S_T = \sum_{i=1}^n (y_i - m)(y_i - m)^t, \quad m = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^k n_i m_i, \quad m_i = \frac{1}{n_i} \sum_{y_j \in C_i} y_j$$

$$S_W = \sum_{i=1}^k S_i, \quad S_i = \sum_{y_j \in C_i} (y_j - m_i)(y_j - m_i)^t$$

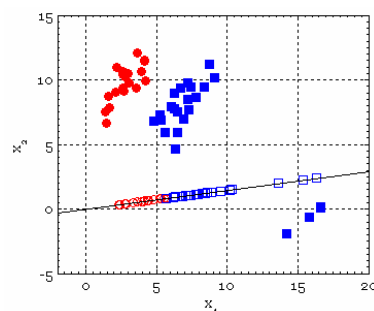
$$S_T = S_W + S_B \Rightarrow S_B = S_T - S_W = \sum_{i=1}^k n_i (m_i - m)(m_i - m)^t$$

معیار بهینه‌سازی نیز به صورت برابر خواهد بود با:

$$J(w) = \frac{\det(w S_B w^t)}{\det(w S_W w^t)} \quad (۱۲)$$

رابطه ۰ نشان می‌دهد که S_B مجموع k ماتریس از درجه یک می‌باشد و فقط $k-1$ از آنها مستقل می‌باشند (k مین ماتریس با توجه به این $k-1$ ماتریس قابل بدست آوردن است، چگونه؟). در نتیجه S_B از درجه $k-1$ بوده و دارای حداکثر $k-1$ مقدار ویژه غیر صفر می‌باشد. این نشان می‌دهد که تابش داده‌ها در فضای $k-1$ بعدی حاصل از بردارهای ویژه S_B تغییری در مقدار $J(w)$ نخواهد داد. و این یعنی اینکه با این روش فضایی با حداکثر $k-1$ بعد به عنوان فضای نگاشت شده خواهیم داشت. پس در حالت کلی جداساز خطی فیشر در یک مساله دسته‌بندی k کلاسه، فضا را به $k-1$ بعد کاهش می‌دهد.

استفاده از جداساز خطی فیشر همیشه مطلوب نیست. در جداسازی فیشر فرض می‌شود که بردارهای ویژگی تاییده شده دارای توزیع نرمال می‌باشند (چرا؟). اگر این فرض برقرار نباشد، جداساز فیشر، جهت مناسبی را برای تابش بر نمی‌گرداند (شکل ۹).

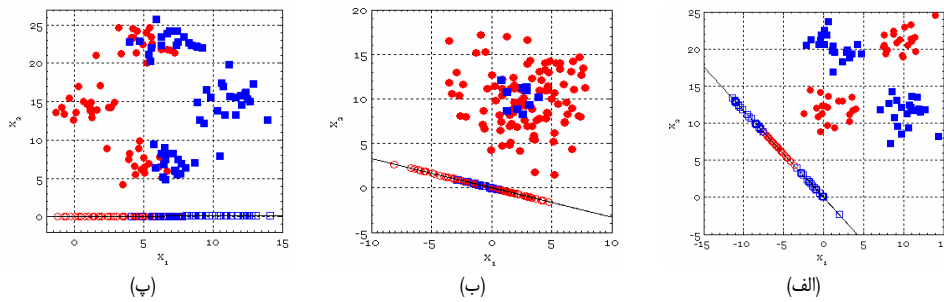


شکل ۹: جداساز فیشر وقتی که بردارهای ویژگی تاییده شده دارای توزیع نرمال نباشند، جهت مناسبی را بر نمی‌گرداند.

جداساز خطی فیشر، در حالتی که $J(w)$ صفر بوده یا همیشه دارای مقدار بزرگی باشد، با شکست مواجه شده و قادر نیست جواب مسایل را پیدا کند. $J(w) = 0$ در حالتی رخ می‌دهد که میانگین دو کلاس مساوی باشند (حالت‌های الف و ب از شکل ۱۰). اگر داده‌ها در هر جهتی که تابانده شوند، دارای تداخل زیادی باشند، نیز حالتی است که $J(w)$ همیشه مقادیر کوچکی را بر می‌گرداند (حالت پ از شکل ۱۰). در شکل‌ها توجه کنید که جهت‌های فیشر بدست آمده نامناسب بوده و خطای بالایی را در

¹⁰ Multiple Discriminant Analysis

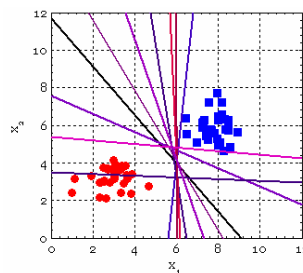
جداسازی بوجود می‌آورند.



شکل ۱۰: حالات نامناسب برای استفاده از جداساز خطی فشر، با جهت‌های نامناسب بدست آمده.

۲- جداسازهای با بیشترین حاشیه

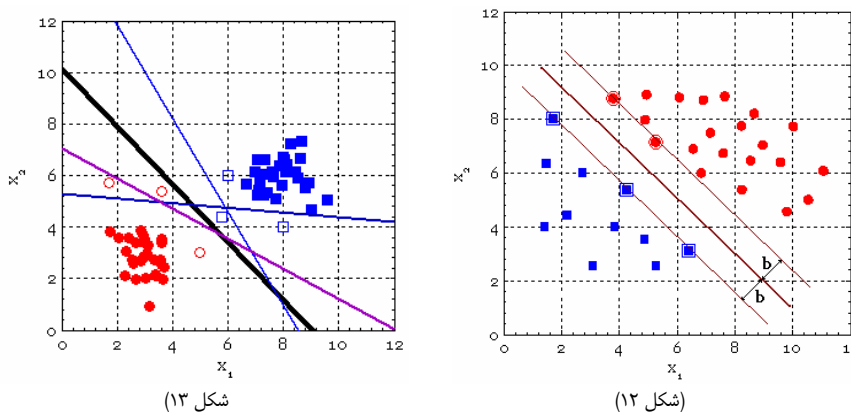
برای هر مساله دسته‌بندی بصورت خطی جداپذیر می‌توان بی‌نهایت ابر صفحه ارائه نمود که به عنوان جداساز عمل نمایند (شکل ۱۱). سوالی که پیش می‌آید این است که کدام یک از این جداسازها بهترین می‌باشد؟



شکل ۱۱: بهترین جداساز کدام است؟

ایده اصلی در جداسازهای با بیشترین حاشیه^{۱۱} همانطور که از نامشان پیداست این است که ابرصفحه‌هایی را با بیشترین حاشیه بدست آورند که کلاس‌ها را از هم جدا کند. حاشیه به فاصله‌ای گفته می‌شود که دو خط موازی جداساز بطور مساوی از دو طرف طی می‌کنند، تا یکی از آن دو به یکی از داده‌ها برخورد نماید (شکل ۱۲). از بین جداسازهای خطی، جداسازی که بیشترین حاشیه را داشته باشد (بیشترین فاصله را از داده‌ها داشته باشد)، خطای تعمیم را حداقل خواهد کرد (شکل ۱۳).

جداساز با بیشترین حاشیه، دارای کمترین خطا در تعمیم خواهد بود.



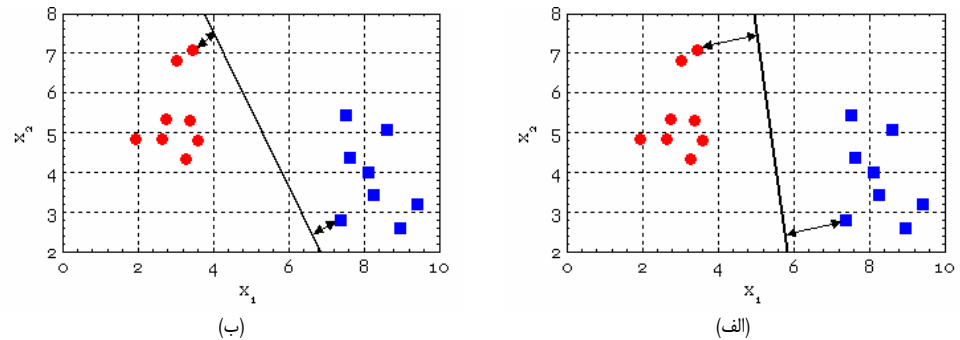
شکل ۱۲: یک مجموعه داده با مرز جداسازی. حاشیه جداساز با ضخامت b نیز نشان داده شده است..

شکل ۱۳: جداسازهای با حاشیه بیشتر، در هنگام تعمیم خطای کمتری خواهند داشت. اشکال توپر داده‌های آموزشی و اشکال توخالی داده‌های جدید (آزمایشی) هستند. به رابطه عملکرد صحیح در مواجهه با داده‌های جدید و اندازه حاشیه توجه کنید.

^{۱۱} Maximum Margin

پیدا کردن جداساز با بیشترین
حاشیه، چگونه؟

حال که مفهوم حاشیه و جداساز با بیشترین حاشیه مشخص شد باید به دنبال راهی برای پیدا کردن جداساز با بیشترین حاشیه باشیم. ایده اصلی برای انجام این کار این است که سعی کنیم جداسازی را بیابیم که فاصله آن از نزدیک‌ترین داده بیشینه باشد. از آنجایی که در دو طرف یک خط (یا در حالت کلی یک ابر صفحه) داده‌های دو کلاس قرار دارند، فاصله گرفتن خط از داده‌های یکی از کلاس‌ها مترادف خواهد بود با نزدیک شدن خط به داده‌های کلاس دیگر. در نتیجه، خط (ابر صفحه) بهینه در حالتی بدست خواهد آمد که فاصله خط جداساز از نزدیک‌ترین داده‌های دو طرف مساوی و بیشینه باشد (شکل ۱۴).



شکل ۱۴: جداسازی بیشترین حاشیه را خواهد داشت که فاصله آن از نزدیک‌ترین داده‌های دو طرف مساوی و بیشینه باشد (الف).

پیدا کردن بردارهای پشتیبان
مستلزم شناخت جداساز با
بیشترین حاشیه می‌باشد و پیدا
کردن جداساز با بیشترین
حاشیه نیز مستلزم شناخت
بردارهای پشتیبان است!

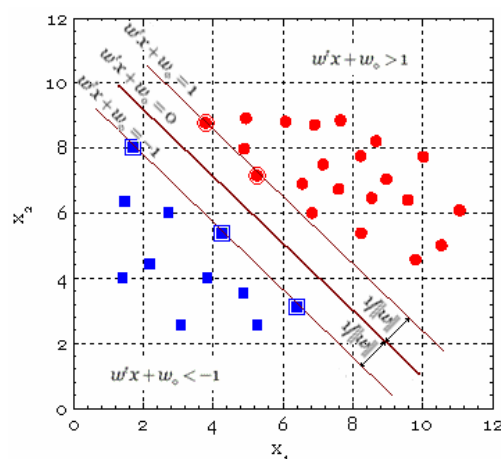
نزدیک داده‌های دو طرف، به جداساز با بیشترین حاشیه را بردارهای پشتیبان می‌نامند. پیدا کردن بردارهای پشتیبان کار مشکلی است، چون پیدا کردن آنها مستلزم شناخت جداساز با بیشترین حاشیه می‌باشد و پیدا کردن جداساز با بیشترین حاشیه نیز مستلزم شناخت بردارهای پشتیبان است! در ادامه سعی می‌کنیم با شناخت بیشتر از مساله و فرموله کردن و ساده‌سازی آن بر این حلقه غلبه کنیم.

جداساز یا بیشترین حاشیه یک جداساز خطی بوده و همانطور که قبلاً مشاهده کردیم دارای رابطه $g(x) = w^T x + w_0 = 0$ می‌باشد. فاصله هر داده x از این ابر صفحه برابر است با:

$$d = \frac{|g(x)|}{\|w\|} = \frac{|w^T x + w_0|}{\|w\|}$$

فرض کنید x_i بردار پشتیبان بوده و تابع جداساز $g(x)$ نیز، با تغییر اندازه w و w_0 ، طوری تنظیم شده باشد که:

$$|g(x_i)| = |w^T x_i + w_0| = 1$$



شکل ۱۵: رابطه جداساز و حاشیه‌ها

در این صورت فاصله x_i تا جداساز، $1/\|w\|$ ، اندازه حاشیه، $2/\|w\|$ و معادلات خطوط حاشیه، $w^T x + w_0 = \pm 1$ خواهند بود (شکل ۱۵). صورت مساله نیز به این صورت درمی‌آید: برای $g(x) = w^T x + w_0$ جفت مقدار (w, w_0) را طوری بیابید که

$\|w\|/2$ بیشینه بوده و برای داده‌های آموزشی داده شده، داشته باشیم:

$$\begin{cases} g(x) \geq 1 & \text{if } x \in C_1 \\ g(x) \leq -1 & \text{if } x \in C_2 \end{cases}$$

با تعریف متغیر Z بصورت

$$\begin{cases} z = 1 & \text{if } x \in C_1 \\ z = -1 & \text{if } x \in C_2 \end{cases}$$

رابطه فوق به صورت $z.g(x) \geq 1$ در می‌آید. همچنین بیشینه کردن $\|w\|/2$ معادل کمینه‌سازی $\frac{1}{2}\|w^2\|$ می‌باشد. با این تغییرات مساله به یک مساله استاندارد بهینه‌سازی با قیدهای نامساوی، تبدیل خواهد شد:

$$\begin{aligned} & \text{minimize } J(w) = \frac{1}{2}\|w^2\| \\ & \text{subject to } z_i.(w^t x_i + w_o) - 1 \geq 0, \quad \forall i \end{aligned} \quad (13)$$

مساله بهینه‌سازی معادل پیدا

کردن جدا ساز با بیشترین

حاشیه

مسایل بهینه‌سازی، ضرایب لاگرانژ و شرایط KKT

یک مساله بهینه‌سازی با قیدهای نامساوی دارای حالت کلی زیر می‌باشد:

$$\begin{aligned} & \text{minimize } J(\theta) \\ & \text{subject to } f_i(\theta) \geq 0, \quad i = 1, 2, \dots, m \end{aligned}$$

لاگرانژ نشان داد که می‌توان قیدها را در تابع اصلی بهینه شونده ادغام نموده و مساله را به صورت زیر در آورد:

$$\text{minimize } L(\theta, \lambda), \quad L(\theta, \lambda) = J(\theta) - \sum_{i=1}^m \lambda_i f_i(\theta)$$

که در آن λ را ضرایب لاگرانژ گویند که مجهول بوده و طوری باید بدست آیند که تابع $L(\theta, \lambda)$ را بیشینه نمایند. پس خواهیم داشت:

$$\min_{\theta} J(\theta) = \min_{\theta} \max_{\lambda} L(\theta, \lambda) = \max_{\lambda} \min_{\theta} L(\theta, \lambda)$$

همچنین نشان داده شده است برای اینکه θ^* یک کمینه‌ساز محلی برای مساله فوق باشد، باید شرایط زیر را دارا باشد (شرایط Karush-kuhn-Tucker):

$$\begin{aligned} (1) \quad & \frac{\partial}{\partial \theta} L(\theta^*, \lambda) = 0 \\ (2) \quad & \lambda_i \geq 0, \quad i = 1, 2, \dots, m \\ (3) \quad & \lambda_i f_i(\theta^*) = 0, \quad i = 1, 2, \dots, m \end{aligned}$$

صورت لاگرانژی این مساله به فرم زیر می‌باشد:

$$L(w, w_o, \lambda) = \frac{1}{2} w^t w - \sum_{\forall i} \lambda_i [z_i.(w^t x_i + w_o) - 1] \quad (14)$$

که برای (w, w_o) بهینه، باید شرایط زیر (شرایط KKT) برقرار باشد:

$$\begin{aligned} (1) \quad & \frac{\partial}{\partial w} L(w, w_o, \lambda) = 0, \quad \frac{\partial}{\partial w_o} L(w, w_o, \lambda) = 0 \\ (2) \quad & \lambda_i \geq 0, \quad \forall_i \\ (3) \quad & \lambda_i [z_i(w^t x_i + w_o) - 1] = 0, \quad \forall_i \end{aligned} \quad (15)$$

با ترکیب روابط (14) و (15) خواهیم داشت:

صورت لاگرانژی مساله و بیان

شرایط KKT برای آن

$$\begin{aligned}\frac{\partial}{\partial w} L(w, w_0, \lambda) &= w - \sum_{\forall i} \lambda_i z_i x_i = \mathbf{0} \Rightarrow w = \sum_{\forall i} \lambda_i z_i x_i \\ \frac{\partial}{\partial w_0} L(w, w_0, \lambda) &= -\sum_{\forall i} \lambda_i z_i = \mathbf{0} \Rightarrow \sum_{\forall i} \lambda_i z_i = \mathbf{0}\end{aligned}$$

پس صورت نهایی مساله با استفاده از ضرایب لاگرانژ بصورت زیر در می آید:

$$\max_{\lambda} L(w, w_0, \lambda)$$

subject to:

$$(1) w = \sum_{\forall i} \lambda_i z_i x_i \quad (16)$$

$$(2) \sum_{\forall i} \lambda_i z_i = \mathbf{0}$$

$$(3) \lambda \geq \mathbf{0}$$

با ترکیب قید (۱) از رابطه (۱۶) در تابع $L(w, w_0, \lambda)$ خواهیم داشت:

$$\begin{aligned}L(w, w_0, \lambda) &= \frac{1}{2} w^t w - \sum_{\forall i} \lambda_i [z_i \cdot (w^t x_i + w_0) - 1] \\ &= \frac{1}{2} \sum_{\forall i} \lambda_i z_i x_i^t \sum_{\forall j} \lambda_j z_j x_j - \sum_{\forall i} \lambda_i z_i \sum_{\forall j} \lambda_j z_j x_j^t x_i - w_0 \sum_{\forall i} \lambda_i z_i + \sum_{\forall i} \lambda_i \quad (17) \\ &= -\frac{1}{2} \sum_{\forall i} \lambda_i \lambda_j z_i z_j x_i^t x_j - w_0 \sum_{\forall i} \lambda_i z_i + \sum_{\forall i} \lambda_i\end{aligned}$$

و با توجه به قید (۲) از رابطه (۱۶) و همچنین توجه به این نکته که در رابطه (۱۷) (w, w_0) ظاهر نشده اند، رابطه (۱۷) به صورت زیر در می آید:

$$L(\lambda) = -\frac{1}{2} \sum_{\forall i} \sum_{\forall j} \lambda_i \lambda_j z_i z_j x_i^t x_j + \sum_{\forall i} \lambda_i \quad (18)$$

در نتیجه صورت نهایی مساله با استفاده از ضرایب لاگرانژ بصورت زیر ساده می شود:

$$\max_{\lambda} L(\lambda) = \sum_{\forall i} \lambda_i - \frac{1}{2} \sum_{\forall i} \sum_{\forall j} \lambda_i \lambda_j z_i z_j x_i^t x_j$$

subject to:

$$(1) \sum_{\forall i} \lambda_i z_i = \mathbf{0}$$

$$(2) \lambda \geq \mathbf{0}$$

$L(\lambda)$ را می توان بصورت ماتریسی نیز بیان نمود:

$$L(\lambda) = \sum_{\forall i} \lambda_i - \frac{1}{2} \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{bmatrix}^t H \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{bmatrix} = \sum_{\forall i} \lambda_i - \frac{1}{2} \lambda^t H \lambda \quad (20)$$

که در آن $\lambda = [\lambda_1 \cdots \lambda_n]^t$ و $H_{ij} = z_i z_j x_i^t x_j$.

بیشینه سازی $L(\lambda)$ می تواند با روش های تکراری محاسبات عددی (نظیر برنامه ریزی درجه دوم) حل شود.

با حل این مساله، مقادیر λ_i ها بدست می آیند. تمامی بردارهای ویژگی ای که به ازای آنها λ_i غیر صفر باشد، بردارهای پشتیبان می باشند (چرا؟). از رابطه $w = \sum_{\forall i} \lambda_i z_i x_i$ بدست می آید. w_0 نیز از هر کدام از (x_i, z_i) هایی که λ_i معادل آنها صفر نباشد، از رابطه $w_0 = 1/z_i - w^t x_i$ قابل حصول است.

تابع جداساز نیز بصورت

$$g(x) = \left(\sum_{x_i \in S} \lambda_i z_i x_i \right)^t x + w_0 \quad (21)$$

صورت نهایی مساله

بهینه سازی معادل با پیدا کردن
جداساز با بیشترین حاشیه

مشخصات بردارهای پشتیبان،
مرز جداسازی و تابع جداساز با
بیشترین حاشیه

می‌باشد که در آن $S = \{x_i | \lambda_i \neq 0\}$ یعنی پشتیبان می‌باشد،

مثال ۷: برای یک مساله دسته‌بندی دو کلاسه در فضای دو بعدی با بردارهای ویژگی زیر، جداساز با بیشترین حاشیه را بیابید.

$$C_1 = \{(1,6), (1,10), (4,11)\}, C_2 = \{(5,2), (7,6), (10,4)\}$$

ماتریس‌های X و Z برای این مثال بصورت زیر خواهند بود (۱ برای کلاس اول و -۱ برای کلاس دوم):

$$X^t = \begin{bmatrix} 1 & 1 & 4 & 5 & 7 & 10 \\ 6 & 10 & 11 & 2 & 6 & 4 \end{bmatrix}, Z^t = [1 \quad 1 \quad 1 \quad -1 \quad -1 \quad -1]$$

برای محاسبه ماتریس H از رابطه $H = XX^t.ZZ^t$ استفاده می‌کنیم (زیرا دارایه این ماتریس $H_{ij} = z_i z_j x_i^t x_j$ می‌باشند):

$$H = \begin{bmatrix} 37 & 61 & 70 & -17 & -43 & -34 \\ 61 & 101 & 114 & -25 & -67 & -50 \\ 70 & 114 & 137 & -42 & -94 & -84 \\ -17 & -25 & -42 & 29 & 47 & 58 \\ -43 & -67 & -94 & 37 & 85 & 94 \\ -34 & -50 & -84 & 58 & 84 & 116 \end{bmatrix}$$

با داشتن ماتریس H باید بردار λ با اعداد مثبت را طوری بیابیم که $L(\lambda) = \sum_{i=1}^6 \lambda_i - \frac{1}{2} \lambda^t H \lambda$ بیشینه شده و $Z^t \lambda$ برابر صفر شود. برای حل این مساله بهینه‌سازی از ابزارهای مربوطه می‌توان بهره برد. به عنوان مثال مطلب، قادر به حل مسائل بهینه‌سازی استاندارد (به فرم زیر) می‌باشد:

$$\text{minimize } L_D(\lambda) = \frac{1}{2} \lambda^t H \lambda + f^t \lambda \quad \text{constraint to } A\lambda \leq a \text{ and } B\lambda = b$$

که در آن A ، B و H ماتریس‌های $n \times n$ بوده و f ، a و b بردارهای n تایی هستند.

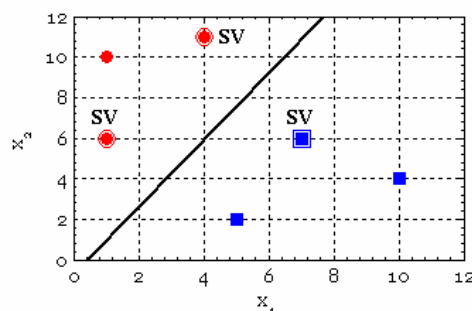
در مثال ما با در نظر گرفتن بردار f بصورت برداری از -۱‌ها، بردارهای a و b بردارهایی از صفرها، ماتریس A بصورت ماتریس قطری با درایه‌های قطری -۱ و ماتریس B با سطر اول برابر Z و صفر برای سایر سطرها، مساله به فرم استاندارد تبدیل شده و با دستور زیر در مطلب قابل حل می‌باشد:

$$\lambda = \text{quadprog}(H, f, A, a, B, b)$$

با حل مساله بهینه‌سازی فوق، λ به صورت $\lambda = [0.036 \quad 0 \quad 0.039 \quad 0 \quad 0.076 \quad 0]^t$ بدست می‌آید. در نتیجه (w, w_0) نیز برابر خواهند بود با:

$$w = \sum_{i=1}^n \lambda_i z_i x_i = (\lambda.z) x = [-0.33 \quad 0.20]^t$$

$$w_0 = \frac{1}{z_1} - w^t x_1 = 0.13$$



□

۴- روش‌های بر پایه کمینه‌سازی خطا

روش‌های بر پایه کمینه‌سازی خطا، در حالتی که داده‌ها به صورت خطی جدایی‌پذیر نباشند، خطای موجود را کمینه می‌کنند.

یکی از بزرگترین مزایای استفاده از جداسازهای خطی ساده استفاده از آنهاست. تا جایی که در برخی از موارد با وجود اینکه می‌دانیم که بردارهای ویژگی داده شده بصورت خطی از هم جدایی‌پذیر نیستند ولی باز هم علاقمند به استفاده از توابع جداساز خطی هستیم. در این موارد تلاش بر این است که جداسازی یافت شود که خطای موجود را کمینه نماید.

در این بخش به بررسی روش‌هایی می‌پردازیم که بر اساس کمینه‌سازی خطا اقدام به یافتن بهترین جداساز می‌نمایند. در این قسمت ابتدا به عنوان یک مثال روش‌های بر پایه کمینه‌سازی خطا را برای حالتی بررسی می‌کنیم که داده‌های خطی شده دارای توزیع نرمال باشند. سپس در بخش‌های بعدی به بررسی روش‌های کلی‌تر بر پایه مربعات خطا خواهیم پرداخت.

۴.۱- روش‌های بر پایه کمینه‌سازی خطا و توزیع‌های نرمال

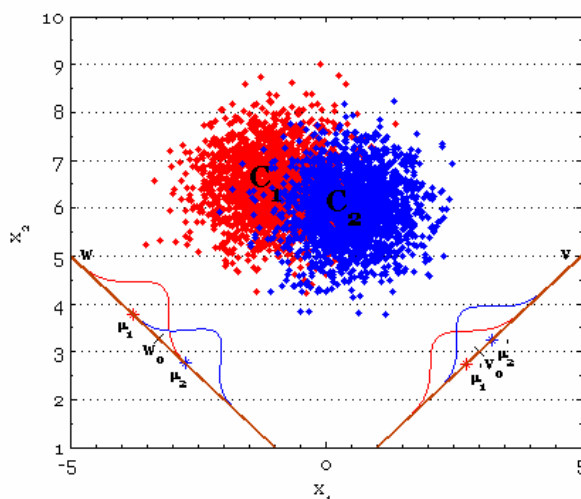
[Fukunaga] همانطور که پیش‌تر ذکر شد، در رابطه $g(x) = w^T x + w_0$ جزء $w^T x$ بیان می‌کند که داده‌ها باید بر روی بردار w تابانده شوند. و بر روی این ابرصفحه (خط در حالت دو بعدی) که داده‌ها قرار دارند، w_0 مرز جداسازی را تعیین می‌کند (شکل ۱۶).

اگر بردارهای ویژگی بصورت نرمال توزیع شده باشند یا تعداد آنها زیاد باشد، نگاشت آنها نیز نرمال خواهد بود (شکل ۱۶).

اگر بردارهای ویژگی داده شده برای دسته‌بندی، بصورت نرمال توزیع شده باشند ولی تعداد آنها زیاد باشد، نیز، حاصل نگاشت خطی بردارهای ویژگی، نزدیک به نرمال خواهد بود. زیرا در نگاشت خطی، هر بردار ویژگی نگاشت شده، مجموع وزن‌دار ابعاد بردار ویژگی معادلش در فضای چند بعدی بوده و چون حاصلجمع چند ترم می‌باشد، قضیه حد مرکزی^{۱۲} در مورد آن صادق است.

تمرین ۱۰: توابع جداساز خطی برای حالتی که داده‌ها با مقایسه میانگین‌های داده‌های دو کلاس جدایی‌پذیر هستند بهتر عمل می‌کند یا برای حالتی که داده‌ها با مقایسه پراکندگی (واریانس) داده‌های دو کلاس جدایی‌پذیر هستند یا هر دو؟

در شکل ۱۶، مساحت فضاهای برهم‌افتادگی نمودارهای گاوسی نمایانگر میزان خطاها هستند. خطای جداسازی، پس از تابش داده‌ها در جهت w ، کمتر از خطای جداسازی، پس از تابش داده‌ها در جهت v ، می‌باشد.



شکل ۱۶: نگاشت خطی داده‌ها، جداسازی و خطا

جداساز خوب، جداسازی است که جهت (w, w_0) را به قسمی پیدا کند که کمترین میزان خطا را برای جداسازی، پس از تابش داده‌ها در آن جهت (در فضای g)، داشته باشیم.

اگر داده‌ها از یک جمعیت نرمال انتخاب شده باشند، حاصل تابش آنها نیز دارای توزیعی نرمال خواهد بود. میزان خطا نیز بوسیله

¹² Central Limit Theorem

پارامترهای این توزیع نرمال $(\eta_i = E[g(x)|C_i], \sigma_i^2 = Var[g(x)|C_i])$ تعیین خواهد شد که (η_i, σ_i^2) تابعی از (w, w_0) هستند. در این حالت داریم:

$$\begin{aligned}\eta_i &= E[g(x)|C_i] = W^t E[x|C_i] + w_0 = w^t \mu_i + w_0 \\ \sigma_i^2 &= Var[g(x)|C_i] = w^t E[(x - \mu_i)(x - \mu_i)^t | C_i] w = w^t \Sigma_i w\end{aligned}\quad (22)$$

فرض کنید $J(\eta_1, \eta_2, \sigma_1^2, \sigma_2^2)$ معیاری باشد که برای تعیین بهترین (w, w_0) باید بهینه شود. در این صورت می‌بایستی

$$\begin{cases} \frac{\partial J}{\partial w} = \frac{\partial J}{\partial \sigma_1^2} \frac{\partial \sigma_1^2}{\partial w} + \frac{\partial J}{\partial \sigma_2^2} \frac{\partial \sigma_2^2}{\partial w} + \frac{\partial J}{\partial \eta_1} \frac{\partial \eta_1}{\partial w} + \frac{\partial J}{\partial \eta_2} \frac{\partial \eta_2}{\partial w} = 0 \\ \& \frac{\partial J}{\partial w_0} = \frac{\partial J}{\partial \sigma_1^2} \frac{\partial \sigma_1^2}{\partial w_0} + \frac{\partial J}{\partial \sigma_2^2} \frac{\partial \sigma_2^2}{\partial w_0} + \frac{\partial J}{\partial \eta_1} \frac{\partial \eta_1}{\partial w_0} + \frac{\partial J}{\partial \eta_2} \frac{\partial \eta_2}{\partial w_0} = 0 \end{cases}\quad (23)$$

با ترکیب روابط (۲۲) و (۲۳) خواهیم داشت:

$$2 \left[\frac{\partial J}{\partial \sigma_1^2} \Sigma_1 + \frac{\partial J}{\partial \sigma_2^2} \Sigma_2 \right] w + \left[\frac{\partial J}{\partial \eta_1} \mu_1 + \frac{\partial J}{\partial \eta_2} \mu_2 \right] = 0 \quad (24)$$

$$\frac{\partial J}{\partial \eta_1} + \frac{\partial J}{\partial \eta_2} \frac{\partial \eta_2}{\partial w} = 0 \quad (25)$$

با جایگذاری رابطه (۲۵) در رابطه (۲۴) و توجه به این نکته که خطا فقط به جهت w بستگی دارد و نه به سایر آن و همچنین با ساده‌سازی ترم‌های ثابت، خواهیم داشت:

$$\begin{aligned}w &= [s\Sigma_1 + (1-s)\Sigma_2]^{-1} (\mu_2 - \mu_1) \quad ; s = \frac{\partial J / \partial \sigma_1^2}{\partial J / \partial \sigma_1^2 + \partial J / \partial \sigma_2^2} \\ w_0 &= \text{the Solution of } \left\langle \frac{\partial J}{\partial \eta_1} + \frac{\partial J}{\partial \eta_2} = 0 \right\rangle\end{aligned}\quad (26)$$

مثال ۸: با کمینه‌سازی خطا، (w, w_0) را برای حالتی بدست آورید که معیار کمینه‌سازی خطا برابر $J = \frac{(\eta_1 - \eta_2)^2}{\sigma_1^2 + \sigma_2^2}$ (معیار فیشر) باشد.

با توجه به رابطه (۲۶) ابتدا مقدار s را محاسبه می‌کنیم. داریم:

$$\frac{\partial J}{\partial \sigma_1^2} = \frac{\partial J}{\partial \sigma_2^2} = \frac{-(\eta_1 - \eta_2)^2}{(\sigma_1^2 + \sigma_2^2)^2} \Rightarrow s = \frac{1}{2}$$

در نتیجه $w = \left[\frac{1}{2}\Sigma_1 + \frac{1}{2}\Sigma_2 \right]^{-1} (\mu_2 - \mu_1)$ و w_0 نیز باید با توجه به داده‌ها بدست آید (جواب بدست آمده را با جوابی که قبلاً برای فیشر بدست آورده بودیم مقایسه نمایید).

□

تمرین ۱۱: نشان دهید که با کمینه‌سازی خطا، جداساز بهینه برای حالتی که معیار کمینه‌سازی به صورت $J = \frac{p_1 \eta_1^2 - p_2 \eta_2^2}{p_1 \sigma_1^2 + p_2 \sigma_2^2}$ باشد، دارای پارامترهای $\langle w = [p_1 \Sigma_1 + p_2 \Sigma_2]^{-1} (\mu_2 - \mu_1), w_0 = -w^t (p_1 \mu_1 + p_2 \mu_2) \rangle$ خواهد بود.

تمرین ۱۲: با کمینه‌سازی خطا، (w, w_0) را برای حالتی بدست آورید که معیار کمینه‌سازی خطا برابر تابع خطای بیز در فضای تابش (فضای g) باشد. خطای بیز برابر است با:

$$. \varepsilon = p_1 \int_{-\eta_1/\sigma_1}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{\zeta^2}{2}} d\zeta + p_2 \int_{-\infty}^{-\eta_2/\sigma_2} \frac{1}{\sqrt{2\pi}} e^{-\frac{\zeta^2}{2}} d\zeta$$

۴.۲- روش میانگین مربعات خطا

[theodiridis] این بخش را با تمرکز بر حالت دو کلاسه شروع می‌کنیم. برای یک مساله دسته‌بندی دو کلاسه، تابع جداساز خطی همگن $g(x)$ را در نظر بگیرید، که برای کلاس C_1 مقدار $+1$ و برای کلاس C_2 مقدار -1 را بر می‌گرداند. از آنجایی که این تابع جداساز دارای خطا می‌باشد، جواب صحیح را به ازای هر x بر نمی‌گرداند. y را بقسمی در نظر بگیرید که تعلق واقعی داده x را به کلاس‌های C_1 و C_2 به ترتیب با مقادیر $+1$ و -1 نشان دهد (خروجی‌های مطلوب). در این صورت تابع میانگین مربع خطا (MSE^{13}) برابر خواهد بود با:

$$J(w) = E\left[\|y - g(x)\|^2\right] = E\left[\|y - x^t w\|^2\right] \quad (27)$$

w باید طوری تعیین شود که $J(w)$ کمینه شود یا $\hat{w} = \arg \min_w J(w)$. برای کمینه کردن $J(w)$ نیز باید داشته باشیم:

$$\frac{\partial J(w)}{\partial w} = 2E[x(y - x^t w)] = 2E[xy - xx^t w] = 2(E[xy] - E[xx^t]w) = 0$$

در نتیجه:

$$\hat{w} = \frac{E[xy]}{E[xx^t]} = \frac{E[xy]}{R_x} = R_x^{-1} E[xy] \quad (28)$$

که در آن R_x ماتریس همبستگی^{۱۴} بردارهای ویژگی (معادل ماتریس کوواریانس) بوده و $E[xy]$ همبستگی متقابل^{۱۵} بردارهای ویژگی و خروجی مطلوب می‌باشد:

$$R_x = E[xx^t] = \begin{bmatrix} E[x_1 x_1] & \dots & E[x_1 x_n] \\ E[x_2 x_1] & \dots & E[x_2 x_n] \\ \vdots & \vdots & \vdots \\ E[x_n x_1] & \dots & E[x_n x_n] \end{bmatrix}, \quad E[xy] = E \begin{bmatrix} x_1 y \\ x_2 y \\ \vdots \\ x_n y \end{bmatrix}$$

بحث فوق را می‌توان به حالت k کلاسه نیز بسط داد. در این حالت باید k تابع جداساز خطی $g_i(x)$ بر اساس معیار میانگین مربعات خطا طراحی شوند. برای خروجی‌های مطلوب نیز اگر $x \in C_i$ ، باید $y_i = 1$ باشد، در غیر اینصورت $y_i = 0$. در نتیجه برای هر بردار ویژگی x بردار $y^t = [y_1, \dots, y_k]$ را به عنوان خروجی مطلوب خواهیم داشت. در این مساله، هدف پیدا کردن k تابع جداساز یا در اصل یافتن $W = [w_1, \dots, w_k]$ می‌باشد. در این صورت خواهیم داشت:

$$\hat{W} = \arg \min_W E\left[\|y - W^t x\|^2\right] = \arg \min_W E\left[\sum_{i=1}^k (y_i - w_i^t x)^2\right]$$

ک معادل k مساله مستقل کمینه‌سازی میانگین مربعات خطا می‌باشد. و این یعنی اینکه برای طراحی مجموعه توابع جداساز برای حالت k کلاسه بر اساس میانگین مربعات خطا، کافی است که مساله بصورت k زیر مساله کوچکتر دو کلاسه در نظر گرفته شود که برای هر زیر مساله i ام، می‌بایستی داده‌های کلاس C_i به عنوان یک کلاس و سایر داده‌ها، به عنوان کلاس دیگر در نظر گرفته شوند (ماشین جدایز خطی کامل).

۴.۳- روش کمترین میانگین مربعات (الگوریتم Widrow-Hoff)

دیدیم که برای بدست آوردن تحلیلی w از طریق رابطه (۲۸) نیاز به ماتریس همبستگی بردارهای ویژگی و بردار همبستگی متقابل ورودی‌ها و خروجی‌های مطلوب است. بدست آوردن این اطلاعات همیشه راحت نیست، بخصوص اگر تعداد بردارهای ویژگی ورودی بسیار زیاد باشد. در نتیجه مایلیم روشی داشته باشیم که بتوان بدون استفاده از این اطلاعات آماری به کمینه‌سازی رابطه (۲۷) اقدام نمود. برای این منظور می‌توانیم از روش کاهش گرادیان بهره بگیریم.

قانون بروز رسانی متغیرها (روش تک نمونه‌ای) برای این مساله بصورت زیر در می‌آید:

¹³ Mean Square Error

¹⁴ Correlation Matrix

¹⁵ Cross-Correlation

$$w_{(k+1)} = w_{(k)} + \eta_{(k)} x_{(i)} (y_{(i)} - w_{(k)}^t x_{(i)}) \quad (29)$$

که در آن (x_k, y_k) ها جفت‌های بردارهای ویژگی آموزشی و خروجی‌های مطلوب معادل آنها می‌باشند که بصورت متوالی به این رابطه معرفی می‌شوند. این الگوریتم به عنوان کمترین میانگین مربعات (LMS¹⁶) مشهور است که به میانگین مربعات خطا همگرا می‌شود. الگوریتم‌های LMS دیگری نیز با شرایط مختلف پیشنهاد شده‌اند.

توجه کنید که رابطه بازگشتی (29) می‌تواند به عنوان الگوریتم تعیین سیستماتیک وزن‌های یک نرون خطی نیز بکار گرفته شود. این نوع از الگوریتم‌های آموزشی که خروجی مطلوب را به مجموع وزن‌دار ورودی‌های آموزشی نرون اعمال می‌کنند، اولین بار توسط Widrow و Hoff بکار گرفته شدند (الگوریتم Widrow-Hoff) و این ساختار که با نام $adaline^{17}$ شناخته می‌شود، در واقع نرونی است که در آموزش آن به جای پرسپترون از LMS استفاده می‌شود.

۴.۴- روش مجموع مربعات خطا (روش ماتریس شبه معکوس)

یک مساله دسته‌بندی به فرم نرمال را در نظر بگیرید. در این مساله برای تمامی بردارهای ویژگی باید $w^t x_i > 0$ باشد. یا اگر به ازای هر بردار ویژگی x_i عدد مثبت b_i (با هر مقدار دلخواهی) را داشته باشیم، باید برای تمامی بردارهای ویژگی $w^t x_i = b_i$ باشد. از آنجایی که $\|w\| w^t x_i / \|w\|$ فاصله بردار ویژگی x_i از بردار w می‌باشد، می‌توان گفت که b_i ها متناظر فواصل بردارهای ویژگی معادلشان از بردار w (با ضریب ثابت $1/\|w\|$) می‌باشند. پس حل معادله $Xw = b$ و بدست آوردن w ، معادل بدست آوردن $g(x) = w^t x$ می‌باشد.

اگر بخواهیم مساله را به طریق تحلیلی حل کنیم به نتیجه $w = X^{-1}b$ می‌رسیم. برای محاسبه X^{-1} ، ماتریس X باید یک ماتریس مربعی بوده و غیر تکیه^{۱۸} باشد (دترمینان آن مخالف صفر باشد). از آنجایی که در هر سطر یکی از بردارهای ویژگی قرار می‌گیرد، تعداد سطرهای X برابر تعداد بردارهای ویژگی بوده و تعداد ستون‌ها نیز برابر تعداد ابعاد بردارهای ویژگی می‌باشد. مربعی بودن ماتریس X یعنی اینکه تعداد بردارهای ویژگی با تعداد ابعاد برابر باشد. اتفاقی که در عمل به ندرت رخ می‌دهد. اگر مربعی نبودن X بدلیل بیشتر بودن تعداد ابعاد بردارهای ویژگی از تعداد بردارهای ویژگی می‌بود، می‌توانستیم با اعمال یکی از الگوریتم‌های کاهش ابعاد، ماتریس را مربعی کنیم، ولی در این حالت، نمی‌توانیم اقدامی برای رفع مشکل انجام دهیم. در نتیجه علیرغم اینکه روش تحلیلی تضمین می‌کند که ابر صفحه جداساز را بیابد، در عمل قابل استفاده نیست.

بنابراین باید به دنبال یک حل تقریبی برای این معادله باشیم. هر چند که امکان دارد که این حل تقریبی لزوماً ابر صفحه کاملاً جداسازی را بدست ندهد، ولی می‌توان امیدوار بود که جوابی که بدست می‌آید، جواب تقریباً مناسبی باشد. برای این منظور می‌توان به جای پیدا کردن جواب صفر $Xw - b$ سعی نمود که آن را تا حد امکان کمینه نمود. در این صورت مساله به یک مساله بهینه‌سازی با معیار کمینه‌سازی زیر تبدیل می‌شود:

$$J(w) = \|Xw - b\|^2 = \sum_{i=1}^n (w^t x_i - b_i)^2 = \sum_{i=1}^n e_i^2$$

از آنجایی که $Xw - b$ نمایانگر خطا می‌باشد (چرا؟)، معیار فوق را معیار کمینه‌سازی مربعات خطا^{۱۹} می‌نامند. و همچنین بدلیل اینکه مجموع مربعات خطا را نیز کمینه می‌کند به روش تخمین مجموع مربعات خطا نیز مشهور است. در اصل در روش تخمین مجموع مربعات خطا، همانند روش میانگین مربعات خطا عمل می‌کنیم با این تفاوت که تابع هدف به صورت فوق در نظر گرفته می‌شود.

برای بدست آوردن w بهینه داریم:

$$\frac{\partial J(w)}{\partial w} = \sum_{i=1}^n 2x_i (w^t x_i - b_i) = 2X^t (Xw - b) = 0$$

در نتیجه خواهیم داشت:

¹⁶ Least Mean Squares

¹⁷ Adaptive linear element

¹⁸ Singular

¹⁹ Minimum Squared Error Criterion

روش ماتریس شبه معکوس
برای بدست آوردن تخمین
مجموع مربعات خطا

$$X^t X w - X^t b = 0 \Rightarrow X^t X w = X^t b \Rightarrow w = \frac{X^t b}{X^t X} = (X^t X)^{-1} X^t b$$

$X^t X$ ماتریس همبستگی بردارهای ویژگی (معادل ماتریس کوواریانس) می باشد که یک ماتریس مربعی می باشد که دارای $d+1$ سطر و ستون بوده و اغلب غیر تکین می باشد. $X^t X$ را نیز ماتریس شبه وارون X^+ می گویند (چون داریم،
 $((X^t X)^{-1} X^t) X = (X^t X)^{-1} (X^t X) = I$).

در تعیین b_i ها در یک مساله، توجه داشته باشید که اگر برای همه آنها $b_i = 1$ باشد، به ازای همه بردارهای ویژگی $g(x) = w^t x_i = 1$ خواهد شد. w بدست آمده نیز معادل ابر صفحه ای خواهد بود که قادر به جداسازی تمامی بردارهای ویژگی می باشد. پس در یک مساله برای بدست آوردن تابع جداساز، اگر اطلاعاتی از مساله نداشته باشیم، برای تعیین مقادیر b_i ها، می توان در ساده ترین حالات تمامی آنها را برابر ۱ در نظر گرفت.

توضیحاتی که برای حالت چند کلاسه در MSE داده شد، در اینجا نیز صادق است. همچنین لازم به ذکر است که تخمین مجموع مربعات خطا، در n های بزرگ، به سمت MSE میل می کند (چرا؟).

مثال ۹: کلاس C_1 حاوی بردارهای ویژگی $\{[0.2, 0.7]^t, [0.3, 0.3]^t, [0.4, 0.5]^t, [0.6, 0.5]^t, [0.1, 0.4]^t\}$ و کلاس C_2 حاوی بردارهای ویژگی $\{[0.4, 0.6]^t, [0.6, 0.2]^t, [0.7, 0.4]^t, [0.8, 0.6]^t, [0.7, 0.5]^t\}$ را در نظر بگیرید. مرزهای جداساز بهینه را با روش مجموع مربعات خطا طراحی نمایید.

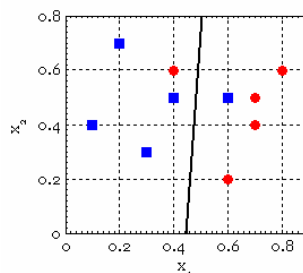
مرز جداساز بهینه به فرم $w_2 x_2 + w_1 x_1 + w_0 = 0$ خواهد بود. در نتیجه هدف پیدا کردن $w = [w_2, w_1, w_0]^t$ بهینه خواهد بود. برای ساختن ماتریس X ابتدا ۱۰ بردار ویژگی داده شده را با در نظر گرفتن ۱ به عنوان بعد سومشان به یک ماتریس 10×3 تبدیل می کنیم. b_i ها را نیز برای بردارهای ویژگی متعلق به کلاس اول برابر ۱ و برای بردارهای ویژگی متعلق به کلاس دوم برابر -۱ در نظر می گیریم. در نتیجه بردار b از ۵ عدد ۱ و ۵ عدد -۱ ساخته می شود. سپس ماتریس همبستگی بردارهای ویژگی و را به صورت زیر محاسبه می کنیم:

$$X^t X = \begin{bmatrix} 2.8 & 2.24 & 4.8 \\ 2.24 & 2.41 & 4.7 \\ 4.8 & 4.7 & 10 \end{bmatrix} \quad \& \quad X^t b = \begin{bmatrix} -1.6 \\ 0.1 \\ 0.0 \end{bmatrix}$$

پس،

$$w = (X^t X)^{-1} X^t b = [-3.218, 0.241, 1.431]^t$$

در شکل زیر، داده های دو کلاس و مرز جداساز بهینه بدست آمده نمایش داده شده اند.



□

روش ماتریس شبه معکوس در صورتی که ابعاد $X^t X$ بالا باشد یا بردارهای ویژگی دارای همبستگی بالایی باشند، دچار مشکل می شود.

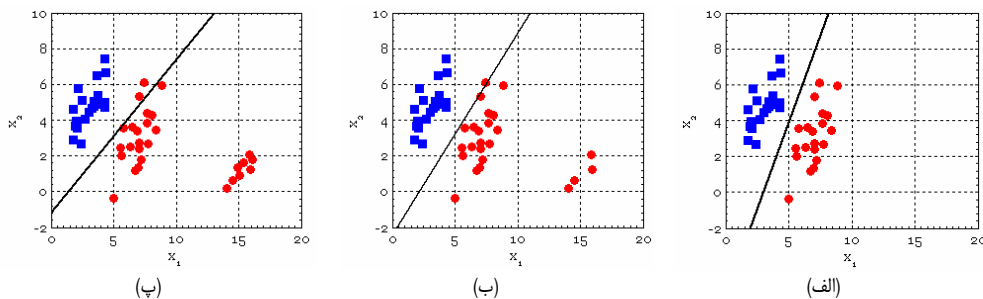
اگر تعداد ابعاد $X^t X$ بالا باشد، محاسبه وارون آن هزینه بر است. همچنین اگر بردارهای ویژگی دارای همبستگی بالایی باشند (اگر سطرهای X ترکیبی خطی از سایر سطرها باشند)، دترمینان $X^t X$ نزدیک به صفر خواهد شد و محاسبه وارون آن مشکل ساز خواهد شد. در این حالت برای حل مساله می توان از روش های بهینه سازی تکراری استفاده نمود. مثلاً در استفاده از الگوریتم کاهش گرادینان، نظر به اینکه $\nabla J(w) = 2X^t(Xw - b)$ می باشد، قانون بروز رسانی تک نمونه ای به صورت زیر در

خواهد آمد:

$$w_{(k+1)} = w_{(k)} - \eta_{(k)} x_{(i)} (w_{(k)}^t x_{(i)} - b_{(i)})$$

که معادل Widrow-Hoff در تخمین میانگین مربعات خطا است.

[Bishop] روش مجموع مربعات خطا پایدار نبوده و نسبت به تغییرات داده‌ها بسیار حساس است. در شکل ۱۷ اثرات اضافه نمودن داده‌ها، بر مرز بدست آمده بوسیله روش مربعات خطا نشان داده شده است دقت کنید که داده‌های دور از مرز، که بیشترین حاشیه امنیت را دارند، منجر به بیشترین تأثیرات بر روی جداساز می‌شوند.



شکل ۱۷: روش مجموع مربعات خطا بسیار حساس به داده می‌باشد. شکل‌های فوق، به ترتیب از الف تا پ، اثرات اضافه نمودن داده‌ها بر مرز جداپذیری را نشان می‌دهند.

متد Ho-Kashyap

در روش تخمین مجموع مربعات خطا، راهکاری برای تعیین مقادیر b_i های بهینه ارائه ندادیم. سوالی که ممکن است پیش بیاید این است که چه مقداری برای خروجی‌های مطلوب دو کلاس در نظر گرفته شوند تا جواب بهتری بدست آید؟

در روش تخمین مجموع مربعات خطا قادر به حل $Xw = b$ نبودیم و با کمینه‌سازی $Xw - b$ فقط تقریبی از آن راه دست آوردیم. یعنی $Xw \approx b$ یا $w^t x_i = b_i + \varepsilon_i$ مقادیر ε_i ها ممکن است منفی نیز باشد. در این حالت اگر ε_i ها اعداد منفی کوچکی باشند، باز هم $w^t x_i$ بزرگتر از صفر بوده و دسته‌بندی به درستی انجام می‌شود. ولی چنانچه تقریب بدست آمده مناسب نبوده و ε_i ها اعداد منفی بزرگی باشند، $w^t x_i$ منفی شده و دسته‌بندی اشتباه انجام می‌شود. نتیجه اینکه تقریب مجموع مربعات خطا ممکن است همیشه جواب بهینه‌ای بدست ندهد، هر چند که اکثر اوقات جواب قابل قبولی را بدست می‌دهد.

مثال ۱۰: آیا اگر b_i ها اعداد بزرگی در نظر گرفته شوند، می‌توان بر مشکل عدم بهینگی روش تخمین مجموع مربعات خطا غلبه کرد؟

خیر. فرض کنید به جای b_i ها اعداد βb_i در نظر گرفته شوند که در آن β عدد مثبت دلخواهی باشد. در این صورت w برابر خواهد بود با:

$$w = \arg \min_w \|Xw - \beta b\|^2 = \arg \min_w \beta^2 \|X(w/\beta) - b\|^2 = \arg \min_w \|X(w/\beta) - b\|^2 = \beta w_{old}$$

در نتیجه اگر برای یکی از داده‌ها و تابع جداساز قبلی داشته بوده باشیم $w_{old}^t x_i < 0$ در تابع جداساز فعلی نیز خواهیم داشت:

$$w^t x_i = \beta w_{old}^t x_i < 0$$

و این یعنی اینکه اندازه b_i ها به تنهایی مهم نیست، بلکه اختلاف آنها نسبت به هم تأثیرگذار است.

□

فرض کنید که در یک مساله که با تخمین مجموع مربعات می‌خواهیم ابر صفحه جداساز را بدست آوریم، بردارهای ویژگی بصورت خطی قابل جداسازی باشند. در این صورت w^* و b^* وجود دارند به طوری که $Xw^* = b^* > 0$ می‌دانیم که با داشتن b^* قادر به یافتن w^* خواهیم بود. اما خود b^* نیز مجهول می‌باشد. پس باید سعی کنیم هم‌زمان هر دو مجهول را بیابیم. یعنی تابع زیر را باید کمینه کنیم:

روش مجموع مربعات خطا پایدار نبوده و نسبت به تغییرات داده‌ها بسیار حساس است.

مجموع مربعات خطا ممکن است همیشه جواب بهینه‌ای بدست ندهد، هر چند که اکثر اوقات جواب قابل قبولی را بدست می‌دهد.

$$J(w, b) = \|Xw - b\|^2; b > 0$$

بهینه‌سازی تابع هدف نسبت به دو متغیر به کمک روش‌های بهینه‌سازی تکراری

حل این معادله با روش‌های تحلیلی ممکن نیست. به همین دلیل از الگوریتم‌های دیگری برای حل این مساله بهره می‌گیریم. برای این منظور باید دو گام زیر تا رسیدن به همگرایی برداشته شوند:

۱. ثابت نگه داشتن b و کمینه کردن $J(w, b)$ نسبت به w .

۲. ثابت نگه داشتن w و کمینه کردن $J(w, b)$ نسبت به b .

گام اول با روش ماتریس شبه وارون قابل انجام است. برای یک b ثابت، مقدار w برابر خواهد بود با:

$$\nabla_w J(w, b) = 2X^t(Xw - b) = 0 \Rightarrow w = (X^t X)^{-1} X^t b$$

انجام گام دوم کمی مشکل‌تر است. زیرا $\nabla_b J(w, b) = -2(Xw - b) = 0$ نتیجه $b = Xw$ را در بر خواهد داشت که چون محدودیت مثبت بودن اجزای b در آن لحاظ نشده است قابل استفاده نیست. سعی می‌کنیم الگوریتم کاهش گرادیان را با اعمال تغییراتی برای حل این مساله به کار ببریم. قانون بروز رسانی تک نمونه‌ای در حالت معمولی برای این مساله بصورت زیر خواهد بود:

$$b_{(k+1)} = b_{(k)} + \eta_{(k)} 2(Xw_{(k)} - b_{(k)})$$

هر کدام از درایه‌های $Xw_k - b_k$ در صورت منفی شدن می‌تواند منجر به منفی شدن درایه‌ای از b شود. برای رفع این مشکل تمامی درایه‌های منفی $Xw_k - b_k$ را صفر می‌کنیم. یعنی،

$$b_{(k+1)} = b_{(k)} + \eta_{(k)} [(Xw_{(k)} - b_{(k)}) + |(Xw_{(k)} - b_{(k)})|]$$

با انجام این عمل کماکان با کاهش گرادیان به سمت مقدار بهینه نزدیک می‌شویم ولی سرعت کاهش گرادیان و رسیدن به جواب بهینه کم شده است. حالت نهایی روش فوق که به روال Ho-Kashyap معروف است، به صورت زیر می‌باشد:

روال Ho-Kashyap

۱. مقادیر اولیه‌ای برای $w_{(0)}$ و $b_{(0)} > 0$ در نظر بگیر و مقدار اولیه k را نیز برابر صفر قرار بده

۲. مقدار $b_{(k+1)}$ را از قانون $b_{(k+1)} = b_{(k)} + \eta_{(k)} [(Xw_{(k)} - b_{(k)}) + |(Xw_{(k)} - b_{(k)})|]$ بدست بیاور.

۳. مقدار $w_{(k+1)}$ را از رابطه $w_{(k+1)} = (X^t X)^{-1} X^t b_{(k+1)}$ بدست بیاور.

۴. k را برابر $k+1$ قرار بده.

قدم‌های فوق را تا وقتی $Xw_{(k)} - b_{(k)} \geq 0$ یا $k > k_{\max}$ یا $b_{(k+1)} = b_{(k)}$ تکرار کن.

توجه کنید که برای همگرایی باید نرخ یادگیری دارای مقادیری با شرط $0 < \eta < 1$ باشد. چرا؟

تمرین ۱۳: شرایط پایانی روال Ho-Kashyap را بررسی کرده و برای هر مورد بیان کنید که چرا این شرط قرار داده شده است و در چه وضعیتی این شرایط برقرار خواهند شد.

تمرین ۱۴: روال Ho-Kashyap در حالتی که داده‌ها به صورت خطی جداپذیر نباشند، چه رفتاری از خود نشان خواهد داد؟

تمرین ۱۵: آیا شرایطی وجود دارد که در آن روال Ho-Kashyap همگرا نشود؟

تمرین ۱۶: نشان دهید که روال Ho-Kashyap برای حالت k کلاسه، متناظر k مساله دو کلاسه خواهد بود.

۵- جمع‌بندی و مباحث تکمیلی

در این قسمت به مقایسه روش‌های مطرح شده می‌پردازیم. همچنین برخی از موضوعات را که به مباحث فوق مرتبط بوده ولی فرصت اشاره به آنها در طی مباحث مطرح شده بوجود نیامد را در اینجا مطرح می‌کنیم.

۵.۱- مقایسه روش‌ها

در جدول زیر روش‌هایی را که بررسی کردیم، با معیارهایی که بر اساس آنها عمل می‌کنند و همچنین جواب نهایی‌ای که به دست می‌آورند را آورده‌ایم:

روش	معیار	الگوریتم یافتن پاسخ
پرسپترون	$J(w) = \sum_{x \in X_m(w)} -w^t x$	$w_{(k+1)} = w_{(k)} + \eta_{(k)} x_i$
Relaxation	$J(w) = \frac{1}{2} \sum_{x \in X_b(w)} \frac{(w^t x - b)^2}{ x ^t}$	$w_{(k+1)} = w_{(k)} - \eta_{(k)} \frac{w_{(k)}^t x}{ x ^2} x$
فیشر	$J(w) = \frac{w^t S_b w}{w^t S_w w}$	$\bar{w} = S_w^{-1}(\mu_1 - \mu_2)$
ماشین‌های بردار پشتیبان	$\max_{\lambda} L(\lambda) = \sum_{\forall i} \lambda_i - \frac{1}{2} \sum_{\forall i} \sum_{\forall j} \lambda_i \lambda_j z_i z_j x_i^t x_j$ subject to: (1) $\sum_{\forall i} \lambda_i z_i = 0$ & (2) $\lambda \geq 0$	بدست آوردن λ_i های نا صفر: برنامه‌ریزی درجه دوم $w = \sum_{\forall i} \lambda_i z_i x_i$ & $w_0 = 1/z_i - w^t x_i$
کمینه‌سازی خطا برای داده‌های با توزیع نرمال	$J(\eta_1, \eta_2, \sigma_1^2, \sigma_2^2)$	$s = \frac{\partial J / \partial \sigma_1^2}{\partial J / \partial \sigma_1^2 + \partial J / \partial \sigma_2^2}$ $w = [s \Sigma_1 + (1-s) \Sigma_2]^{-1} (\mu_2 - \mu_1)$ $w_0 = \text{the Solution of } \left\langle \frac{\partial J}{\partial \eta_1} + \frac{\partial J}{\partial \eta_2} = 0 \right\rangle$
میانگین مربعات خطا	$J(w) = E \left[\ y - x^t w\ ^2 \right]$	$w = R_x^{-1} E[xy]$
کمترین میانگین مربعات (Widrow-Hoff)	$J(w) = E \left[\ y - x^t w\ ^2 \right]$	$w_{(k+1)} = w_{(k)} + \eta_{(k)} x_{(i)} (y_{(i)} - w_{(k)}^t x_{(i)})$
مجموع مربعات خطا (ماتریس شبه معکوس)	$J(w) = \sum_{i=1}^n (w^t x_i - b_i)^2 = \sum_{i=1}^n e_i^2$	$w = (X^t X)^{-1} X^t b$
مجموع مربعات خطا (Widrow-Hoff)	$J(w) = \sum_{i=1}^n (w^t x_i - b_i)^2 = \sum_{i=1}^n e_i^2$	$w_{(k+1)} = w_{(k)} - \eta_{(k)} x_{(i)} (w_{(k)}^t x_{(i)} - b_i)$
Ho-Kashyap	$J(w, b) = \ Xw - b\ ^2 ; b > 0$	$b_{(k+1)} = b_{(k)} + \eta_{(k)} [(Xw_{(k)} - b_k) + (Xw_{(k)} - b_{(k)})]$ $w_{(k+1)} = (X^t X)^{-1} X^t b_{(k+1)}$

پرسپترون در در حالتی که بردارهای ویژگی به صورت خطی جداپذیر باشند، همگرا شده و ابرصفحه جداساز را می‌یابد. ولی وقتی داده‌ها بصورت خطی جدا پذیر نباشند، همگرا نمی‌شود. با ترکیب چند لایه‌ای از نرون‌های متصل به هم می‌توان مرزهای جداسازی پیچیده‌تری را نیز بدست آورد و به حل مسایل جداپذیر غیر خطی نیز پرداخت.

روش‌های بر پایه کمینه‌سازی خطا در حالتی که داده‌ها به صورت خطی جداپذیر نباشند نیز همگرا می‌شوند، ولی تضمین نمی‌کنند که جواب بهینه را بیابند. ولی متد Ho-Kashyap که فاصله داده‌ها تا مرز جداپذیری را نیز خود تعیین می‌کند، اگر چه از لحاظ زمان اجرایی هزینه‌برتر است، ولی همیشه به یک جواب بهینه همگرا می‌شود. این روش‌ها در مقابل نویزهای تکین

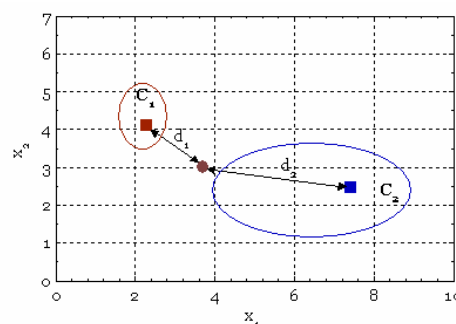
بسیار حساس می‌باشند.

روش ماشین بردار پشتیبان نیز بهترین جواب ممکن را برای تعمیم به دست می‌دهد. این روش می‌تواند برای جداسازی داده‌های غیر خطی نیز به کار گرفته شود که در فصل‌های بعدی به آن خواهیم پرداخت.

۵.۲- مباحث متفرقه

یک کلاس و چند نماینده

[Bow&Bow] در روش دسته‌بندی بر اساس کمترین فاصله، یک بردار ویژگی از هر کلاس (یا یک نقطه در آن فضا) به عنوان نماینده آن کلاس در نظر گرفته شده و داده‌های آزمایشی بر اساس مقایسه فاصله با این بردار ویژگی (نقطه) دسته‌بندی می‌شوند. این کار در صورتی معتبر می‌باشد که کلاس‌ها دارای توزیعی نرمال با کوواریانس یکسان باشند. اگر این شرط برقرار نباشد، دسته‌بندی دچار اشتباه خواهد شد (شکل ۱۸).



شکل ۱۸: در صورت عدم برقراری شرایطی خاص، دسته‌بندی بر اساس کمترین فاصله، نتایج اشتباهی ارائه خواهد داد.

برای رفع این مشکل می‌توان از هر کلاس چندین بردار ویژگی (نقطه) را به عنوان نماینده کلاس در نظر گرفت. در این حالت فاصله یک بردار ویژگی از یک کلاس، برابر فاصله آن بردار ویژگی از نزدیک‌ترین نماینده کلاس به خودش خواهد بود. یعنی:

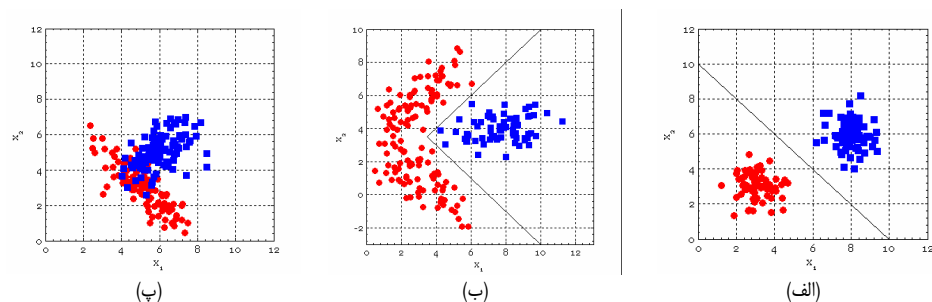
$$d(x, w_k) = \min_{m=1, \dots, N_k} \{d(x, z_k^m)\} \quad (30)$$

که در آن k نشان دهنده کلاس k ام، m نشان دهنده نماینده m ام و N_k تعداد نماینده‌های کلاس k ام می‌باشد.

تمرین ۱۷: یک مساله دسته‌بندی با k کلاس و دسته‌بندی بر اساس کمترین فاصله را در نظر بگیرید. برای جلوگیری از خطاهای احتمالی از هر کلاس چندین نماینده انتخاب می‌شود. رابطه مرزهای جداساز را برای این مساله بیان نمایید.

جداسازی خطی تکه‌ای

[Bow&Bow] در برخی حالات بردارهای ویژگی داده شده به صورت خطی جداپذیر نیستند ولی بصورت خطی تکه‌ای جداپذیرند. یک تابع جداپذیر خطی تکه‌ای، تابعی است که در زیر نواحی‌های مختلف از فضای ویژگی، خطی می‌باشد، اما در کل فضا خطی نیست. این تابع مرزهای خطی تکه‌ای را در زیرنواحی‌های مختلف بدست می‌دهد (شکل ۱۹).



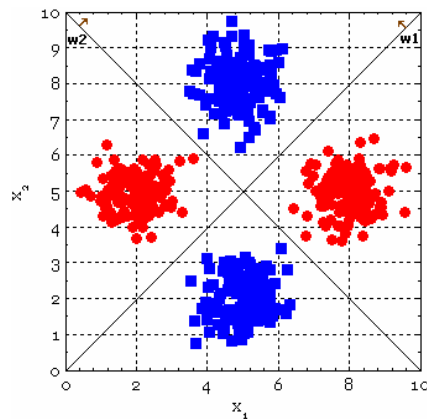
شکل ۱۹: جداپذیری خطی (الف)، جداپذیری خطی تکه‌ای (ب) و جداناپذیری (پ)

روش‌های دسته‌بندی بر اساس کمترین فاصله در صورتی معتبر می‌باشند که کلاس‌ها دارای توزیعی نرمال با کوواریانس یکسان باشند.

یک تابع جداپذیر خطی تکه‌ای، تابعی است که در زیر نواحی‌های مختلف از فضای ویژگی، خطی می‌باشد، اما در کل فضا خطی نیست

جداسازی و ماشین‌های چند لایه‌ای

[Bow&Bow] دو خط متقاطع در یک صفحه، آن را به ۴ ناحیه مختلف تقسیم می‌کنند. ممکن است این ۴ ناحیه دو بدو طوری به دو کلاس مختلف تعلق داشته که بصورت خطی جداپذیر نباشند (شکل ۲۰).



شکل ۲۰: تقسیم فضا با دو خط به دو کلاس خطی جدا ناپذیر

این داده‌ها را می‌توان با یک جداساز دو لایه (ماشین دو لایه) از هم تفکیک نمود. لایه اول در این ماشین، ارائه کننده خطوط جداساز بوده و لایه دوم در مورد نحوه تقسیم‌بندی نواحی بدست آمده تصمیم‌گیری می‌کند.

با توجه به توضیحات ارائه شده، لایه اول ماشین دولایه حاوی تعدادی ماشین جداساز خطی معمولی می‌باشد. فرض کنید $g_1(x)$ و $g_2(x)$ توابع جداساز در این لایه باشند، که w_1 و w_2 مرزهای جداساز متناظر آنها می‌باشند که در شکل فوق نشان داده شده‌اند. همچنین فرض کنید برای هر مرز جداساز، ناحیه‌ای که با پیکان علامت خورده است، ناحیه‌ای باشد که تابع جداساز متناظر آن، در آن ناحیه مقدار مثبتی را برمی‌گرداند. در این صورت برای چهار ناحیه تشکیل شده خواهیم داشت:

	A	B	C	D
$g_1(x)$	+	-	-	+
$g_2(x)$	+	+	-	-
$x \in$	C_1	C_2	C_1	C_2

در این مثال، برای لایه دوم نیز که وظیفه تصمیم‌گیری را بر عهده دارد، به راحتی می‌توان یک ضرب کننده را قرار داد. این ضرب کننده خروجی‌های دو تابع جداساز را در هم ضرب نموده و بر اساس نتیجه بدست آمده تصمیم‌گیری می‌کند. مقدار برگشتی مثبت توسط ضرب کننده نمایانگر تعلق داده به کلاس C_1 و مقدار منفی نشان دهنده تعلق داده به کلاس C_2 می‌باشد.

در یک ماشین جداساز دولایه،
لایه اول ارائه کننده خطوط
جداساز بوده و لایه دوم در
مورد نحوه تقسیم‌بندی نواحی
بدست آمده تصمیم‌گیری
می‌کند.