

فصل هفتم

انتخاب مشخصه‌ها

6-1: روشهای مبتنی بر انتخاب ویژگی

مساله انتخاب ویژگی، یکی از مسائلی است که در مبحث یادگیری ماشین و همچنین شناسایی آماری الگو مطرح است. این مساله در بسیاری از کاربردها (مانند طبقه‌بندی) اهمیت بسزایی دارد، زیرا در این کاربردها تعداد زیادی ویژگی وجود دارند که بسیاری از آنها یا بلا استفاده هستند و یا بار اطلاعاتی چندانی ندارند. حذف نکردن این ویژگی‌ها مشکلی از لحاظ اطلاعاتی ایجاد نمی‌کند ولی بار محاسباتی را برای کاربرد مورد نظر بالا می‌برد و علاوه بر این باعث می‌شود که اطلاعات غیر مفید زیادی را به همراه داده‌های مفید ذخیره کنیم.

برای مساله انتخاب ویژگی، راه حل‌ها و الگوریتم‌های فراوانی ارائه شده است که بعضی از آنها از قدمت سی یا چهل ساله برخوردارند. مشکل بعضی از الگوریتم‌ها در زمانی که ارائه شده بودند، بار محاسباتی زیاد آنها بود، اگر چه امروزه با ظهور کامپیوترهای سریع و منابع ذخیره‌سازی بزرگ این مشکل به چشم نمی‌آید ولی از طرف دیگر، مجموعه داده‌های بسیار بزرگ برای مسائل جدید باعث شده است که هم‌چنان پیدا کردن یک الگوریتم سریع برای این کار مهم باشد.

در این بخش ما در ابتدا تعاریفی که برای انتخاب ویژگی ارائه شده‌اند و همچنین، تعاریف مورد نیاز برای درک این مساله را ارائه می‌دهیم. سپس روش‌های مختلف برای این مساله را بر اساس نوع و ترتیب تولید زیرمجموعه ویژگی‌های کاندید و همچنین نحوه ارزیابی این زیرمجموعه‌ها دسته‌بندی می‌کنیم. سپس تعدادی از روش‌های معرفی شده در هر دسته را معرفی و بر اساس اهمیت، تا جایی که مقدور باشد، آنها را تشریح و الگوریتم برخی از آنها را ذکر می‌کنیم. لازم به ذکر است که به دلیل اینکه مبحث انتخاب ویژگی به مبحث طبقه‌بندی بسیار نزدیک است، بعضی از مسائلی که در اینجا مطرح می‌شود مربوط به مبحث طبقه‌بندی می‌باشد. توضیحات ارائه‌شده برای الگوریتم‌های مختلف در حد آشنائی است. شما می‌توانید برای کسب اطلاعات بیشتر به منابع معرفی شده مراجعه کنید.

6-2: تعاریف

مساله انتخاب ویژگی به وسیله نویسندگان مختلف، از دیدگاه‌های متفاوتی مورد بررسی قرار گرفته است. هر نویسنده نیز با توجه به نوع کاربرد، تعریفی از آن ارائه داده است. در ادامه چند مورد از این تعاریف را بیان می‌کنیم:

1. **تعریف ایده‌آل:** پیدا کردن یک زیرمجموعه با حداقل اندازه ممکن برای ویژگی‌ها است، که برای هدف مورد نظر اطلاعات لازم و کافی را در بر داشته باشد. بدیهی است که هدف تمام الگوریتم‌ها و روش‌های انتخاب ویژگی همین زیرمجموعه است.
2. **تعریف کلاسیک:** انتخاب یک زیرمجموعه M عنصری از میان N ویژگی، به طوریکه $M < N$ باشد و تابع معیار¹ برای زیرمجموعه مورد نظر، نسبت به سایر زیرمجموعه‌های هم‌اندازه دیگر بهینه باشد. این تعریفی است که Narenda و Fukunaga در سال 1977 ارائه داده‌اند.
3. **افزایش دقت پیشگویی:** هدف انتخاب ویژگی این است که یک زیرمجموعه از ویژگی‌ها برای افزایش دقت پیشگویی انتخاب شوند. به عبارت دیگر، هدف انتخاب ویژگی، کاهش اندازه ساختار بدون کاهش قابل ملاحظه در دقت پیشگویی طبقه‌بندی‌کننده‌ای که با استفاده از ویژگی‌های داده شده بدست می‌آید، می‌باشد.
4. **تخمین توزیع کلاس اصلی:** هدف از انتخاب ویژگی این است که یک زیرمجموعه کوچک از ویژگی‌ها انتخاب شوند، توزیع ویژگی‌هایی که انتخاب می‌شوند، بایستی تا حد امکان به توزیع کلاس اصلی با توجه به تمام مقادیر ویژگی‌های انتخاب شده نزدیک باشد.

¹ Criterion Function

روش‌های مختلف انتخاب ویژگی، تلاش می‌کنند تا از میان 2^N زیر مجموعه کاندید، بهترین زیرمجموعه را پیدا کنند. در تمام این روشها بر اساس کاربرد و نوع تعریف، زیر مجموعه‌ای به عنوان جواب انتخاب می‌شود، که بتواند مقدار یک تابع ارزیابی را بهینه کند. با وجود اینکه هر روشی سعی می‌کند که بتواند، بهترین ویژگی‌ها را انتخاب کند، اما با توجه به وسعت جواب‌های ممکن، و اینکه این مجموعه‌های جواب بصورت توانی با N افزایش پیدا می‌کنند، پیدا کردن جواب بهینه مشکل و در N های متوسط و بزرگ بسیار پر هزینه است.

به طور کلی روش‌های مختلف انتخاب ویژگی را بر اساس نوع جستجو به دسته‌های مختلفی تقسیم‌بندی می‌کنند. در بعضی روش‌ها تمام فضای ممکن جستجو می‌گردد. در سایر روش‌ها که می‌تواند مکاشفه‌ای و یا جستجوی تصادفی باشد، در ازای از دست دادن مقداری از کارائی، فضای جستجو کوچکتر می‌شود.

برای اینکه بتوانیم تقسیم بندی درستی از روش‌های مختلف انتخاب ویژگی داشته باشیم، به این صورت عمل می‌کنیم که فرآیند انتخاب ویژگی در تمامی روش‌ها را به بخش‌های زیر تقسیم می‌کنیم:

1. **تابع تولید کننده²:** این تابع زیر مجموعه‌های کاندید را برای روش مورد نظر پیدا می‌کند.

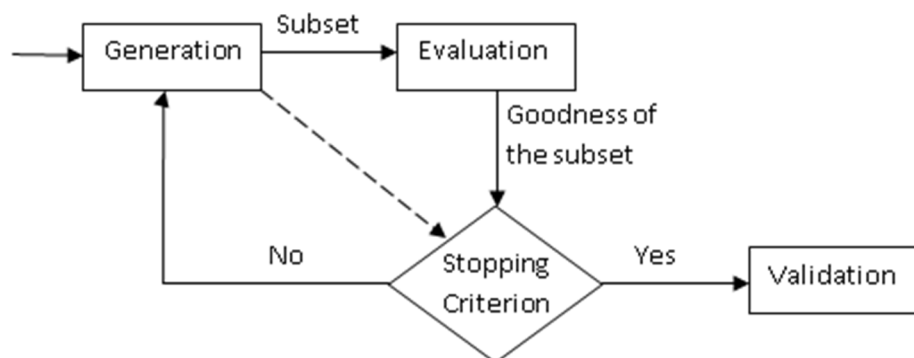
2. **تابع ارزیابی³:** زیرمجموعه مورد نظر را بر اساس روش داده شده، ارزیابی و یک عدد به عنوان میزان خوبی روش باز می‌گرداند. روش‌های مختلف سعی در یافتن زیرمجموعه‌ای دارند که این مقدار را بهینه کند.

² Generation procedure

³ Evaluation function

3. شرط خاتمه: برای تصمیم‌گیری در مورد زمان توقف الگوریتم.

4. تابع تعیین اعتبار⁴: تصمیم می‌گیرد که آیا زیر مجموعه انتخاب شده معتبر است یا خیر؟



شکل 6-1 فرآیند انتخاب ویژگی

تابع تولیدکننده در واقع تابع جستجو است. این تابع زیرمجموعه‌های مختلف را به ترتیب تولید می‌کند، تا به وسیله تابع ارزیابی، مورد ارزیابی قرار بگیرد. تابع تولیدکننده از یکی از حالت‌های زیر شروع به کار می‌کند:

- (1) بدون ویژگی
- (2) با مجموعه تمام ویژگی‌ها
- (3) با یک زیرمجموعه تصادفی

در حالت اول ویژگی‌ها به ترتیب به مجموعه اضافه می‌شوند و زیرمجموعه‌های جدید را تولید می‌کنند. این عمل آنقدر تکرار می‌شود تا به زیر مجموعه مورد نظر برسیم. به این گونه روش‌ها، روش‌های پائین به بالا می‌گویند.

⁴ Validation procedure

در حالت دوم از یک مجموعه شامل تمام ویژگی‌ها، شروع می‌کنیم و به مرور و در طی اجرای الگوریتم، ویژگی‌ها را حذف می‌کنیم، تا به زیرمجموعه دلخواه برسیم. روش‌هایی که به این صورت عمل می‌کنند، روش‌های بالا به پائین نام دارند.

یک تابع ارزیابی، میزان خوب بودن یک زیرمجموعه تولید شده را بررسی کرده و یک مقدار به عنوان میزان خوب بودن زیرمجموعه موردنظر بازمی‌گرداند. این مقدار با بهترین زیرمجموعه قبلی مقایسه می‌شود. اگر زیر مجموعه جدید، بهتر از زیرمجموعه‌های قدیمی باشد، زیرمجموعه جدید به عنوان زیرمجموعه بهینه، جایگزین قبلی می‌شود.

باید توجه داشت که بدون داشتن یک شرط خاتمه مناسب، فرآیند انتخاب ویژگی ممکن است، برای همیشه درون فضای جستجو، برای یافتن جواب سرگردان بماند. شرط خاتمه می‌تواند بر پایه تابع تولیدکننده باشد، مانند:

(1) هر زمان که تعداد مشخصی ویژگی انتخاب شدند.

(2) هر زمان که به تعداد مشخصی تکرار رسیدیم.

و یا اینکه بر اساس تابع ارزیابی انتخاب شود، مانند:

(1) هر زمان که اضافه یا حذف کردن ویژگی، زیرمجموعه بهتری را تولید نکند

(2) هر زمان که به یک زیرمجموعه بهینه بر اساس تابع ارزیابی برسیم.

تابع تعیین اعتبار، جزئی از فرآیند انتخاب ویژگی نیست، اما در عمل بایستی یک زیرمجموعه ویژگی را در شرایط مختلف تست کنیم تا ببینیم که آیا شرایط مورد نیاز، برای حل مساله مورد نظر ما را دارد یا نه؟ برای این کار می‌توانیم از داده‌های نمونه‌برداری شده و یا مجموعه داده‌های شبیه‌سازی شده استفاده کنیم.

6-3: روش‌های مختلف انتخاب ویژگی

در این بخش ابتدا روش‌های مختلف انتخاب ویژگی را بر اساس دو معیار تابع تولید کننده و تابع ارزیابی طبقه‌بندی می‌کنیم. سپس آنها را بر اساس عملکرد دسته‌بندی و نحوه اجرای هر دسته را به اختصار شرح می‌دهیم.

6-3-1: توابع تولید کننده

اگر تعداد کل ویژگی‌ها برابر N باشد، تعداد کل زیرمجموعه‌های ممکن برابر 2^N می‌شود. این تعداد حتی برای N های متوسط خیلی زیاد است. بر اساس نحوه جستجو در میان این تعداد زیر مجموعه، روش‌های مختلف انتخاب ویژگی را می‌توان به سه دسته زیر تقسیم‌بندی نمود:

- (1) جستجوی کامل
- (2) جستجوی مکاشفه‌ای
- (3) جستجوی تصادفی

در ادامه به معرفی هر کدام از این دسته‌ها می‌پردازیم.

6-3-1-1: جستجوی کامل

در روش‌هایی که از این نوع جستجو استفاده می‌کنند، تابع تولید کننده بر اساس تابع ارزیابی استفاده شده، تمام فضای جواب (زیرمجموعه‌های ممکن) را برای یافتن جواب بهینه جستجو می‌کند. البته Schlimmer استدلال آورده است که "کامل بودن جستجو به این معنی نیست که جستجو باید جامع باشد".

توابع مکاشفه‌ای مختلف زیادی طراحی شده‌اند تا جستجو را بدون از دست دادن شانس پیدا کردن جواب بهینه، کاهش دهند. اما با توجه به بزرگی فضای جستجو که از مرتبه $O(2^N)$ می‌باشد، این روش‌ها باعث می‌شوند که فضای کمتری جستجو شود. روش‌ها و تکنیک‌های مختلفی برای این کار استفاده شده‌اند، بعضی از آنها از تکنیک بازگشت به عقب⁵ نیز در جریان کار استفاده کرده‌اند، مانند: beam search و best first search، branch and bound.

2-1-3-6: جستجوی مکاشفه‌ای

در روش‌هایی با این نوع جستجو، در هر بار اجرای الگوریتم، یک ویژگی به مجموعه ویژگی انتخاب شده، اضافه و یا از آن حذف می‌شود. به همین دلیل پیچیدگی زمانی آنها محدود و کمتر از $O(N^2)$ می‌باشد. در این‌گونه موارد، اجرای الگوریتم خیلی سریع می‌باشد و پیاده‌سازی آن‌ها نیز بسیار ساده است.

3-1-3-6: جستجوی تصادفی

روش‌هایی که از این نوع جستجو استفاده می‌کنند، محدوده کمتری از فضای کل حالات را جستجو می‌کنند که اندازه این محدوده به حداکثر تعداد تکرار الگوریتم بستگی دارد. در این روش‌ها پیدا شدن جواب بهینه به اندازه منابع موجود و زمان اجرای الگوریتم بستگی دارد. در هر بار تکرار، تابع تولیدکننده تعدادی از زیرمجموعه‌های ممکن از فضای جستجو را به صورت تصادفی انتخاب می‌کند و در اختیار تابع ارزیابی قرار می‌دهد. تابع تولیدکننده تصادفی پارامتر-هایی دارد که بایستی تنظیم شوند، تنظیم مناسب این پارامترها در سرعت رسیدن به جواب و یافتن جواب‌های بهتر موثر است.

⁵ Backtracking

6-3-2: تابع ارزیابی

یافتن یک زیرمجموعه بهینه از مجموعه ویژگی‌ها، به صورت مستقیم به انتخاب تابع ارزیابی بستگی دارد. چرا که اگر تابع ارزیابی به زیرمجموعه ویژگی بهینه یک مقدار نامناسب نسبت دهد، این زیرمجموعه هیچ‌گاه به عنوان زیرمجموعه بهینه انتخاب نمی‌شود. مقادیری که توابع ارزیابی مختلف به یک زیرمجموعه نسبت می‌دهند، با هم متفاوت است. به طور کلی توابع ارزیابی را می‌توان به 5 گروه تقسیم کرد:

1. معیارهای مبتنی بر فاصله⁶: در این معیارها، مثلاً برای یک مساله دو کلاسه، یک ویژگی یا یک مجموعه ویژگی مانند X بر یک ویژگی یا یک مجموعه ویژگی دیگر مانند Y ارجحیت دارد، اگر با آن مجموعه ویژگی مقادیر بزرگتری برای اختلاف بین احتمالات شرطی دو کلاس داشته باشیم. نمونه‌ای از این معیارها همان معیار فاصله اقلیدسی می‌باشد.

2. معیارهای مبتنی بر اطلاعات⁷: این معیارها میزان اطلاعاتی را که به وسیله یک ویژگی به دست می‌آید، در نظر می‌گیرند. ویژگی X در این روش‌ها بر ویژگی Y اولویت دارد، اگر اطلاعات به دست آمده از ویژگی X بیشتر از اطلاعاتی باشد، که از ویژگی Y به دست می‌آید. نمونه‌ای از این معیارها همان معیار آنتروپی می‌باشد.

3. معیارهای مبتنی بر وابستگی⁸: این معیارها که با عنوان معیارهای همبستگی⁹ نیز شناخته می‌شوند، قابلیت پیشگویی مقدار یک متغیر به وسیله یک متغیر دیگر را اندازه‌گیری می‌کنند. ضریب¹⁰ یکی از معیارهای وابستگی کلاسیک است که می‌توان از

⁶ Distance Measures

⁷ Information Measures

⁸ Dependence Measures

⁹ Correlation

¹⁰ Coefficient

آن برای یافتن همبستگی بین یک ویژگی و یک کلاس استفاده کرد. اگر همبستگی ویژگی X با کلاس C بیشتر از همبستگی ویژگی Y با کلاس C باشد، در این صورت ویژگی X بر ویژگی Y برتری دارد. با یک تغییر کوچک، می‌توانیم وابستگی یک ویژگی با ویژگی‌های دیگر را اندازه‌گیری کنیم. این مقدار درجه افزونگی این ویژگی را نشان می‌دهد. همه توابع ارزیابی بر پایه معیار وابستگی را می‌توانیم به دو دسته معیارهای مبتنی بر فاصله و اطلاعات تقسیم کنیم. اما به خاطر این که این روش‌ها از یک وجه دیگر به مساله نگاه می‌کنند، این کار را انجام نمی‌دهیم.

4. معیارهای مبتنی بر سازگاری¹¹: این معیارها جدیدتر هستند و اخیراً توجه بیشتری به آنها شده است. این معیارها خصوصیات متفاوتی نسبت به سایر معیارها دارند، زیرا به شدت به داده‌های آموزشی تکیه دارند و برای انتخاب یک زیرمجموعه از ویژگی‌ها تمایل دارند که مجموعه ویژگی‌های کوچکتری را انتخاب کنند. این روش‌ها زیرمجموعه‌هایی با کمترین اندازه را بر اساس از دست دادن یک مقدار قابل قبول سازگاری که توسط کاربر تعیین می‌شود، پیدا می‌کنند.

5. معیارهای مبتنی بر خطای طبقه‌بندی‌کننده¹²: روش‌هایی که این نوع از تابع ارزیابی را استفاده می‌کنند، با عنوان "**wrapper methods**" شناخته می‌شوند. دقت عملکرد در این روش‌ها برای تعیین کلاسی که نمونه داده شده متعلق به آن است، برای نمونه‌های دیده نشده بسیار بالا است، اما هزینه‌های محاسباتی در آنها نیز نسبتاً زیاد است.

در جدول زیر مقایسه‌ای بین انواع مختلف تابع ارزیابی، صرف نظر از نوع تابع تولیدکننده مورد استفاده، انجام شده است. پارامترهایی که برای مقایسه استفاده شده‌اند به صورت زیر می‌باشند:

¹¹ Consistency Measures

¹² Classifier Error Rate Measures

1. **عمومیت¹³**: اینکه بتوان زیرمجموعه انتخاب شده را برای طبقه‌بندی‌کننده‌های متفاوت به کار ببریم.

2. **پیچیدگی زمانی**: زمان لازم برای پیدا کردن زیرمجموعه ویژگی جواب.

3. **دقت**: دقت پیشگویی با استفاده از زیرمجموعه انتخاب شده.

علامت "---" که در ستون آخر آمده است، به این معنی است که در مورد میزان دقت حاصل نمی‌توانیم مطلبی بگوئیم. به جز خطای طبقه‌بندی‌کننده، دقت سایر توابع ارزیابی به مجموعه داده مورد استفاده و طبقه‌بندی‌کننده‌ای که بعد از انتخاب ویژگی برای طبقه‌بندی کلاس‌ها استفاده می‌شود، بستگی دارد.

نوع تابع ارزیابی	عمومیت	پیچیدگی زمانی	دقت
معیار فاصله	دارد	پائین	---
معیار اطلاعات	دارد	پائین	---
معیار وابستگی	دارد	پائین	---
معیار سازگاری	دارد	متوسط	---
خطای طبقه بندی کننده	ندارد	بالا	خیلی زیاد

جدول 6-1 مقایسه توابع ارزیابی مختلف

6-3-3: دسته بندی و تشریح الگوریتم های مختلف انتخاب ویژگی

در این قسمت بر اساس تابع ارزیابی و تابع تولیدکننده، روش‌های مختلف انتخاب ویژگی را به چند دسته تقسیم‌بندی می‌کنیم و سپس تعدادی از روش‌ها را شرح داده و الگوریتم کار را به صورت شبهه کد، ذکر می‌کنیم.

¹³ Generality

قبل از این که بحث را ادامه دهیم، لازم است که متغیرهای به کار رفته در شبه کدها را معرفی کنیم. این متغیرها و شرح آنها به صورت زیر می باشد:

D: مجموعه آموزشی

S: مجموعه ویژگی اصلی (شامل تمام ویژگی ها)

N: تعداد ویژگی ها

T: زیرمجموعه ویژگی انتخاب شده

M: تعداد ویژگی های انتخاب شده یا تعداد ویژگی هایی که لازم است انتخاب شوند.

6-3-3-1: تابع ارزیابی مبتنی بر فاصله - تابع تولیدکننده مکاشفه ای

مهم ترین روش در این گروه Relief است. در اینجا ما ابتدا این روش را به عنوان نماینده این گروه شرح می دهیم، سپس یک مرور مختصر بر سایر روش ها خواهیم داشت.

روش Relief از یک راه حل آماری برای انتخاب ویژگی استفاده می کند، همچنین یک روش مبتنی بر وزن است که از الگوریتم های مبتنی بر نمونه الهام گرفته است. روش کار به این صورت است که از میان مجموعه نمونه های آموزشی، یک زیرمجموعه نمونه انتخاب می کنیم. کاربر بایستی تعداد نمونه ها در این زیرمجموعه را مشخص کرده باشد و آن را به عنوان ورودی به الگوریتم ارائه دهد. الگوریتم به صورت تصادفی یک نمونه از این زیرمجموعه را انتخاب می کند، سپس برای هر یک از ویژگی های این نمونه، نزدیک ترین برخورد¹⁴ و نزدیک ترین شکست¹⁵ را بر اساس معیار اقلیدسی پیدا می کند.

¹⁴ Nearest Hit

¹⁵ Nearest Miss

نزدیکترین برخورد نمونه‌ای است که کمترین فاصله اقلیدسی را در میان سایر نمونه‌های هم-کلاس با نمونه انتخاب شده دارد. نزدیکترین شکست نیز نمونه‌ای است که کمترین فاصله اقلیدسی را در میان نمونه‌هایی که هم‌کلاس با نمونه انتخاب شده نیستند، دارد.

ایده اصلی در این الگوریتم این است که هر چه اختلاف بین اندازه یک ویژگی در نمونه انتخاب شده و نزدیک‌ترین برخورد کمتر باشد، این ویژگی بهتر است و به‌علاوه یک ویژگی خوب آن است که اختلاف بین اندازه آن ویژگی و نزدیکترین شکست وی بیشتر باشد. دلیل کار خیلی ساده است، ویژگی‌هایی که به خوبی دو کلاس (یا یک کلاس از سایر کلاس‌ها) را از هم تمیز می‌دهند، برای نمونه‌های متعلق به دو کلاس متفاوت، مقادیری نزدیک به هم نمی‌دهند بلکه یک فاصله معنی‌دار، بین مقادیری که به نمونه‌های یک کلاس می‌دهند و مقادیری که به سایر کلاس‌ها می‌دهند وجود دارد.

الگوریتم پس از تعیین نزدیکترین برخورد و نزدیکترین شکست، وزن‌های ویژگی‌ها را به روزرسانی می‌کند، این به‌روزرسانی به این صورت است که مربع اختلاف بین مقدار ویژگی مورد نظر در نمونه انتخاب شده و نمونه نزدیکترین برخورد از وزن ویژگی کم می‌شود و در عوض مربع اختلاف بین مقدار ویژگی در نمونه انتخاب شده و نزدیکترین شکست به وزن ویژگی اضافه می‌شود. هر چه مقدار این وزن بزرگتر باشد، ویژگی مورد نظر، بهتر می‌تواند نمونه‌های یک کلاس را از دیگران جدا کند.

بعد از تعیین فاصله برای تمام نمونه‌های موجود در مجموعه نمونه‌ها، الگوریتم ویژگی‌هایی را که وزن آنها کمتر یا مساوی با یک حد آستانه است، حذف می‌کند و سایر ویژگی‌ها را به‌عنوان زیرمجموعه ویژگی جواب باز می‌گرداند. مقدار حد آستانه توسط کاربر تعیین می‌گردد. البته ممکن است که بصورت اتوماتیک بوسیله یک تابعی از تعداد کل ویژگی‌ها تعیین شود و یا این که با سعی و خطا تعیین گردد. همچنین می‌توان ویژگی‌هایی که وزن آنها منفی است را حذف کرد.

```

Relief( $D, S, NoSample, Threshold$ )
(1)  $T = \phi$ 
(2) Initialize all weights,  $W_i$ , to zero.
(3) For  $i = 1$  to  $NoSample/*$  Arbitrarily chosen  $*/$ 
    Randomly choose an instance  $x$  in  $D$ 
    Finds its  $nearHit$  and  $nearMiss$ 
    For  $j = 1$  to  $N$ 
         $W_j = W_j - diff(x_j, nearHit_j)^2 + diff(x_j, nearMiss_j)^2$ 
(4) For  $j = 1$  to  $N$ 
    If  $W_j \geq Threshold$ 
        Append feature  $f_j$  to  $T$ 
(5) Return  $T$ 

```

شکل 2-6 الگوریتم Relief

الگوریتم Relief برای ویژگی‌های دارای نویز یا ویژگی‌های دارای همبستگی، خوب کار می‌کند و پیچیدگی زمانی آن به صورت خطی و تابعی از تعداد ویژگی‌ها و تعداد نمونه‌های مجموعه نمونه می‌باشد و هم برای داده‌های پیوسته و هم برای داده‌های صوری خوب کار می‌کند.

یکی از محدودیت‌های اساسی این الگوریتم این است که ویژگی‌هایی که دارای افزونگی باشند را پیدا نمی‌کند و بنابراین مجموعه‌های غیر بهینه را پیدا می‌کند که دارای افزونگی هستند. این مشکل را می‌توان با یک جستجوی تعیین جامعیت برای زیرمجموعه‌های تولید شده توسط الگوریتم حل کرد. علاوه بر این مشکل دیگر این الگوریتم این است که با مسائل دو کلاسه خوب کار می‌کند. این محدودیت نیز با الگوریتم Relief-F مرتفع شده است، با الگوریتم جدید مشکل داده‌های غیر کامل (نمونه‌های آموزشی غیرکامل) نیز حل شده است.

روش دیگری که در این زمینه برای انتخاب ویژگی مطرح است، از یک تابع ارزیابی استفاده می‌کند که مجموع یک معیار اختلاف آماری و یک معیار پیچیدگی ویژگی را محاسبه کرده و آن را مینیمم می‌کند. این الگوریتم، اولین ویژگی را که بهتر بتواند کلاس‌ها را از هم تمیز دهد پیدا می‌کند. سپس ویژگی‌هایی را پیدا می‌کند که در ترکیب با ویژگی‌های انتخاب شده، جدائی-

پذیری کلاس‌ها را افزایش دهند. این فرآیند زمانی متوقف می‌شود که به حداقل معیار بازنمایی مورد انتظار برسیم.

6-3-3-2: تابع ارزیابی مبتنی بر فاصله - تابع تولیدکننده کامل

استفاده از این ترکیب در روش‌های قدیمی نظیر B&B¹⁶ یافت می‌شود. سایر روش‌های این گروه، نسخه‌های متفاوتی از B&B هستند. به این ترتیب که از یک تابع تولیدکننده دیگر استفاده کرده‌اند (BFF) و یا این که از یک تابع ارزیابی متفاوت استفاده کرده‌اند (Bobrowski's method). در اینجا ابتدا به شرح B&B می‌پردازیم و سپس یک شرح مختصر در مورد دو روش دیگر ارائه می‌دهیم.

تعریف کلاسیک ارائه شده در این زمینه از انتخاب ویژگی، احتیاج دارد که تابع ارزیابی یکنوا باشد. یعنی اگر دو زیرمجموعه ویژگی A و B با اندازه‌های M و N موجود و $B \subset A$ باشد در این صورت مقدار تابع ارزیابی برای A نباید بیشتر از مقدار تابع برای B باشد. این تعریف باعث ایجاد مشکل در مسائل دنیای واقعی می‌شود، زیرا اندازه تخمینی زیرمجموعه ویژگی بهینه در حالت عمومی ناشناخته است.

البته به سادگی می‌توان این تعریف را تغییر داد تا با مسائل عمومی سازگار شود، به این صورت که می‌گوئیم: الگوریتم‌های مشابه B&B تلاش می‌کنند که دو شرط زیر را هم‌زمان ارضاء کنند:

1. زیرمجموعه ویژگی جواب تا حد امکان کوچک باشد.
2. یک کران برای مقدار تابع ارزیابی را در نظر بگیرد. (یا یک اندازه مینیمم برای تعداد ویژگی‌های انتخاب شده در نظر بگیرد. مانند بهترین زیرمجموعه ویژگی سه عنصری)

¹⁶ (Branch and Bound)

به وسیله کران تعیین شده، فضای جستجو تا حد امکان کوچک می‌شود. به این ترتیب الگوریتم B&B از یک زیرمجموعه شامل تمام ویژگی‌های موجود شروع می‌کند و درخت جستجو را تشکیل می‌دهد. در این درخت در ریشه تمام ویژگی‌ها قرار دارند و فرزندان وی، زیرمجموعه‌هایی هستند که زیرمجموعه گره پدر هستند و از حذف تنها یکی از عناصر پدرشان تشکیل شده‌اند. این روند برای سایر گره‌های درخت تکرار می‌شود تا به مجموعه‌های تک‌عنصری (با تعداد ویژگی‌های تعیین شده به عنوان کران) برسیم. یعنی برگ‌های درخت مجموعه‌های تک‌عنصری هستند و ریشه درخت یک مجموعه شامل همه ویژگی‌های موجود می‌باشد.

با توجه به این خاصیت که تمام زیرمجموعه‌های یک مجموعه مقدار کمتری برای تابع ارزیابی دارند، در حین جستجو اگر یک گره به واسطه کم بودن مقدار تابع ارزیابی انتخاب نشد، زیرشاخه‌های آن را برای یافتن جواب جستجو نمی‌کنیم، چون قطعاً تابع ارزیابی مقدار کمتری را برای آن‌ها باز می‌گرداند. عموماً توابع ارزیابی زیر برای این کار استفاده می‌شوند:

• فاصله ماهالانوبیس¹⁷

• تابع جداساز¹⁸

• معیار فیشر¹⁹

• فاصله باتاچاریا²⁰

• Divergence

یک الگوریتم مشابه برای انتخاب ویژگی، BFF است، در این الگوریتم، تابع جستجو مشابه حل مساله جستجوی یک مسیر بهینه در یک درخت وزن دار، و با یک استراتژی تغییر یافته از Best

¹⁷ Mahalanobis Distance

¹⁸ Discriminant Function

¹⁹ (Fisher Criterion

²⁰ Bhattacharya

first search عمل می‌کند. این الگوریتم تضمین می‌کند که بهترین هدف، (زیرمجموعه بهینه) بدون از دست دادن جامعیت و با ارضای معیار یکنوا بودن تابع ارزیابی مساله، پیدا شود.

```

B&B( $D, S, M$ )
  if ( $card(S) <> M$ ) {
    /* subset generation */
     $j = 0$ 
    For all features  $f$  in  $S$  {
       $S_j = S - f$  /* remove one feature at a time */
      if ( $S_j$  is legitimate) /*
        if  $IsBetter(S_j, T)$ 
           $T = S_j$ 
        /* recursion */
        B&B( $S_j, M$ )
      }
    }
    return  $T$ 
  }

```

شکل 3-6 الگوریتم Branch and Bound

3-3-3-6: تابع ارزیابی مبتنی بر اطلاعات - تابع تولید کننده مکاشفه ای

در این دسته دو روش وجود دارد:

1) روش درخت تصمیم

در روش درخت تصمیم، نمونه‌ها به یک الگوریتم C4.5، که یکی از درخت‌های تصمیم‌گیری است اعمال می‌شوند، سپس درخت هرس شده حاصل از الگوریتم C4.5 را گرفته و کلیه ویژگی‌هایی که در آن وجود دارد را به‌عنوان جواب مساله باز می‌گرداند.

الگوریتم C4.5، از یک تابع مکاشفه بر پایه اطلاعات استفاده می‌کند، یک فرم ساده این توابع برای مسائل دو کلاسه به صورت زیر است:

$$(,) = \frac{-}{+} \log \frac{-}{+} - \frac{-}{+} \log \frac{-}{+}$$

که در آن p تعداد نمونه‌های کلاس اول و n تعداد نمونه‌های کلاس دوم است. فرض کنید که صفت F_1 به‌عنوان ریشه درخت در نظر گرفته شده است و مجموعه آموزشی را به دو زیرمجموعه T_0 و T_1 تقسیم کرده است. آنتروپی ویژگی F_1 برابر است با:

$$() = \frac{+}{+} (,) + \frac{+}{+} (,)$$

الگوریتم درخت تصمیم به صورت زیر است :

- (1) $() = \emptyset$
- (2) 4.5 h , h
- (3) h
- (4)

شکل 4-6 الگوریتم درخت تصمیم

2) روش استفاده شده توسط Sahami و Koller

روش استفاده شده توسط Sahami و Koller که اخیراً ارائه شده است، بر این پایه استوار است که ویژگی‌هایی که داده مفید چندانی را در بر ندارند و یا اصلاً داده مفیدی را در اختیار قرار نمی‌دهند و می‌توان آن‌ها را با سایر ویژگی‌ها نمایش داد، یا ویژگی‌هایی که بی‌ربط هستند و یا داده اضافی هستند، را شناسایی و حذف می‌کنیم. برای

پیاده سازی این مطلب، تلاش می‌کنیم تا با پوشش مارکوف آنها را پیدا کنیم، به این صورت که یک زیرمجموعه مانند T ، یک پوشش مارکوف برای ویژگی f_i است، اگر f_i برای زیرمجموعه T به صورت مشروط هم از کلاس و هم از سایر ویژگی‌هایی که در T نیستند، مستقل باشد.

6-3-3-4: تابع ارزیابی مبتنی بر اطلاعات - تابع تولید کننده کامل

مهمترین روشی که در این گروه می‌توانیم پیدا کنیم، روش MDLM²¹ است. محققان در این روش تلاش می‌کنند تا همه ویژگی‌های بدون استفاده (بی‌ربط یا اضافی) را حذف نمایند، در این صورت اگر ویژگی‌های زیرمجموعه V را بتوانیم به صورت یک تابع ثابت مانند F که وابسته به کلاس نیست و بر اساس یک زیرمجموعه ویژگی دیگر مانند U ، بیان کنیم آن‌گاه وقتی که مقادیر ویژگی‌های زیرمجموعه U شناخته شده باشند، ویژگی‌های موجود در زیرمجموعه V بدون استفاده هستند.

از دیدگاه انتخاب ویژگی، اجتماع دو زیرمجموعه U و V ، یک مجموعه کامل شامل تمام ویژگی‌ها را تشکیل می‌دهد. و کاری که ما باید در انتخاب ویژگی انجام دهیم این است که این دو زیرمجموعه را جدا کنیم. برای انجام این کار، نویسندگان MDLM، از معیار MDLC²² استفاده کرده‌اند. این روش شامل تعداد بیت‌های لازم برای انتقال کلاس‌ها، پارامترهای بهینه‌سازی، ویژگی‌های مفید و ویژگی‌های غیرمفید است. الگوریتم تمام زیرمجموعه‌های ممکن (2^N) را جستجو می‌کند و به عنوان خروجی زیرمجموعه‌ای را باز می‌گرداند که معیار MDLC را ارضا کند. این روش می‌تواند تمام ویژگی‌های مفیدی را پیدا کند که دارای توزیع نرمال باشند.

²¹ Minimum Description Length Method

²² Minimum Description Length Criterion

برای حالت‌های غیر نرمال این روش قادر نیست، ویژگی‌های مفید را پیدا کند. الگوریتم زیر روش کار و فرمول‌های استفاده شده را نشان می‌دهد.

MDLM(D)

(1) $MDL = \infty$

(2) For all feature subsets L

1.1 Compute $Length_L = \sum_{i=1}^{i=q} \frac{p_i}{2} \log \frac{|D_L(i)|}{|D_L|} + h_L$

where $h_L = \frac{1}{2}(N - M)(N + M + 3) \log P + \sum_{i=1}^{i=q} M(M + 3) \log P_i$,

N – total number of features,

M – number of features in the candidate subset,

P – total number of instances in D ,

P_i – number of instances with class label i ,

q – total number of class labels,

D_L – covariance matrix formed from all the useful feature vectors,

$D_L(i)$ – covariance matrix formed from the useful feature vectors,

of class i ,

$|\cdot|$ – denotes determinant.

1.2 If $Length_L < MDL$ then

$T = L, MDL = Length_L$

(3) Return T

شکل 5-6 الگوریتم روش (MDLM)

5-3-3-6: تابع ارزیابی مبتنی بر وابستگی-تابع تولید کننده مکاشفه‌ای

در این روش اولین ویژگی به این صورت تعیین می‌شود که احتمال خطا را برای تمام ویژگی‌ها محاسبه می‌کنیم، ویژگی با کمترین احتمال خطا (P_e)، به عنوان اولین ویژگی انتخاب می‌شود. ویژگی بعدی، ویژگی است که مجموع وزن دار P_e و میانگین ضریب همبستگی (ACC) با ویژگی(های) انتخاب‌شده را مینیمم کند. سایر ویژگی‌ها به همین ترتیب انتخاب می‌شوند. میانگین ضریب همبستگی به این صورت است که میانگین ضریب همبستگی ویژگی کاندید با ویژگی‌های انتخاب شده در آن نقطه محاسبه می‌شوند.

POE+ACC(D, M, w_1, w_2)

(1) $T = \phi$

(2) Find the feature with minimum P_e and append it to T

(3) For $i = 1$ to $M - 1$

Find the next feature with minimum $w_1(P_e) + w_2(ACC)$

Append it to T

(4) Return T

شکل 6-6 الگوریتم (POE+ACC) Probability Of Error & Average Correlation Coefficient

این روش می‌تواند تمام ویژگی‌ها را بر اساس مجموع وزن‌دار درجه‌بندی کند. شرط خاتمه نیز در این روش تعداد ویژگی‌های مورد نیاز خواهد بود.

6-3-3-6: تابع ارزیابی مبتنی بر سازگاری - تابع تولیدکننده کامل

روش‌هایی که در این گروه قرار دارند، در سالهای اخیر ارائه شده‌اند. ما به صورت مختصر سه روش در این گروه را بررسی می‌کنیم ولی بحث اصلی ما بر روی روش اول است.

Focus (1)

این روش یک حداقل‌گرا است، به این معنی که سعی می‌کند که حداقل تعداد ویژگی ممکن را برای ارائه پیدا کند. این روش فرضیه‌های قابل تعریف را بررسی کرده و فرضیه‌ای را که بتواند سازگاری را با حداقل تعداد ویژگی ممکن برقرار کند به‌عنوان جواب بازمی‌گرداند.

در ساده‌ترین پیاده‌سازی برای این روش، برای یافتن یک ناسازگاری با یک زیرمجموعه از ویژگی‌های انتخاب شده، درخت جستجو را به صورت سطح به سطح (جستجو در پهنا)، پیمایش می‌کنیم. در این جریان که از مجموعه‌های کوچکتر شروع می‌شود، در صورتی که به یک ناسازگاری برسیم، مجموعه انتخاب شده رد می‌شود و جستجو با

مجموعه بعدی ادامه می‌یابد. به محض این‌که به یک مجموعه برسیم که ناسازگاری نداشته باشد، جستجو متوقف شده و مجموعه یادشده به عنوان جواب انتخاب می‌شود. می‌توان گفت که این روش به نويز حساس است و نمی‌تواند نويز را مدیریت کند، زیرا در صورتی که نويز وجود داشته باشد، هیچ زیرمجموعه‌ای را نمی‌توان پیدا کرد که ناسازگاری نداشته باشد و الگوریتم تمام ویژگی‌ها را به عنوان جواب بازمی‌گرداند. با یک تغییر کوچک می‌توانیم این مساله را حل کنیم، به این‌صورت که اجازه می‌دهیم یک میزان معینی از ناسازگاری در مجموعه انتخاب شده وجود داشته باشد.

$$\begin{aligned}
 (1) \quad & \left(\begin{matrix} \cdot \\ \cdot \end{matrix} \right) = \\
 (2) \quad & = 0 \\
 & h \\
 & h \qquad h \\
 & =
 \end{aligned}$$

شکل 6-7 الگوریتم روش Focus

Schlimmer (2)

این روش از یک شمارش سیستماتیک برای تابع تولیدکننده و یک معیار ناسازگاری نیز به عنوان تابع ارزیابی استفاده می‌کند. همچنین با استفاده از یک تابع مکاشفه‌ای سرعت جستجو را برای یافتن زیرمجموعه بهینه افزایش می‌دهد. این تابع مکاشفه‌ای یک معیار قابلیت اعتماد است، بر پایه این ایده که احتمال مشاهده یک ناسازگاری مشاهده شده، نسبتی از درصد مقادیری است که زیاد مشاهده شده‌اند.

MIFES1 (3)

این روش در انتخاب ویژگی شباهت زیادی به روش Focus دارد. در اینجا مجموعه نمونه‌ها را به شکل یک ماتریس ارائه می‌دهیم، هر عنصر نماینده یک ترکیب یکتا از یک نمونه منفی و یک نمونه مثبت است. یک ویژگی مانند f ، یک پوشش برای یک نمونه از ماتریس نامیده می‌شود، اگر برای نمونه‌های منفی و نمونه‌های مثبت، مقادیر عکسی وجود داشته باشد. این روش از یک پوشش با همه ویژگی‌ها شروع می‌کند، و تکرار می‌شود تا وقتی که هیچ کاهشی نتوانیم برای پوشش داشته باشیم. مشکل اساسی این روش این است که فقط برای مسائل دو کلاسه و ویژگی‌های منطقی قابل استفاده است.

7-3-3-6: تابع ارزیابی مبتنی بر سازگاری - تابع تولیدکننده تصادفی

نماینده این گروه جدیدتر، روش LVF است. این روش فضای جستجو را به صورت تصادفی با استفاده از یک الگوریتم Las Vegas جستجو می‌کند، که یک سری انتخاب‌های احتمالی انجام می‌دهد تا با استفاده از آنها و یک معیار سازگاری که با معیار استفاده شده در الگوریتم Focus متفاوت است، سریع‌تر به جواب بهینه برسیم.

این روش برای هر زیرمجموعه کاندید، تعداد ناسازگاری را محاسبه می‌کند و بر این ایده استوار است که کلاس محتمل‌تر آن است که در میان نمونه‌های این زیرمجموعه ویژگی، تعداد بیشتری متعلق به آن کلاس باشند. یک حد آستانه برای ناسازگاری در نظر گرفته می‌شود، که در ابتدا ثابت و به صورت پیش فرض صفر است و هر زیرمجموعه‌ای که مقدار ناسازگاری آن بیشتر باشد، رد می‌شود.

```

LVF( $D, S, MaxTries, \delta$ )
(1)  $T = S$ 
(2) For  $i = 1$  to  $MaxTries$  {
    randomly choose a subset of features,  $S_j$ 
    if  $card(S_j) \leq card(T)$ 
        if  $inConCal(S_j, D) \leq \delta$ 
            if  $card(S_j) < card(T)$ 
                 $T = S_j$ 
            output  $S_j$ 
        else
            append  $S_j$  to  $T$ 
            output  $S_j$  as 'yet another solution' }
(3) return  $T$ 

```

شکل 6-8 الگوریتم روش LVF

این روش می‌تواند هر زیرمجموعه بهینه را حتی برای داده‌های دارای نویز پیدا کند، به شرط آنکه سطح نویز درست را در ابتدا مشخص کنیم. یک مزیت این روش اینست که احتیاجی نیست که کاربر مدت زیادی را برای به‌دست آوردن یک زیرمجموعه بهینه منتظر بماند، زیرا الگوریتم هر زیرمجموعه‌ای که بهتر از بهترین جواب قبلی باشد (هم از لحاظ اندازه زیرمجموعه انتخاب‌شده و هم از لحاظ نرخ سازگاری)، را به عنوان جواب باز می‌گرداند.

این الگوریتم کارا است، زیرا تنها مجموعه‌هایی برای ناسازگاری تست می‌شوند که تعداد ویژگی‌های درون آن کمتر یا مساوی بهترین زیرمجموعه‌ای است که تا کنون پیدا شده است. پیاده‌سازی آن آسان است و پیدا شدن زیرمجموعه بهینه را اگر منابع موجود اجازه دهند، تضمین می‌کند. یکی از مشکلات این الگوریتم این است که برای پیدا کردن جواب بهینه زمان بیشتری نسبت به الگوریتم‌هایی که از توابع تولیدکننده مکاشفه‌ای استفاده می‌کنند، نیاز دارد، چون این الگوریتم از دانش مربوط به زیرمجموعه‌های قبلی استفاده نمی‌کند.

6-3-3-8: تابع ارزیابی مبتنی بر خطای طبقه‌بندی کننده-تابع تولید کننده مکاشفه ای

همانطور که قبلاً نیز اشاره کردیم، به مجموعه روش‌هایی که از تابع ارزیابی مبتنی بر نرخ خطای طبقه‌بندی کننده استفاده می‌کنند، (بدون توجه به نوع تابع تولیدکننده استفاده شده) روش‌های wrapper می‌گویند. در این گروه روش‌های مشهور زیر را می‌توانیم ببینیم:

23 SFS (1)

این روش، با یک مجموعه خالی شروع می‌کند، سپس در هر تکرار یک ویژگی با استفاده از تابع ارزیابی مورد استفاده، به مجموعه جواب اضافه می‌کند، این کار را تکرار می‌کند تا زمانی که تعداد ویژگی لازم انتخاب شود. مشکلی که این روش با آن روبروست، این است که ویژگی اضافه‌شده در صورتی که مناسب نباشد، از مجموعه جواب حذف نمی‌شود.

24 SBS (2)

این روش برعکس SFS با مجموعه‌ای شامل تمام ویژگی‌ها شروع می‌کند و در هر بار تکرار الگوریتم، ویژگی که به‌وسیله تابع ارزیابی انتخاب می‌شود، را از مجموعه مورد نظر حذف می‌کند. این کار را تا زمانی ادامه می‌دهد که تعداد ویژگی‌ها برابر یک تعداد معینی شود. مانند روش قبل مشکل این روش این است که ویژگی حذف شده را دیگر به مجموعه اضافه نمی‌کند، حتی اگر مناسب باشد.

روش‌های دیگری که در این گروه وجود دارند، نسخه‌های متفاوتی از دو روش قبلی یا ترکیب آنها هستند.

²³ Sequential Forward Selection

²⁴ Sequential Backward Selection

SBS-Slash (3)

این روش بر پایه این مشاهده است که هنگامی که تعداد زیادی ویژگی وجود دارد، بعضی از طبقه‌بندی‌کننده‌ها (مانند ID3 یا C4.5) مکرراً تعداد زیادی از آنها را استفاده نمی‌کنند. الگوریتم با یک مجموعه ویژگی شروع می‌کند (مانند SBS)، اما بعد از یک مرحله تمام ویژگی‌هایی را که در این مرحله یاد گرفته است و استفاده نشده‌اند، را حذف²⁵ می‌کند.

²⁶PQSS (4)

در اینجا از بعضی از خواص بازگشت به عقب استفاده می‌کنیم. عملکرد الگوریتم به این صورت است که در هر مرحله p ویژگی را به مجموعه اضافه و q ویژگی از آن را حذف می‌کند. حال اگر الگوریتم از مجموعه خالی شروع کرده باشد، بایستی اندازه p بزرگتر از اندازه q باشد. ولی اگر از مجموعه تمام ویژگی‌ها شروع شده باشد، بایستی اندازه p کوچکتر از q باشد.

²⁷ BDS (5)

مانند روش‌های قبل است با این تفاوت که جستجو را از دو طرف انجام می‌دهد.

Schemata Search (6)

الگوریتم با مجموعه خالی و یا مجموعه تمام ویژگی‌ها شروع می‌کند و در هر تکرار، بهترین زیرمجموعه را با حذف یا اضافه تنها یک ویژگی به مجموعه ویژگی، پیدا می‌کند.

²⁵ Slashes

²⁶ (p,q) Sequential Search

²⁷ Bi-Directional Search

کند. برای اینکه هر زیرمجموعه را ارزیابی کند، از تعیین اعتبار $LOOCV^{28}$ استفاده می‌کند. در هر تکرار زیرمجموعه‌ای انتخاب می‌شود که کمترین خطای $LOOCV$ را داشته باشد. کار به این صورت ادامه می‌یابد تا هیچ تغییر با تک ویژگی نتواند باعث بهتر شدن زیرمجموعه شود.

7 RC ²⁹

در اینجا این واقعیت تشریح شده است که بعضی از ویژگی‌ها فقط به قسمتی از فضای کار مربوط هستند. روش کار مشابه SBS است، اما دارای تغییرات عمده‌ای است که باعث محلی شدن آن شده است (به عبارت دیگر انتخاب ویژگی‌های مرتبط بر اساس تصمیم‌گیری بوسیله نمونه‌ها می‌باشد).

8 Queiros and Gelsema

این الگوریتم شبیه SFS است اما پیشنهاد می‌کند که در هر تکرار، هر ویژگی با تنظیمات متفاوتی بوسیله اثرات متفاوت ناشی از ویژگی‌های قبلی ارزیابی شود. دو نمونه از این تنظیمات به این صورت هستند:

(i) همیشه فرض کنیم که ویژگی‌ها مستقل هستند (ویژگی‌های قبلی را در نظر نمی‌گیریم).

(ii) هیچگاه فرض نمی‌کنیم که ویژگی‌ها مستقل هستند (ویژگی‌های قبلی را در نظر می‌گیریم).

در این روش و تعدادی از روش‌های قبلی در این گروه از نرخ خطای بیز به عنوان خطای طبقه‌بندی‌کننده استفاده می‌کنیم.

²⁸ Leave-One-Out Cross Validation (

²⁹ Relevance in Context

9-3-3-6: تابع ارزیابی مبتنی بر خطای طبقه بندی کننده-تابع تولیدکننده کامل

در این گروه چهار روش وجود دارد که دو روش اول آن بوسیله Ichino و Sklansky ارائه شده است.

Linear Classifier (1)

Box Classifier (2)

در دو روش فوق، مساله انتخاب ویژگی به وسیله برنامه نویسی صفر و یک حل شده است.

³⁰AMB&B (3)

این روش برای حل مشکلات B&B ارائه شده است، به این صورت که به تابع ارزیابی اجازه داده می شود که غیر یکنوا باشد. در اینجا به تابع تولیدکننده اجازه داده می شود که زیرمجموعه هایی تولید کند که محدودیت تعیین شده را نقض می کنند، اما زیرمجموعه ای که به عنوان جواب انتخاب می شود، نباید محدودیت ذکر شده را نقض کرده باشد.

³¹BS (4)

این روش یک نمونه از جستجوی Best-First است که از صف محدود شده برای محدود کردن فضای جستجو استفاده می کند. صف از بهترین به بدترین مرتب می شود، در این صورت، بهترین زیرمجموعه در ابتدای صف قرار داده می شود. تابع تولیدکننده به این صورت عمل می کند که زیرمجموعه موجود در ابتدای صف را انتخاب و کلیه زیرمجموعه های ممکن با اضافه کردن یک ویژگی به آن را تولید می کند و آنها را در

³⁰ Approximate monotonic branch and bound

³¹ Beam Search

محل مناسب خود در صف قرار می‌دهد. در صورتی که هیچ محدودیتی در اندازه صف نداشته باشیم، این روش یک جستجوی جامع است. در حالتی که محدودیت طول برابر یک را برای صف داشته باشیم، این روش با SFS برابر است.

10-3-3-6: تابع ارزیابی مبتنی بر خطای طبقه بندی کننده-تابع تولید کننده تصادفی

در این گروه پنج روش وجود دارد، که به شرح ذیل می‌باشند.

LVW (1)

این روش زیرمجموعه‌هایی به صورت کاملاً تصادفی با استفاده از الگوریتم Las Vegas تولید می‌کند.

³²GA (2)

در این روش یک جمعیت از زیرمجموعه‌های کاندید تولید می‌کنیم. در هر بار تکرار الگوریتم، با استفاده از عملگرهای جهش و بازترکیبی بر روی عناصر جمعیت قبلی، عناصر جدیدی تولید می‌کنیم. با استفاده از یک تابع ارزیابی، میزان شایستگی عناصر جمعیت فعلی را مشخص کرده و عناصر بهتر را به عنوان جمعیت نسل بعد انتخاب می‌کنیم. پیدا شدن بهترین جواب در این روش تضمین نمی‌شود ولی همیشه یک جواب خوب به نسبت مدت زمانی که به الگوریتم اجازه اجرا داده باشیم، پیدا می‌کند.

³³SA (3)

در اینجا نیز مانند الگوریتم ژنتیک، تابع تولیدکننده آن از تولید تصادفی استفاده می‌کند ولی در تولید تصادفی از یک جریان خاصی پیروی می‌نماید.

³² Genetic Algorithm

³³ Simulated Annealing

34RGSS (4)

این روش مشابه SBS و SFS است. با این تفاوت که یک زیرمجموعه تصادفی تولید می‌کند و سپس SBS و SFS را با استفاده از این زیرمجموعه تولید شده اجرا می‌کند. در واقع فاکتور تصادف را به دو روش ذکر شده تزریق می‌کند، تا کارایی آنها را افزایش دهد.

35 RMHC-PF1 (5)

نمونه‌های اولیه و ویژگی‌ها در اینجا به‌طور همزمان برای استفاده در طبقه‌بندی‌کننده نزدیک‌ترین همسایه، انتخاب می‌شوند، همچنین برای ثبت نمونه‌های اولیه و ویژگی‌ها از یک بردار شرطی استفاده می‌شود. تابع ارزیابی نیز، نرخ خطای طبقه‌بندی‌کننده نزدیک‌ترین همسایه می‌باشد. در هر تکرار، به‌صورت تصادفی یکی از بیت‌های بردار جهش داده می‌شوند تا یک بردار جدید برای تکرار بعدی تولید شود.

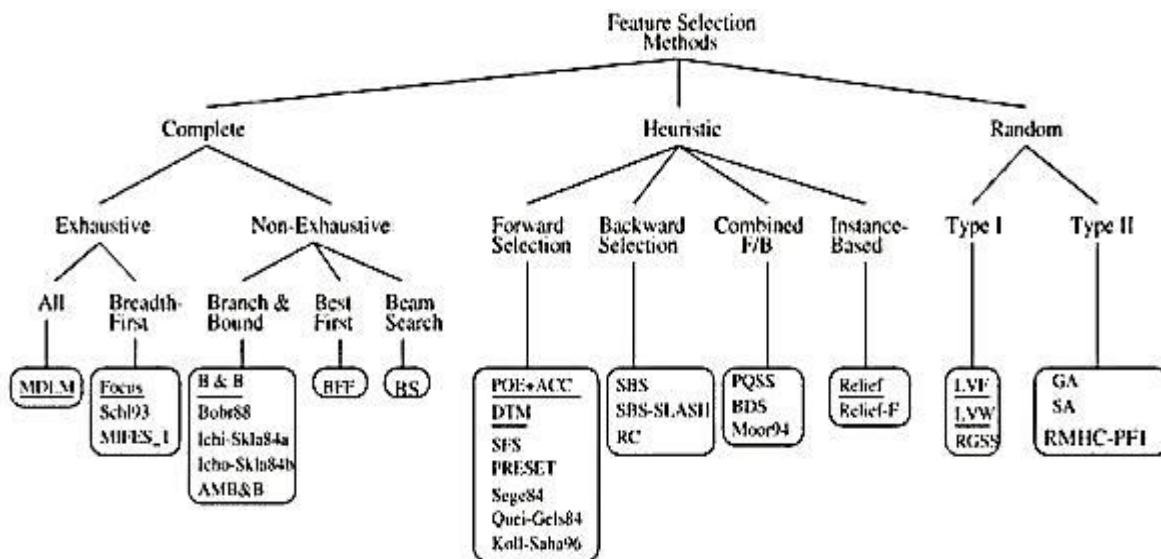
تمام روش‌های این گروه پارامترهای زیادی دارند که بایستی تنظیم شود، مثلاً در LVW حد آستانه برای نرخ ناسازگاری، در الگوریتم‌های ژنتیک اندازه جمعیت اولیه، نرخ بازترکیبی و نرخ جهش و یا در SA، تعداد تکرار حلقه، دمای اولیه و احتمال جهش بایستس تنظیم شوند. تنظیم دقیق این پارامترها عملکرد این الگوریتم‌ها را بهبود می‌بخشد.

³⁴ Random Generation plus Sequential Selection

³⁵ Random Mutation Hill Climbing-Prototype and Feature selection

4-3-6: جمع بندی روش های انتخاب ویژگی

برای اینکه یک جمع بندی از کلیه روش های انتخاب ویژگی داشته باشیم، نمودار آنها را برحسب سه نوع تابع تولیدکننده در شکل زیر نشان داده ایم.



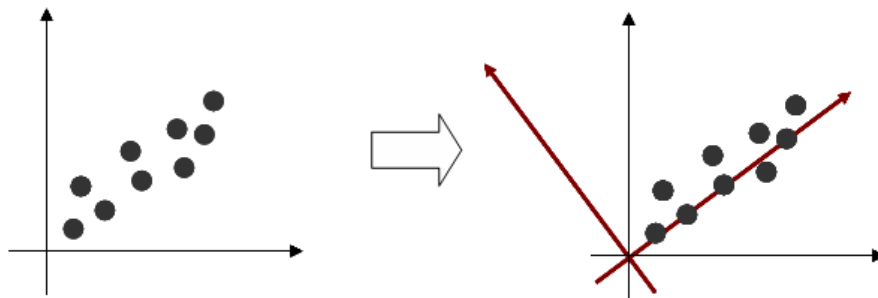
شکل 9-6 طبقه بندی روش های مختلف انتخاب ویژگی

4-6: PCA³⁶

تکنیک PCA بهترین روش برای کاهش ابعاد داده به صورت خطی می باشد. یعنی با حذف ضرایب کم اهمیت به دست آمده از این تبدیل، اطلاعات از دست رفته نسبت به روش های دیگر کمتر است. البته کاربرد PCA محدود به کاهش ابعاد داده، نمی شود و در زمینه های دیگری مانند شناسایی الگو و تشخیص چهره نیز مورد استفاده قرار می گیرد. در این روش محورهای

³⁶ Principal Component Analysis

مختصات جدیدی برای داده‌ها تعریف شده و داده‌ها براساس این محورهای مختصات جدید بیان می‌شوند. اولین محور باید در جهتی قرار گیرد که واریانس داده‌ها ماکزیمم شود (یعنی در جهتی که پراکندگی داده‌ها بیشتر است). دومین محور باید عمود بر محور اول و به همین ترتیب محورهای بعدی عمود بر تمامی محورهای قبلی به گونه‌ای قرار می‌گیرند که داده‌ها در آن جهت دارای بیشترین پراکندگی باشند. در شکل زیر این مطلب برای داده‌های دو بعدی نشان داده شده است.



شکل 6-10 انتخاب محورهای جدید برای داده‌های دو بعدی

روش PCA به نام‌های دیگری نیز معروف است. مانند:

➤ ³⁷SVD

➤ ³⁸KLT

➤ Hotelling Transform

➤ ³⁹EOF

³⁷ Singular Value Decomposition

³⁸ Karhunen Loeve Transform

³⁹ Empirical Orthogonal Function

قبل از این که به جزئیات این روش بپردازیم ابتدا مفاهیم ریاضی و آمار مرتبط با این روش را به طور مختصر بیان می کنیم. این مفاهیم شامل انحراف از معیار استاندارد، کواریانس، بردارهای ویژه و مقادیر ویژه می باشد.

6-4-1: مفاهیم مقدماتی مورد نیاز در PCA

6-4-1-1: مفاهیم آماری

فرض کنید X رشته ای از مقادیر است. میانگین این مقادیر از رابطه زیر به دست می آید.

$$= \frac{\sum}{n}$$

انحراف از معیار نیز از رابطه زیر محاسبه می شود.

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)}$$

علت اینکه در مخرج رابطه فوق از عبارت $n-1$ استفاده شده (و نه n) این است که فرض شده X شامل تمام مقادیر موجود نیست. بلکه تعدادی از این مقادیر انتخاب شده اند و در X قرار گرفته اند. یعنی X مجموعه نمونه است و شامل کل داده ها نمی باشد. با این فرض اگر از $n-1$ در رابطه فوق استفاده شود، انحراف از معیار به دست آمده به انحراف از معیار داده های واقعی نزدیک تر خواهد بود نسبت به این که از n استفاده شود. با بتوان 2 رساندن انحراف از معیار، واریانس به دست می آید.

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)}$$

معیارهایی که در بالا ذکر شد فقط اطلاعات مربوط به یک بُعد را ارائه می کنند و دانشی در مورد ارتباط بین ابعاد مختلف به ما نمی دهند. با استفاده از کواریانس می توانیم ارتباط بین ابعاد

مختلف داده‌ها را پیدا کنیم. فرض کنید یک رشته دیگر از اعداد داریم که آن را با Y نشان می‌دهیم. کواریانس بین X و Y از رابطه زیر بدست می‌آید.

$$(\sigma_{XY}) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{(n - 1)}$$

مقداری که از رابطه بالا بدست می‌آید در بازه $[-1, 1]$ قرار خواهد داشت که یکی از سه حالت زیر را به وجود می‌آورد:

- اگر مقدار به دست آمده مثبت باشد آنگاه X و Y با هم افزایش یا کاهش می‌یابند.
- اگر مقدار به دست آمده منفی باشد آنگاه با افزایش X مقدار Y کاهش می‌یابد و بالعکس.
- اگر مقدار به دست آمده صفر باشد آنگاه X و Y از یکدیگر مستقلند.

کواریانس بین تمامی ابعاد داده‌ها را می‌توان دوبه دو محاسبه کرده و در یک ماتریس ذخیره کرد. به این ماتریس، ماتریس کواریانس می‌گویند. ماتریس کواریانس یک ماتریس مربعی متقارن است. مثلاً اگر سه بعد به نامهای X ، Y و Z داشته باشیم، ماتریس کواریانس آنها برابر است با:

$$= \begin{pmatrix} (\sigma_{XX}) & (\sigma_{XY}) & (\sigma_{XZ}) \\ (\sigma_{YX}) & (\sigma_{YY}) & (\sigma_{YZ}) \\ (\sigma_{ZX}) & (\sigma_{ZY}) & (\sigma_{ZZ}) \end{pmatrix}$$

6-4-1-2: مفاهیم جبر ماتریسها

در این بخش مفهوم بردار ویژه و مقادیر ویژه را بیان می‌کنیم. همانطور که می‌دانید برای این که بتوان دو ماتریس را در یکدیگر ضرب کرد، آن دو باید از نظر اندازه با هم سازگار باشند. بردارهای ویژه نوع خاصی از ضرب ماتریسها را ارائه می‌کنند. به مثال زیر توجه کنید.

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \frac{1}{3} = \frac{11}{5}$$

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \frac{3}{2} = \frac{12}{8} = 4 \times \frac{3}{2}$$

در مثال اول بردار به دست آمده مضرب صحیحی از بردار اولیه نیست. اما در مثال دوم، بردار به دست آمده چهار برابر بردار اولیه می باشد. ماتریس 2×2 که در این دو بردار ضرب کرده ایم را می توان یک ماتریس تبدیل⁴⁰ در نظر گرفت که با ضرب آن در یک بردار می توان اندازه و راستای آن بردار را تغییر داد. در میان تمام بردارهایی که می توان ماتریس تبدیل را در آنها ضرب کرد، بردارهایی وجود دارند که پس از تبدیل، راستای آنها تغییر نمی کند و فقط اندازه آنها ممکن است عوض شود، مانند بردار $[3;2]$ در مثال فوق. این بردارها را بردارهای ویژه می نامند. برای هر بردار ویژه یک مقدار ویژه نیز وجود دارد که بیان می کند اندازه آن بردار (و تمام بردارهای دیگر که در راستای آن بردار هستند) پس از تبدیل، چند برابر خواهد شد. در مثال فوق مقدار ویژه برای بردار $[3;2]$ و البته تمام بردارهای هم راستا با آن مانند $[6;4]$ برابر با 4 می باشد.

بردارهای ویژه و مقادیر ویژه فقط برای ماتریس های مربعی معنی پیدا می کنند. یک ماتریس $n \times n$ می تواند دارای n بردار ویژه باشد. به منظور استاندارد کردن بردارهای ویژه، پس از یافتن بردارهای ویژه اندازه آنها را به گونه ای تغییر می دهند تا طول آنها برابر با یک شود. مثلاً برای بردار $[3;2]$ داریم:

$$\sqrt{(3 + 2)} = \sqrt{13} \quad \frac{3}{2} \div 13 = \frac{3/\sqrt{13}}{2/\sqrt{13}}$$

ویژگی مهم بردارهای ویژه این است که آنها برهم عمودند. مثلاً ماتریس تبدیل زیر را در نظر بگیرید.

⁴⁰ Transformation matrix

$$= \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$$

با ضرب ماتریس A در هر بردار می‌توان قرینه آن بردار نسبت به خط $y=0$ را به‌دست آورد. بردارهای ویژه این ماتریس عبارتند از $[1;0]$ و $[0;1]$. مقادیر ویژه این بردارها نیز به ترتیب 1 و 1 می‌باشند. همانطور که گفتیم این دو بردار ویژه برهم عمودند. به‌دست آوردن بردارهای ویژه برای ماتریس‌های بزرگتر از 3×3 کار نسبتاً دشواری است. این کار توسط یک الگوریتم بازگشتی انجام می‌شود.

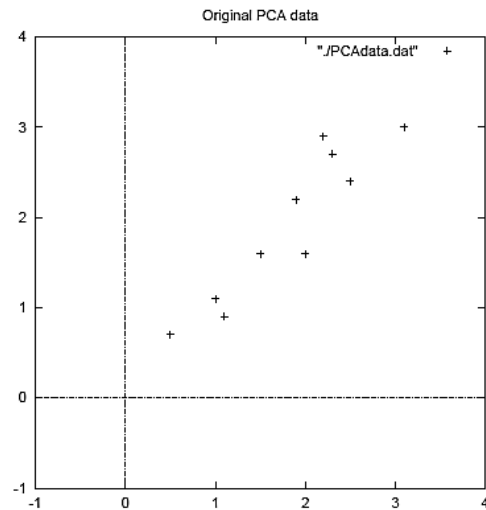
6-4-2: الگوریتم PCA

در این بخش الگوریتم PCA را با ذکر یک مثال توضیح می‌دهیم.

مرحله 1- انتخاب داده

در اینجا ما قصد داریم PCA را بر روی یک مجموعه داده دو بعدی اعمال کنیم. این داده‌ها در شکل زیر نشان داده شده است.

x	y
2.5	2.4
0.5	0.7
2.2	2.9
1.9	2.2
3.1	3.0
2.3	2.7
2	1.6
1	1.1
1.5	1.6
1.1	0.9



شکل 6-8 داده‌های دوبعدی اولیه که قرار است PCA بر روی آنها اعمال شود

مرحله 2- کم کردن میانگین از داده‌ها

در این مرحله، میانگین هر بُعد را از مقادیر آن بُعد کم می‌کنیم تا میانگین داده‌ها در هر بُعد صفر شود.

مرحله 3- محاسبه ماتریس کواریانس

ماتریس کواریانس را به طریقی که در بالا ذکر شد برای داده‌ها به‌دست می‌آوریم. برای مثال ما این ماتریس، یک ماتریس 2×2 است:

$$= \begin{pmatrix} .616555556 & .615444444 \\ .615444444 & .716555556 \end{pmatrix}$$

مرحله 4- محاسبه بردارهای ویژه و مقادیر ویژه

اکنون بردارهای ویژه و مقادیر ویژه ماتریس کواریانس را محاسبه می‌کنیم. ماتریس کواریانس، یک ماتریس نیمه قطعی مثبت متقارن⁴¹ است. طبق قضایای جبر خطی، یک ماتریس متقارن $n \times n$ دارای n بردار ویژه مستقل و n مقدار ویژه می‌باشد. همچنین یک ماتریس نیمه قطعی مثبت، دارای مقادیر ویژه غیر منفی است. این مقادیر برای مثال ما برابر است با:

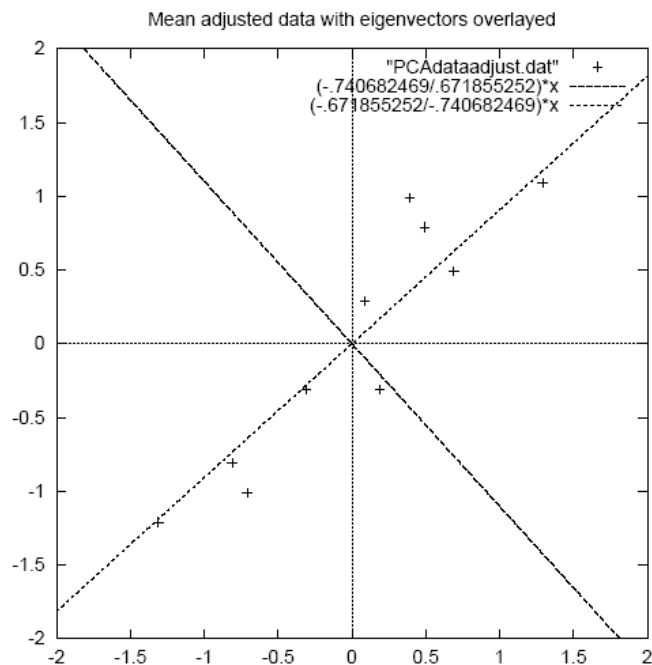
$$= \begin{matrix} -0.490833989 \\ 1.28402771 \end{matrix}$$

$$= \begin{matrix} -0.735178656 & -0.677873399 \\ 0.677873399 & -0.735178656 \end{matrix}$$

توجه داشته باشید که این دو بردار ویژه به گونه‌ای انتخاب شده‌اند که طول هر دو برابر با 1 باشد. اما این دو بردار چه چیزی به ما می‌دهد؟ ما راستای این دو بردار را در شکل 9-6 نشان داده‌ایم. همانطور که می‌بینید یکی از این دو بردار در جهتی قرار گرفته است که داده‌ها در آن جهت بیشترین پراکندگی را دارند. بردار دیگر نیز عمود بر بردار اول است. اما مقادیر ویژه چه چیزی را ارائه می‌دهند؟ در این مثال برداری که در راستای بیشترین پراکندگی داده‌ها قرار گرفته است دارای مقدار ویژه 1.284 و بردار دیگر دارای مقدار ویژه 0.049 می‌باشد. در واقع مقادیر ویژه، میزان پراکندگی داده‌ها در راستای بردار ویژه مربوطه را نشان می‌دهد. می‌توان گفت بردار ویژه‌ای که دارای بزرگترین مقدار ویژه است مؤلفه اصلی⁴² داده‌های موجود می‌باشد.

⁴¹ Symmetric positive semidefinite matrix

⁴² Principle component



شکل 6-11 داده‌های نرمال‌سازی شده (با کم شدن میانگین) به همراه بردارهای ویژگی ماتریس کواریانس

مرحله 5- انتخاب مؤلفه‌ها و ساختن بردار ویژگی

در این مرحله مفهوم کاهش ابعاد داده وارد می‌شود. بردارهای ویژه‌ای که در مرحله قبل به دست آوردیم را بر اساس مقادیر ویژه آنها از بزرگ به کوچک مرتب می‌کنیم (توجه داشته باشید که مقادیر ویژه ماتریس کواریانس همگی بزرگتر یا مساوی صفر هستند). بدین ترتیب مؤلفه‌های داده‌ها از پر اهمیت به کم اهمیت مرتب می‌شوند. در اینجا اگر بخواهیم ابعاد داده‌ها را کاهش دهیم می‌توانیم مؤلفه‌های کم اهمیت را حذف نماییم. البته این کار با از دست دادن مقدار کمی اطلاعات همراه است.

کاری که باید در این مرحله انجام دهیم ایجاد یک بردار ویژگی است که در واقع ماتریسی از بردارها می‌باشد. این ماتریس شامل بردارهای ویژگی است که ما می‌خواهیم آن‌ها را نگه داریم.

$$= (\quad , \quad , \quad , \dots , \quad)$$

اگر همه بردارهای ویژگی را در این ماتریس قرار دهیم، هیچ اطلاعاتی از دست نخواهد رفت و دوباره می‌توانیم دقیقاً همان داده‌های اولیه را به‌دست آوریم. در ادامه مثال فوق، یک بردار ویژگی را برابر با مقدار زیر در نظر می‌گیریم.

$$= \begin{matrix} -0.677873399 \\ -0.735178656 \end{matrix}$$

مرحله 6- بدست آوردن داده‌های جدید

در آخرین مرحله از PCA فقط باید ترانزاده ماتریس بردار ویژگی که در مرحله قبل به‌دست آوردیم را در ترانزاده داده‌های نرمال‌سازی‌شده ضرب نماییم.

$$= \quad \times$$

که RowFeatureVector ماتریسی است که بردارهای ویژه در سطرها و آن به ترتیب مقادیر ویژه از بالا به پایین قرار گرفته‌اند و RowDataAdjust ماتریسی است که شامل داده‌هایی است که میانگین هر بُعد از آن بُعد کم شده است. در این ماتریس، داده‌ها در ستون‌های آن ذخیره شده و هر سطر آن مربوط به یک بُعد است. در مثال ذکر شده به‌دلیل اینکه ما فقط یکی از بردارهای ویژه را انتخاب کردیم داده‌های به‌دست آمده از PCA، داده‌های یک بُعدی می‌باشند.

$$\begin{array}{c} x \\ \hline -.827970186 \\ 1.77758033 \\ -.992197494 \\ -.274210416 \\ -1.67580142 \\ -.912949103 \\ .0991094375 \\ 1.14457216 \\ .438046137 \\ 1.22382056 \end{array}$$

شکل 6-12 داده‌های بدست آمده از تبدیل PCA با انتخاب مهم‌ترین بردار ویژگی

با استفاده از رابطه زیر می‌توانیم مقادیری که از تبدیل PCA بدست آورده‌ایم را به داده‌های اولیه که مقدار میانگین از آنها کم شده بازگردانیم.

$$= \times$$

به دلیل اینکه ماتریس حاوی بردارهای ویژه یکه است، معکوس آن با ترانپوز آن برابر است. بنابراین:

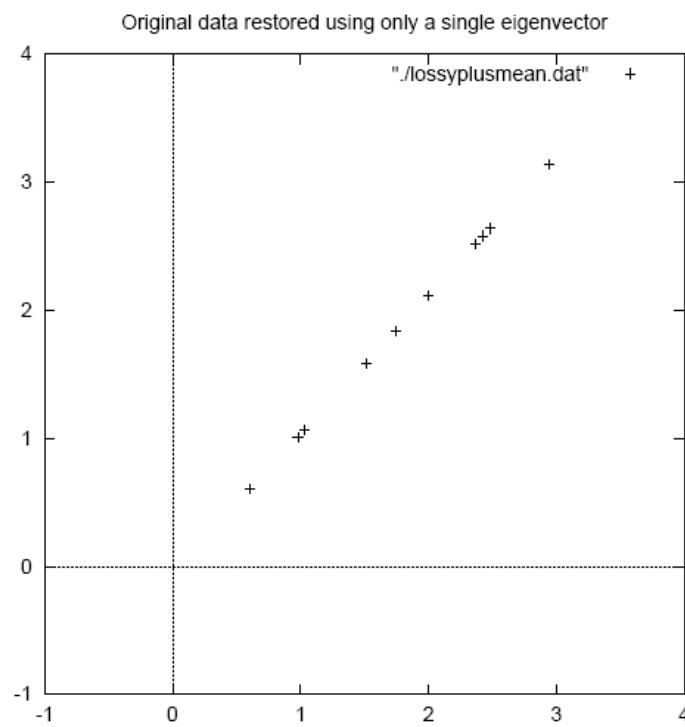
$$= \times$$

با اضافه کردن میانگین، داده‌های اولیه را خواهیم داشت:

$$=$$

$$(\times) +$$

در شکل زیر داده‌هایی که پس از تبدیل PCA بازیابی شده‌اند را مشاهده می‌کنید. همانطور که می‌بینید مقدار کمی اطلاعات از بین رفته است و باعث شده داده‌ها همگی در امتداد یک خط راست قرار گیرند.



شکل 6-13 داده‌های بازیابی شده از تبدیل PCA با انتخاب مهم‌ترین بردار ویژگی

3-4-6: آنالیز تفکیک پذیری خطی (LDA)

بدون تردید یکی از پرطرفدارترین و باسابقه‌ترین گرایش‌های تحقیقاتی در چند سال اخیر استفاده از نسخه‌های مختلف بسط آنالیز تفکیک‌پذیری خطی یا به اختصار LDA است. در بسط LDA هدف تولید فضایی است که در آن پراکندگی درون کلاسی مینیمم و پراکندگی برون

کلاسی ماکزیمم باشد. اما مقصود از پراکندگی درون کلاسی و برون کلاسی چیست؟ بطور شهودی چنانچه داده‌های هر کلاس متمرکز و توده‌های کلاسی تا حد امکان از یکدیگر دور باشند طراحی یک ماشین یادگیری برای طبقه‌بندی داده‌ها با سهولت بیشتری همراه بوده و ماشین یادگیری از پتانسیل بالاتری برای طبقه‌بندی داده‌ها برخوردار است. پراکندگی درون کلاسی (رابطه 4-1) شاخصی برای میزان تمرکز داده‌ها در هر کلاس و پراکندگی برون کلاسی (رابطه 5-1) شاخصی برای میزان دور بودن توده‌های کلاسی از یکدیگر است.

$$= -\sum_{i \in c} (x_i - \mu_c)(x_i - \mu_c) \quad (4-1)$$

$$= -\sum_i (x_i - \mu_c)(x_i - \mu_c) \quad (5-1)$$

در رابطه 4-1 و رابطه 5-1 فرض کرده‌ایم که با یک مسئله شناسایی الگو با N داده و c کلاس مواجه هستیم. در این روابط m_i و m به ترتیب میانگین تمام داده‌ها و میانگین داده‌های درون کلاس i هستند. N_i تعداد داده‌های درون کلاس i را نشان می‌دهد. در بسط LDA برای آن‌که دو شاخص پراکندگی درون و برون کلاسی به طور همزمان و به ترتیب مینیمم و ماکزیمم شوند باید رابطه (6-1) که شاخص فیشر نامیده می‌شود ماکزیمم شود. می‌توان نشان

داد که نگاشت W از کنار هم چیدن بردارهای ویژه ماتریس $S_w^{-1} S_b$ بدست می‌آید.

$$= \arg \max | \text{————} | \quad (6-1)$$

نسخه‌های جدیدتر الگوریتم نیز با تغییر در معیار فیشر حاصل شده‌اند. به عنوان مثال در الگوریتم F-LDA ماتریس پراکندگی برون کلاسی را با ترکیبی وزن‌دار محاسبه می‌کنند. استدلال استفاده از ترکیب وزن‌دار، جریمه کردن کلاس‌هایی است که تفکیک‌پذیری کمتری داشته و مراکز آن‌ها به یکدیگر نزدیک‌تر است. با دادن وزن در محاسبه ماتریس پراکندگی برون کلاسی، نقش این کلاس‌ها در تولید فضای نهایی کمتر می‌شود.

نسخه معروف و جدید دیگری از الگوریتم LDA با نام Regularized LDA توسط لو و همکاران در سال 2005 میلادی با اصلاح معیار فیشر مطابق با رابطه (7-1) معرفی شد.

$$= \frac{| \quad | \quad |}{| \quad | \quad |} \quad (7-1)$$

اگرچه باور عمومی بر این است که کارایی بسط LDA همواره از بسط PCA بهتر است اما مشاهدات مختلفی گزارش شدند که عکس این باور را در برخی سناریوها نشان می دادند. بعدها مارتینز و کاک در یک قالب جامع ریاضی نشان دادند که همواره نمی توان انتظار داشت بسط LDA بهتر از بسط PCA عمل نماید.