

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/276255530>

مروری بر ماشین های بردار پشتیبان

Research · May 2015

DOI: 10.13140/RG.2.1.4981.6487

CITATIONS

0

READS

2,446

1 author:



Mazdak Fatahi

Razi University

9 PUBLICATIONS 3 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Neuromorphic Hardware [View project](#)

All content following this page was uploaded by Mazdak Fatahi on 13 May 2015.

The user has requested enhancement of the downloaded file.



دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

عنوان :

مروری بر ماشین های بردار پشتیبان

Support Vector Machines: A survey

استاد مربوطه:

دکتر جهان شاه کبودیان

ارائه دهنده:

مزدک فتاحی

چکیده

استفاده از ماشین های بردار پشتیبان (SVM) که توسط Vapnik ارائه شد، بعنوان یکی از راه حل ها در یادگیری ماشین و تشخیص الگو گسترش یافته است. SVM پیشگویی های خود را با استفاده از ترکیبی خطی از تابع Kernel که بر روی مجموعه ای از داده های آموزشی با نام بردارهای پشتیبان عمل می کند، انجام می دهد. روش ارائه شده توسط SVM با روش های قابل مقایسه مانند Neural Network متفاوت می باشد، تعلیم SVM همواره می نیمم سراسری را پیدا می کند. ویژگی های یک SVM به مقدار زیادی به انتخاب kernel آن مربوط می شود. تعلیم SVM منجر به یک مسئله برنامه ریزی درجه ی دوم (QP) از نوع مقید می گردد که حل آن برای حجم زیادی از نمونه ها می تواند با روش های عددی بسیار مشکل باشد، لذا برای ساده کردن حل این مسئله بهینه سازی، روش های متعددی ارائه گردیده که متناسب با نیاز قابل استفاده و پیاده سازی می باشند. در این مقاله ابتدا راجع مسئله یادگیری ماشین، طبقه بندی و معرفی SVM بحث شده و سپس معادله ی مربوط به یادگیری در SVM بدست می آید. در پایان نیز تعدادی از روش های حل این معادله ذکر می گردند. در این میان به کاربردها و بسته های کاربردی برای کار با SVM نیز اشاره گردیده است.

کلمات کلیدی: SVM، بردار پشتیبان، بهینه سازی مقید، طبقه بندی، Classification

فهرست مطالب

۱- مقدمه	۱
۲- مروری بر دسته‌بندی کننده‌ی SVM	۴
۲-۱- حالت جدایی پذیر	۶
۳- SVM در حالت‌های غیرخطی و جدایی ناپذیر	۱۳
۳-۱- داده‌های جدایی ناپذیر	۱۳
۳-۲- ماشین‌های بردار پشتیبان غیر خطی	۱۵
۳-۲- ماشین‌های بردار پشتیبان به عنوان جداکننده‌ی چند کلاسه	۱۷
۴- راه‌حل‌های بهینه جهت حل معادله‌ی QP برای SVM	۱۸
۴-۱- Decomposition Method:	۲۰
۴-۲- SVM ^{light} :	۲۴
۴-۳- Chunking	۲۵
۴-۴- SMO(Sequential Minimal Optimization)	۲۵
۴-۵- Core Vector Machines	۳۰
۴-۶- LSVM (Lagrangian SVM)	۳۲
۴-۷- Reduced SVM	۳۲
۴-۸- Relevance Vector Machine	۳۲
۴-۹- انتخاب مجموعه‌ی کاری	۳۳
۵- کاربردهای SVM	۳۵
۵-۱- تشخیص و شناسایی چهره	۳۵
۵-۲- تشخیص و شناسایی اشیاء	۳۵
۵-۳- تشخیص دست‌نوشته و ارقام	۳۶

۳۶	۴-۵- تشخیص صحبت و گوینده.....
۳۶	۵-۵- بازیابی تصاویر و اطلاعات.....
۳۷	۶-۵- پیش بینی.....
۳۷	۷-۵- سایر کاربردها.....
۳۸	۶- ابزارهای پیاده سازی SVM.....
۳۸	۶-۱- MATLAB.....
۴۱	۶-۲- LIBSVM.....
۴۲	۶-۳- svm-struct-matlab.....
۴۲	۷- جمع بندی و نتیجه گیری.....
۴۳	پیوست ۱.....
۴۴	پیوست ۲.....
۴۷	مراجع.....

فهرست تصاویر

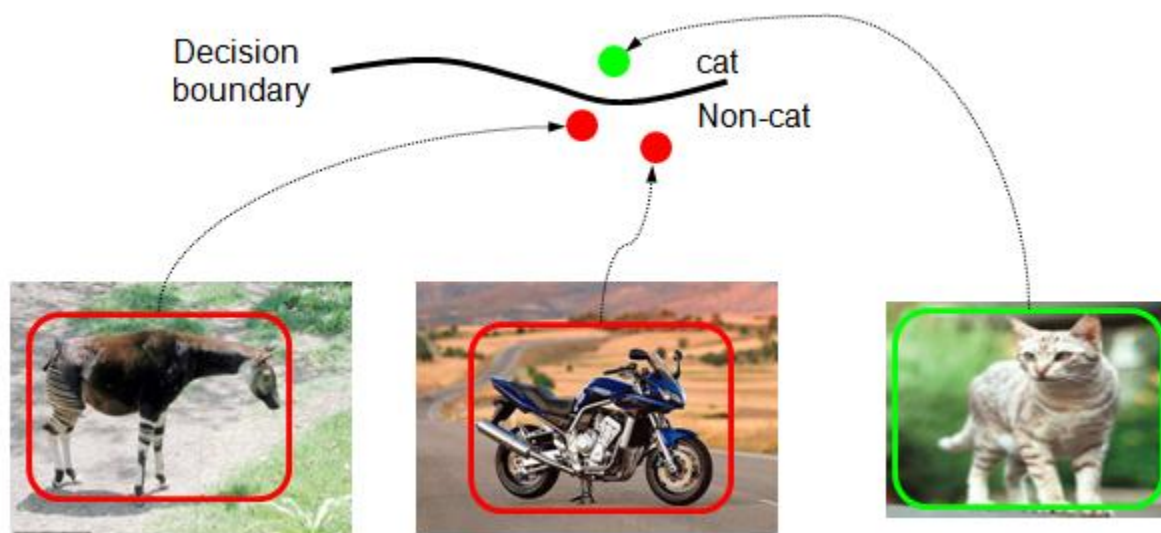
- شکل ۱- مثالی از طبقه بندی تصاویر..... ۱
- شکل ۲ - SVM به عنوان یک ابر صفحه برای جداسازی خطی نمونه ها در فضای داده ها..... ۴
- شکل ۳- تفکیک داده ها در فضا توسط ابر صفحه های مختلف..... ۵
- شکل ۴- با کاهش دقت نقاط آموزش، تفکیک کننده ای که از هر دو دسته بیشترین فاصله را دارد مناسب تر به نظر می رسد..... ۵
- شکل ۵- جداسازی خطی در حالت دوبعدی..... ۷
- شکل ۶- جداسازی خطی در حالت سه بعدی..... ۸
- شکل ۷..... ۹
- شکل ۸ - نمایش بردارهای پشتیبان..... ۱۲
- شکل ۹- داده های آموزشی در حالت جدایی ناپذیر..... ۱۵
- شکل ۱۰ - انتقال مسئله از فضای دو بعدی به سه بعدی..... ۱۶
- شکل ۱۱- الف) یک زیر مسئله ی بهینه سازی که نقاط شرایط بهینه بودن را نقض کرده و داخل ناحیه ی ± 1 قرار می گیرند.
ب) ناحیه ی تصمیم گیری مجدداً تعریف می شود. از آنجایی که هیچ نقطه ای با $\alpha=0$ در فاصله ی ± 1 قرار ندارد، لذا راه حل بهینه است. همانطور که دیده می شود حاشیه نیز کاهش یافته و شکل سطح تصمیم گیری متفاوت خواهد بود. [۱۴]
- ۲۱
- شکل ۱۲ - دو ضریب لاگرانژ باید تمامی قیدهای مطرح شده در مسئله را مرتفع نمایند. قید نامساوی در مسئله سبب می شود که ضرایب در چهارچوب محدود شده و شرط تساوی سبب واقع شدن آنها روی خط قطری می گردد. بنابراین SMO باید مقدار بهینه را روی این قطعه خط قطری بدست آورد..... ۲۷
- شکل ۱۳- مسئله ی MEB..... ۳۱
- شکل ۱۴- داده های آزمایشی تولید شده..... ۳۹
- شکل ۱۵- داده های دسته بندی شده با svm..... ۴۰

۱- مقدمه

در حالت کلی یادگیری ماشین به دو شکل نظارت شده و نظارت نشده انجام می شود. بسیاری از روش های یادگیری ماشین که به صورت نظارت شده عمل می نمایند، به این صورت کار می کنند که مجموعه ای از بردارهای ورودی مانند $X=\{x_n\}$ و بردارهای خروجی متناظر با آنها $T=\{t_n\}$ داده می شود. هدف این است که ماشین قادر باشد با استفاده از این داده های آموزشی برای ورودی x جدید، t را پیش بینی نماید. [6]

در این راستا می توان دو حالت متمایز را در نظر گرفت: رگرسیون (*Regression*)، که در آن t یک متغیر پیوسته است و طبقه بندی کردن (*Classification*)، که در آن t متعلق به یک مجموعه ی گسسته می باشد. بسیاری از مسائل یادگیری ماشین اساساً در دسته دوم واقع می شوند و هدف تعلیم ماشین به گونه ایست که ماشین بتواند بخوبی این دسته بندی را انجام دهد.

بعنوان مثال فرض کنید که هدف طراحی و تعلیم ماشینی است که قادر به تشخیص تصویر چهره (face) و غیر چهره (non-face) باشد. این سیستم عملاً کاری جز دسته بندی انجام نمی دهد. اگر برای هر ورودی (در این مثال تصویر ورودی) یک بردار n بعدی x فرض کنیم (مثلاً بردار ویژگی خاصی که از یک تصویر استخراج شده است) ماشین باید قادر باشد بردارهای x که بعنوان داده های تعلیمی هستند و در فضای n -بعدی به شکل نقاطی دیده می شوند را، دسته بندی نمایند. شکل زیر این دسته بندی را برای ماشینی نمایش می دهد که قادر به تشخیص تصویر گربه از سایر تصاویر می باشد.



شکل ۱- مثالی از طبقه بندی تصاویر

در فرآیند یادگیری نیاز است که ابتدا سیستم تعلیم یافته و سپس به ازای مقادیر ورودی جدید تست شود. به شکل ریاضی مسئله‌ی یادگیری ماشین را می‌توان به شکل یک نگاشت در نظر گرفت که طی آن $x_i \rightarrow y_i$. در واقع ماشین با مجموعه‌ای از نگاشت‌های ممکن به شکل $x \rightarrow f(x, \alpha)$ تعریف می‌شود که در آن خود توابع $f(x, \alpha)$ توسط برچسب α قابل تنظیم می‌باشند. فرض می‌شود که سیستم deterministic بوده و به ازای یک ورودی خاص x و انتخاب α ماشین همواره یک خروجی مشخص برابر با $f(x, \alpha)$ بدهد. انتخاب α ی مناسب همان عملی است که یک ماشین تعلیم دیده (Trained Machine) انجام می‌دهد. بعنوان مثال یک شبکه‌ی عصبی که α متناظر با وزن‌ها و مقدار بایاس در آن باشد، یک ماشین یادگیری می‌باشد.

در اینجا بر حالتی تمرکز خواهیم داشت که پیشگویی $y(x, w)$ توسط ترکیب خطی از تابع پایه‌ی $\Phi_m(x)$ به فرم زیر بیان می‌شود:

$$y(x, w) = \sum_{m=0}^M w_m \Phi_m(x) = w^T \phi$$

(۱)

که در آن w_m پارامترهای مدل می‌باشند که وزن نامیده می‌شوند. در SVM توابع پایه به عنوان توابع Kernel استفاده می‌شوند که به ازای هر x_m در مجموعه‌ی آموزشی داریم $\Phi_m(x) = K(x, x_m)$ که در آن $K(.,.)$ تابع Kernel می‌باشد.

تخمین وزن‌ها در SVM بوسیله‌ی بهینه‌سازی ضوابطی حاصل می‌شود که به صورت همزمان سعی در می‌نیم کردن تابع $y(x, w)$ دارند. در نتیجه تعدادی از وزن‌ها صفر می‌شوند، که سبب می‌شود مدل خلوتی (Sparse Model) حاصل شود که مدیریت پیش‌گویی آن توسط رابطه‌ی (۱) و فقط وابسته به زیرمجموعه‌ای از تابع Kernel باشد. [6]

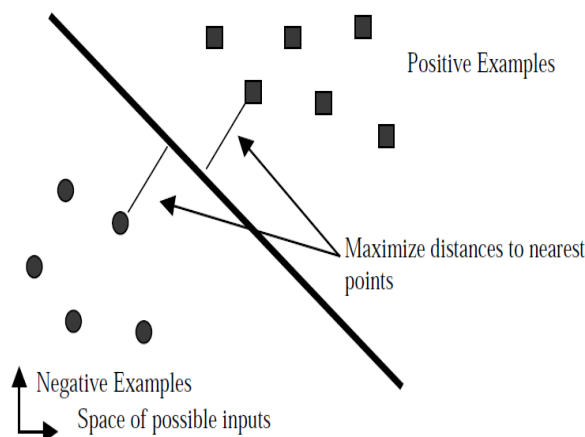
در سال‌های اخیر استفاده از ماشین‌های بردار پشتیبان (SVM)، بسیار مورد توجه قرار گرفته است. بطور تجربی نشان داده شده است که استفاده از SVM در کاربردهایی مانند تشخیص دست‌نوشته‌ها، تشخیص چهره و... نتایج خوبی را حاصل نموده است. تا زمان ارائه‌ی [3] استفاده از ماشین‌های بردار پشتیبان تنها به گروه خاصی از محققین محدود می‌شد. علت این محدودیت در استفاده از SVM، کندی الگوریتم آموزش برای این روش می‌باشد، خصوصاً در مورد مجموعه‌های آموزش بزرگ. به بیان دیگر الگوریتم آموزش SVM برای بسیاری از

مهندسين، پیچیده و دشوار بوده است و به همین دلیل راه حل ها و شیوه های جدیدی برای مواجهه با این مسئله ارائه شده است.

در ادامه ی این مقاله و در بخش دوم مروری بر مفهوم دسته بندی که در SVM انجام می گیرد خواهیم داشت و مثال هایی از کاربردهای این دسته بندی کننده را ذکر می نماییم، سپس در بخش سوم ابتدا راجع مسائل بهینه سازی مقید و حل آنها به روش ضرائب لاگرانژ پرداخته و در ادامه مسئله مربوط به یادگیری در SVM را در ساده ترین حالت به شکل یک مسئله ی برنامه ریزی درجه ی دوم فرموله می نماییم و در بخش سوم راجع به سایر حالت های ممکن صحبت خواهیم کرد. در بخش چهارم تعدادی از راه حل ها که برای محاسبه ی کارآمد مسئله ی QP در SVM ارائه شده اند معرفی می شوند و در بخش پنجم به ابزارهایی که برای پیاده سازی SVM طراحی گردیده اند اشاره می گردد. در بخش ششم به کاربردهای SVM پرداخته و در پایان نیز در بخش هفتم به جمع بندی و نتیجه گیری خواهیم رسید.

۲- مروری بر دسته‌بندی کننده‌ی SVM

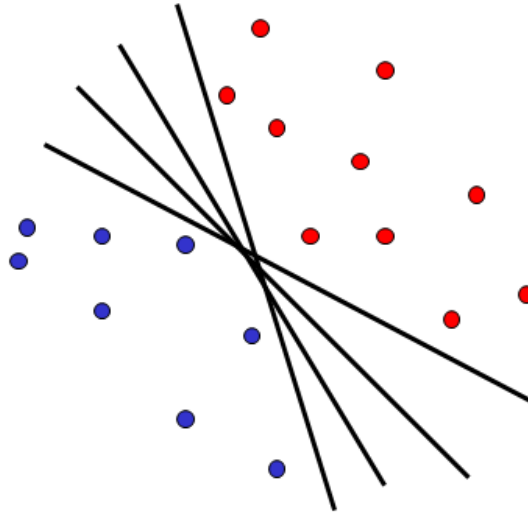
در ۱۹۷۹ ماشین‌های بردار پشتیبان (Support Vector Machines) توسط Vapnik ارائه گردید. در ساده‌ترین فرم آن یعنی SVM خطی، SVM عبارتست از یک ابر صفحه که مجموعه‌ی نمونه‌های مثبت و منفی را با حداکثر فاصله (maximum margin) از هم جدا نموده است. (شکل ۲)



شکل ۲ - SVM به عنوان یک ابر صفحه برای جداسازی خطی نمونه‌ها در فضای داده‌ها

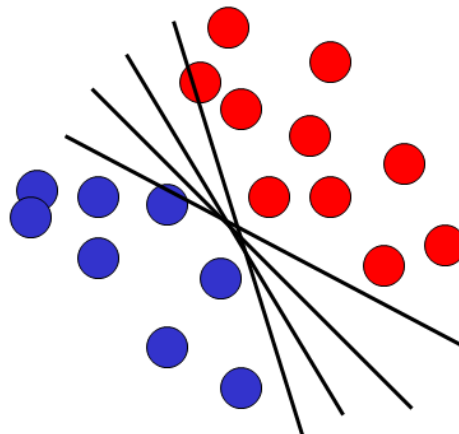
در شکل (۳) نمونه‌ای از داده‌های تفکیک‌پذیر (جدایی‌پذیر) در فضای دوبعدی دیده می‌شوند. (در ادامه گفته خواهد شد که این حالت ساده‌ترین حالت در دسته‌بندی (طبقه‌بندی) داده‌ها می‌باشد که از کمترین پیچیدگی برخوردار می‌باشد و داده‌ها در آن با استفاده از یک ابر صفحه به صورت خطی به دو دسته‌ی کاملاً مجزا تفکیک می‌شوند) همانطور که دیده می‌شود این داده‌ها را می‌توان توسط چندین خط جداکننده (یا ابر صفحه در فضای n-بعدی) تفکیک و طبقه‌بندی نمود و در عین حال تمام این دسته‌بندی‌ها نیز صحیح می‌باشند.

در حالت کلی این مسئله را می‌توان در فضای n-بعدی در نظر گرفت که در آن داده‌ها در دو دسته واقع شده‌اند که این دو دسته به خوبی قابل جداسازی می‌باشند. در این حالت بجای خطوط جداکننده از ابر صفحات جداکننده (Hyperplane) استفاده خواهد شد. در اینجا سؤال این است که در میان این جداکننده‌ها کدامیک بهترین می‌باشد؟ برای این منظور باید به مفهوم ریسک‌پذیری دقت نمود. در واقع در اینجا حالتی را می‌خواهیم انتخاب نماییم که بیشترین دقت را در تفکیک‌کنندگی ارائه دهد، یعنی با وجود مقداری خطا و جابجایی در داده‌ها تفکیک‌کنندگی همچنان صحیح باشد. مثلاً حالتی را در نظر بگیرید که هرکدام از داده‌های دارای حاشیه‌ای از خط باشند، در این



شکل ۳- تفکیک داده‌ها در فضا توسط ابرصفحه‌های مختلف

حالت اگر ابرصفحه‌ی جداکننده به اندازه‌ی کافی از داده‌های هر دسته دور نباشد، سبب خواهد شد که با اندکی خطا داده‌ها از مرز جدا کنندگی عبور نمایند. شکل ۴ این حالت را نمایش می‌دهد.



شکل ۴- با کاهش دقت نقاط آموزش، تفکیک‌کننده‌ای که از هر دو دسته بیشترین فاصله را دارد مناسب‌تر به نظر می‌رسد

همانطور که دیده می‌شود، جدا کننده‌ای مناسب است که بیشترین فاصله را از هر دو دسته از داده‌ها داشته باشد. به عبارت دیگر سبب شود بیشترین حاشیه‌ی ریسک ایجاد شود. بنابراین در مورد این ماشین که خروجی آن

برچسب گروه‌ها (در اینجا ۱- و ۱+) می‌باشد، هدف بدست آوردن ماکزیمم فاصله بین دسته‌های داده می‌باشد. این حالت به داشتن ماکزیمم حاشیه تعبیر می‌شود.

در ارتباط با SVM می‌توان مسئله‌ی طبقه‌بندی داده‌ها را در چند حالت مختلف بررسی نمود:

۱- استفاده از ماشین بردار پشتیبان خطی و وجود دو دسته داده:

الف - داده‌ها دقیقاً در دو دسته قرار گرفته باشند، یعنی جدایی‌پذیر باشند. (Separable Data)

ب - داده‌های دو دسته قابل تفکیک به دو دسته‌ی جدا از هم نباشند. (Nonseparable Data)

۲- ماشین بردار پشتیبان غیرخطی:

۳- ماشین بردار پشتیبان برای تفکیک‌کننده‌های چند کلاسه

۲-۱- حالت جدایی‌پذیر

ساده‌ترین حالت حالتی است که ماشین به صورت خطی روی داده‌های جدایی‌پذیر آموزش دیده باشد. جالب توجه می‌باشد که معادلات حاصل از این حالت قابل تعمیم به حالت غیرخطی و جدایی‌ناپذیر نیز می‌باشد، لذا این حالت اساس تعریف سایر حالات می‌باشد. لازم به ذکر است که در این وضعیت فرض شده که داده‌ها دقیقاً در دو دسته قرار دارند.

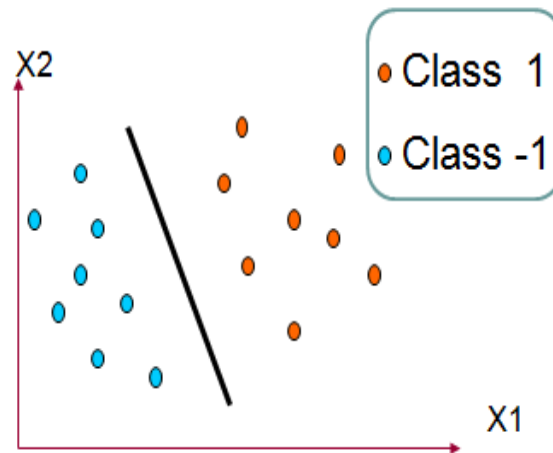
در این حالت فرض می‌شود که مجموعه‌ای از نمونه‌های آموزشی تفکیک‌پذیر وجود دارند که می‌توان آنها را با y_i برچسب زد. در این حالت نمونه‌ها به شکل زوج‌های مرتب (x_i, y_i) بیان می‌شوند که در آن $i=1, \dots, l$ و x_i متعلق به R^n و y_i دارای مقادیر $\{-1, 1\}$ می‌باشند. [1]

همانطور که در شکل ۳ دیده شد داده‌ها را می‌توان توسط چندین خط جداکننده (یا ابرصفحه در فضای n -بعدی) تفکیک و طبقه‌بندی نمود. برای انتخاب بهترین ابرصفحه‌ی جداکننده با توجه به مفهوم تعریف شده در مورد ماکزیمم حاشیه باید عمل نمود. در حالت خطی "حاشیه" به صورت فاصله‌ی ابر صفحه تا نزدیک‌ترین نمونه‌ی منفی و مثبت تعریف می‌شود. [4] هر خط تفکیک‌کننده در فضای ۲ بعدی به شکل معادله‌ی خط زیر نوشته می‌شود:

$$w_1x_1 + w_2x_2 + b = 0$$

(۲)

شکل ۵ این موضوع را برای x_i ها که در دو کلاس با برچسب های -1 و $+1$ واقع شده اند، در فضای دو بعدی نمایش می دهد:



شکل ۵- جداسازی خطی در حالت دوبعدی

خط جداکننده با فرمول ۲ نمایش داده می شود.

و در فضای ۳ بعدی به شکل زیر خواهد بود:

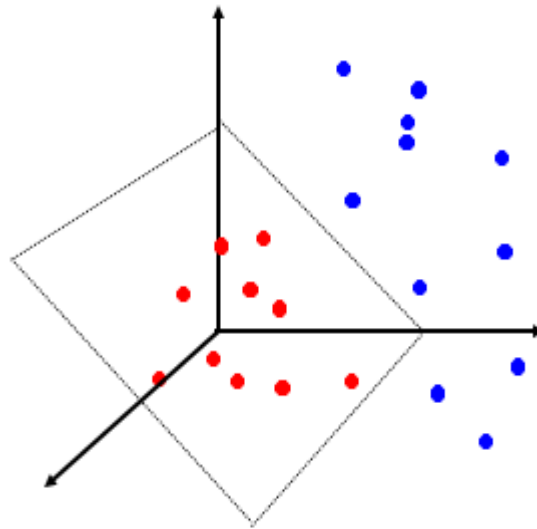
$$w_1x_1 + w_2x_2 + w_3x_3 + b = 0$$

(۳)

که در این حالت فرمول ۳ نمایش دهنده ی یک صفحه خواهد بود (شکل ۶). در حالت کلی و برای یک جداکننده از نوع ابرصفحه خواهیم داشت:

$$\sum_i w_i x_i + b = 0$$

(۴)



شکل ۶- جداسازی خطی در حالت سه بعدی

که در نمایش برداری به شکل زیر قابل بیان می باشد:

$$u = \vec{w} \cdot \vec{x} + b$$

(۵)

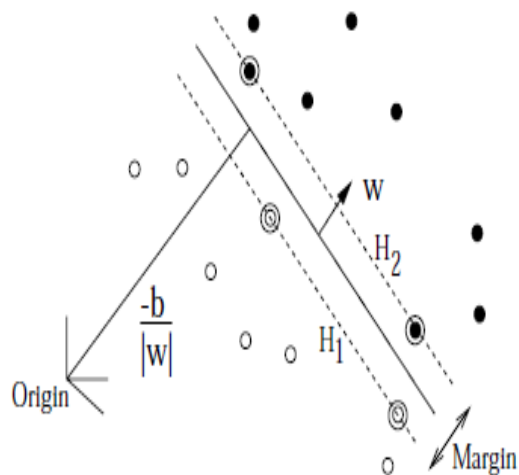
که در آن w بردار وزن بردار عمود بر ابرصفحه و b مقدار bias می باشد. در این نمایش $u=0$ اشاره به خود ابرصفحه جداکننده داشته و نزدیک ترین نقاط بر روی صفحه $u=\pm 1$ قرار دارند. در واقع به فرض تفکیک پذیر بودن دو کلاس داده ی مثبت و منفی، بردارهای مرزی روی ابر صفحه های زیر قرار خواهند گرفت:

$$\vec{w} \cdot \vec{x} + b = \pm 1$$

(۶)

ناحیه ی بین این دو ابرصفحه، حاشیه (Margin) گفته می شود.

شکل ۷ بیانگر حالت دو بعدی با فرض منفی بودن مقدار bias می باشد:



شکل ۷

همانطور که در شکل دیده می شود فضا به دو دسته از نمونه های با ویژگی های زیر تقسیم می شود:

$$\vec{w} \cdot \vec{x} + b \geq +1 \quad \text{for } y_i = 1$$

(۷)

$$\vec{w} \cdot \vec{x} + b \leq -1 \quad \text{for } y_i = -1$$

(۸)

می توان فرمول فوق را به شکل زیر ترکیب نمود:

$$y_i(\vec{w} \cdot \vec{x} + b) - 1 \geq 0 \quad \forall i$$

(۹)

در این حالت فاصله تا مبداء به صورت عمودی (نزدیکترین فاصله) برای نقاطی که روی ابر صفحه ی $x \cdot w + b = +1$ قرار دارند برابر با :

$$\frac{|1 - b|}{\|w\|}$$

(۱۰)

و به شکل مشابه فاصله تا مبدا به صورت عمودی برای نقاطی که روی ابر صفحه‌ی $x.w+b=-1$ قرار دارند برابر با :

$$\frac{|-1 - b|}{\|w\|} \quad (11)$$

خواهد بود. [5] از طرفی فاصله‌ی مبدا تا ابر صفحه‌ی جداکننده برابر با :

$$\frac{|b|}{\|w\|} \quad (12)$$

می‌باشد. بنابراین کمترین فاصله‌ی بین این ابر صفحه تا هر یک از صفحات مذکور به شکل زیر خواهد بود:

$$d_+ = d_- = \frac{1}{\|w\|} \quad (13)$$

بنابراین حاشیه‌ی مذکور یعنی فاصله‌ی بین دو ابر صفحه‌ی مورد نظر به شکل زیر بدست می‌آید:

$$d_+ + d_- = \frac{2}{\|w\|} \quad (14)$$

به این ترتیب حداکثر شدن حاشیه می‌تواند به شکل معادله‌ی بهینه‌سازی مقید زیر بیان شود:

$$\min \frac{1}{2} \|w\|^2 \quad \text{subject to: } y_i(\vec{w} \cdot \vec{x} + b) - 1 \geq 0 \quad \forall i \quad (15)$$

برای حل این مساله ساده‌تر آن است که مساله بصورت دوگان (Dual) مطرح و حل شود. برای بدست آوردن شکل دوگان مسئله ضرایب لاگرانژ مثبت $\alpha_i \geq 0$ در قیود ضرب شده و از تابع هدف کسر می‌شود، در نتیجه به معادله‌ی زیر می‌رسیم که Primal Problem گفته می‌شود:

$$L_p = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i (y_i(\vec{w} \cdot \vec{x} + b) - 1) \quad (16)$$

شرایط Karush-Kuhn-Tucker (KKT) نقش مهمی را در مسائل بهینه سازی مقید، بازی می کنند. این شرایط، شرایط لازم و کافی برای داشتن پاسخ بهینه برای معادلات مقید را بیان می دارند و باید مشتق تابع نسبت به متغیرها برابر صفر باشد. با اعمال شرایط KKT روی L_p و اگر از L_p نسبت به w و b مشتق گرفته شده و برابر صفر قرار داده شود، خواهیم داشت:

$$w = \sum_i \alpha_i y_i x_i$$

(۱۷)

$$0 = \sum_i \alpha_i y_i$$

(۱۸)

در واقع در اینجا و برای Primal Problem (L_p)، شرایط KKT به شکل زیر خواهند بود:

$$\frac{\partial}{\partial w_v} L_p = w_v - \sum_i \alpha_i y_i x_{iv} = 0 \quad v = 1, \dots, d$$

(۱۹)

$$\frac{\partial}{\partial b} L_p = - \sum_i \alpha_i y_i = 0$$

(۲۰)

$$y_i(\vec{w} \cdot \vec{x} + b) - 1 \geq 0 \quad i = 1, \dots, l$$

$$\alpha_i \geq 0 \quad \forall i$$

(۲۱)

$$\alpha_i (y_i(\vec{w} \cdot \vec{x} + b) - 1) \geq 0 \quad \forall i$$

(۲۲)

که از نتیجه ی آن برای بدست آوردن دوگان مسئله (L_D) استفاده شده است. شرایط KKT هر گونه مسئله ی بهینه سازی مقیدی را (چه محدب و چه غیر محدب) با هر نوع قیدی، برآورده می سازد. مسئله ی SVM یک مسئله ی محدب می باشد و برای مسائل محدب شرایط KKT شرایط لازم و کافی می باشند تا w, b, α بتوانند راه حل مسئله باشند. بنابراین حل مسئله ی SVM معادل است با یافتن یک راه حل برای شرایط KKT، که از آن در

روش های مختلفی برای یافتن راه حل SVM استفاده شده است. یک کاربرد از این شرایط این است که اگر w در روال آموزش تعیین شده و b مشخص نشده باشد، می توان b را بدست آورد.

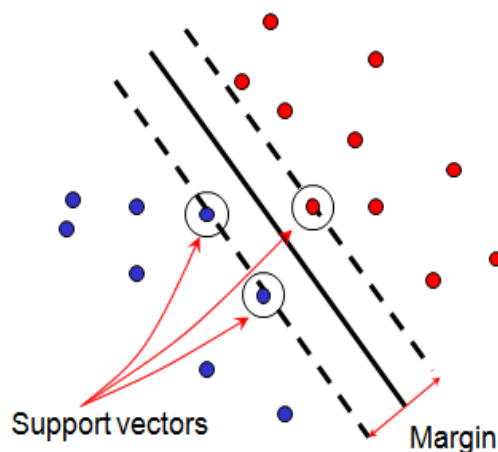
با جایگذاری در معادله ی ۱۶ خواهیم داشت :

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j) \quad (23)$$

که معادله ی دوگان^۱ نامیده می شود. L_D و L_P هر دو از یک شرایط بدست آمده اند، لذا با محاسبه ی می نیمم L_P و یا ماکزیمم L_D (که دوگان L_P می باشد) با شرط $\alpha_i \geq 0$ می توان به جواب رسید.

برای هر نمونه ی آموزشی (x_i) یک ضریب لاگرانژ $\alpha_i \geq 0$ وجود دارد. در حل مساله به نمونه هایی که به ازای آنها $\alpha_i > 0$ می باشد، بردارهای پشتیبان یا "Support Vectors" گفته می شوند که بر روی ابر صفحات معادله ی (۶) قرار می گیرند (شکل ۸).

مسئله ی به دست آمده برای SVM یک مسئله ی بهینه سازی مقید درجه ی دوم و محدب می باشد.



شکل ۸ - نمایش بردارهای پشتیبان

^۱ Dual Problem

قابل توجه می باشد که تمامی کارهایی که انجام می شود برای آن است که مسئله به یک مسئله بهینه سازی تبدیل شود که قیود آن قابل مدیریت تر باشند نسبت به (۷) و (۸). یافتن راه حل برای مسئله مطرح شده مستلزم روش های عددی می باشد که به ازای حجم زیاد داده ها بسیار مشکل و زمانبر خواهد بود. در قسمت های بعدی به راه حل هایی که برای کاهش پیچیدگی این محاسبات ارائه شده است اشاره خواهد گردید.

جهت تکمیل این بحث روش ضرائب لاگرانژ و شرایط KKT برای حل مسائل بهینه سازی مقید در ضمیمه ۲ آمده است.

۳- SVM در حالت های غیرخطی و جدایی ناپذیر

در این بخش سایر حالت های ممکن برای دسته بندی داده ها را بررسی می نماییم.

۳-۱- داده های جدایی ناپذیر

در این مورد فرض بر این است که داده های موجود به سادگی قابل تفکیک به دو دسته نباشند. در واقع در این حالت که به واقعیت نیز نزدیک تر است، داده ها با نویز همراه بوده و لازم است SVM برای غلبه بر این حالت مقداری عمومی سازی (generalization) گردد. این وضعیت با تعریف یک حاشیه نرم (Soft Margin) قابل حل می باشد [9]، زیرا اگر الگوریتم فوق برای حالت جدایی پذیر به مجموعه ای از عناصر جدایی ناپذیر اعمال شود راه حل ممکن پیدا نخواهد شد. برای گسترش ایده ی مطرح شده در بخش قبلی به این حالت، متغیری تحت عنوان متغیر مثبت کاهشی (کمبود)^۲ در قیود مسئله تعریف می شود و به این ترتیب خواهیم داشت:

$$x_i \cdot w + b \geq +1 - \xi_i \quad \text{for } y_i = 1$$

(۲۴)

$$x_i \cdot w + b \leq -1 + \xi_i \quad \text{for } y_i = -1$$

(۲۵)

$$\xi_i \geq 0 \quad \forall i$$

(۲۶)

بنابراین برای هر خطایی که رخ می دهد، هر ξ_i باید صورت جداگانه افزایش می یابد. بنابراین $\sum_i \xi_i$ به عنوان یک کران بالا برای تعداد خطاهای آموزش می باشد. در نتیجه بهترین راه برای نسبت دادن یک هزینه اضافی

^۲ Slack variable

برای خطاها، تغییر تابع هدف برای به حداقل رسانیدن تابع $\frac{\|w\|^2}{2} + C(\sum_i \xi_i)^k$ بجای $\frac{\|w\|^2}{2}$ می باشد. مقدار C توسط کاربر انتخاب می گردد. یک مقدار C بزرگ به معنی اعمال جریمه ی بیشتر برای خطاها می باشد. تابع مطرح شده به ازای هر مقدار مثبت از k تابعی محدب می باشد، خصوصاً به ازای $k=1$ و $k=2$ به یک مسئله ی برنامه ریزی درجه ی دوم تبدیل می شود. با محاسبه ی دوگان تابع مشابه با حالت جدایی پذیر، به ۲۷ می رسم:

$$\text{Maximize: } L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j) \quad (27)$$

subject to :

$$0 \leq \alpha_i \leq C,$$

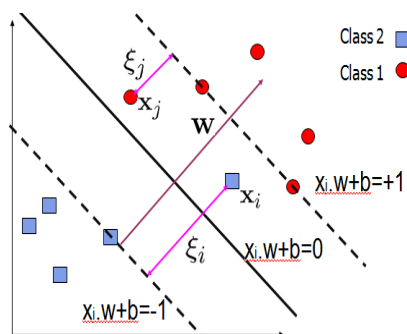
$$\sum_i \alpha_i y_i = 0$$

(28)

در این شرایط راه حل از طریق ۲۹ ارائه می گردد:

$$w = \sum_{i=1}^{N_s} \alpha_i y_i x_i \quad (29)$$

در ۲۹، N_s اشاره به تعداد بردارهای پشتیبان دارد. بنابراین تنها اختلاف با حالت ابرصفحه ی بهینه در این است که در این حالت α_i دارای کران بالای C می باشد. [5]



شکل ۹- داده های آموزشی در حالت جدایی ناپذیر

مسئله ی primal لاگرانژ به شکل ۳۰ خواهد بود:

$$L_p = \frac{\|w\|^2}{2} + C \left(\sum_i \xi_i \right)^k - \sum_i \alpha_i \{y_i(x_i \cdot w + b) - 1 + \xi_i\} - \sum_i \mu_i \xi_i \quad (30)$$

که در آن μ_i ضرائب لاگرانژ می باشند که سبب مثبت شدن ξ_i می گردند.

با اعمال شرایط KKT بر مسئله ی primal حاصل (L_p) خواهیم داشت:

$$\frac{\partial}{\partial w_v} L_p = w_v - \sum_i \alpha_i y_i x_{iv} = 0 \quad v = 1, \dots, d \quad (31)$$

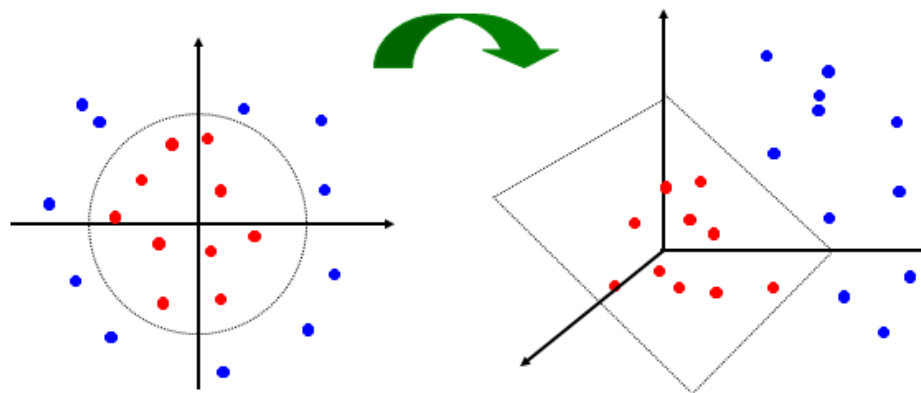
$$\frac{\partial}{\partial b} L_p = - \sum_i \alpha_i y_i = 0 \quad (32)$$

$$\frac{\partial}{\partial \xi_i} L_p = C - \alpha_i - \mu_i = 0 \quad (33)$$

که با اعمال این شرایط در L_p می توان به L_D (۲۷) رسید. پس از محاسبه ی ضرائب لاگرانژ، $0 \leq \alpha_i \leq C$ بردارهای پشتیبان برای این حالت می باشند.

۳-۲- ماشین های بردار پشتیبان غیر خطی

در حالتی که داده ها به سادگی از هم تفکیک نمی شوند، یک تفکیک کننده ی خطی نمی تواند موثر باشد. اما اگر داده ها را به فضایی با ابعاد بیشتر انتقال دهیم، می توان راه حلی برای تفکیک آنها به دست آورد. به عنوان مثال شکل ۱۰ مثالی از این انتقال را نمایش می دهد.



شکل ۱۰ - انتقال مسئله از فضای دو بعدی به سه بعدی

همانطور که دیده می شود پس این انتقال می توان کلاسه بندی را انجام داد. در حالت های ذکر شده در بخش های گذشته، در واقع از حاصلضرب داخلی داده های آموزشی ورودی در بردارهای پشتیبان (۲۷) برای تشکیل جداکننده ی خطی به شکل ابرصفحه استفاده شده است. در اینجا ابتدا داده ها را به یک فضای اقلیدسی با ابعاد بالاتر نگاشت نموده و سپس از ضرب داخلی عناصر بدست آمده استفاده می نماییم. [5]

$$\Phi: R^d \rightarrow \mathcal{H}$$

(۳۴)

این فضای با ابعاد بالاتر و کاملاً برداری فضای هیلبرت خوانده می شود. [12] بنابراین الگوریتم آموزش وابسته به این ضرب داخلی در فضای \mathcal{H} به شکل $\Phi(x_i) \cdot \Phi(x_j)$ خواهد شد. اگر یک تابع هسته^۳ به شکل زیر تعریف کنیم:

$$K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$$

(۳۵)

آنگاه در الگوریتم یادگیری تنها کافیت که بدون اینکه بطور صریح به Φ اشاره کنیم، فقط از K استفاده نماییم. شرایط فوق بیانگر یک جداسازی خطی ولی در ابعاد بالاتر می باشد. Boser, Guyon and Vapnik, 1992 نشان داده اند با این نگاشت و در صورتی که تابع هسته شرایط قضیه ی Mercer^۴ را دارا باشد، هسته ی مورد استفاده مناسب بوده و داده های

^۳ Kernel Function

^۴ Mercer's Condition

جدایی ناپذیر در فضای اصلی در فضای نگاشته شده، جدایی پذیر خواهند بود و لذا جایگزینی هسته می تواند الگوریتمی غیرخطی را از الگوریتم خطی ذکر شده ایجاد نماید. [5,9]

در زیر برخی از توابع Kernel پرکاربرد آمده اند:

$$K(x_i, x_j) = (x_i \cdot x_j + 1)^p \quad (36)$$

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}} \quad (37)$$

$$K(x_i, x_j) = \tanh(\beta x_i \cdot x_j + \delta) \quad (38)$$

به طور مشابه با داشتن هسته های ذکر شده دوگان مسئله ی بهینه سازی به شکل (38) خواهد بود:

$$\text{Maximize: } L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad (39)$$

subject to :

$$0 \leq \alpha_i \leq C,$$

$$\sum_i \alpha_i y_i = 0$$

(40)

و همچنین به شکل مشابه بردارهای پشتیبان نمونه هایی هستند که ضرایب لاگرانژ متناظر آنها در رابطه ی $0 < \alpha_i \leq C$ صادق باشند.

۳-۲- ماشین های بردار پشتیبان به عنوان جداکننده ی چند کلاسه

در دنیای واقعی بسیاری از مسائل شامل دسته بندی های چند کلاسه می باشد و فقط به دسته بندی های دودویی محدود نمی شود. راه حل های مختلفی برای این دسته ارائه شده است. شکل ۱۰ یکی از ساده ترین راه حل ها را که در آن از یک دسته بندی دودویی استفاده می شود نمایش می دهد. [9] عموماً برای این دسته از کلاسه بندی ها از ترکیبی از روش های

دودویی استفاده می شود. یک روش ساده ی دیگر روش یکی در مقابل بقیه^۵ می باشد، در این روش یک دسته به عنوان داده های مثبت و بقیه منفی فرض می شوند. [5]

۴- راه حل های بهینه جهت حل معادله ی QP برای SVM

همانطور که دیده شد، آموزش ماشین بردارهای پشتیبان، منجر به حل مسئله ی بهینه سازی درجه ی دوم زیر می شود:

$$\text{Minimize: } W(\alpha) = - \sum_i \alpha_i + \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad i, j \in \{1, \dots, l\}$$

(۴۱)

subject to :

$$0 \leq \alpha_i \leq C,$$

$$\sum_i \alpha_i y_i = 0$$

(۴۲)

تعداد نمونه های آموزشی، l می باشد. α ها برداری با l متغیر می باشند و α_i ها نظیر به داده های آموزشی (x_i, y_i) می باشند و نیز $K: R^n \times R^n \rightarrow R$ تابع هسته و C ثابتی مثبت است. $\alpha^* = [\alpha_1, \dots, \alpha_l]^T$ با توجه به قیود ذکر شده بعنوان جواب معادله ی فوق می باشد. با تعریف ماتریس Q به شکل $Q_{ij} = y_i y_j k(x_i, x_j)$ می توان مسئله ی فوق را به شکل زیر بازنویسی نمود: [۱۶]

$$\text{minimize } W(\alpha) = -\alpha^T \mathbf{1} + \frac{1}{2} \alpha^T Q \alpha$$

Subject to : $\alpha^T y = 0$

$$0 \leq \alpha \leq C \mathbf{1}$$

(۴۳)

سختی حل مسئله ی QP مطرح شده در اندازه ماتریس ایجاد شده می باشد (ماتریس Q) که عناصر آن عموماً غیر صفر می باشند، بنابراین حجم زیادی از حافظه نیاز است که بتوان این ماتریس و محاسبات مربوطه را در آن گنجاند. همچنین روش های سنتی مانند روش های نیوتنی و یا شبه- نیوتنی^۶ و ... قادر به حل این مسئله به صورت مستقیم نخواهند بود.

^۵ one-versus-rest classifiers

^۶ Quasi-Newton

مسئله‌ی بهینه‌سازی بردارهای پشتیبان فقط در صورتی می‌تواند به شکل تحلیلی حل شود که مجموعه‌ی آموزش عناصر بسیار کمی داشته باشد و یا در مورد تفکیک‌پذیر، زمانی که از قبل بدانیم که کدام عناصر از داده‌های آموزشی Support Vector هستند. SVM از توابع Kernel استفاده می‌نماید. بسیاری از توابع Kernel به شکل برنامه‌ریزی درجه‌ی دوم فرموله می‌شوند. اگر m تعداد داده‌های آموزشی باشد، پیچیدگی مرتبه‌ی زمانی برای یک مسئله‌ی برنامه‌ریزی درجه‌ی دوم، از مرتبه‌ی $O(m^3)$ و پیچیدگی فضایی آن از مرتبه‌ی حداقل درجه‌ی دوم می‌باشد. در مثال دنیای واقعی معادله‌ی بهینه‌سازی SVM از نوع برنامه‌ریزی درجه‌ی دوم، باید به صورت عددی محاسبه گردد. در مواردی با اندازه کوچک، هر راه حلی که برای برنامه‌ریزی خطی درجه‌ی دوم از نوع محدب وجود دارد، قابل استفاده می‌باشد.

با توجه به ابعاد بسیار بزرگ مسئله‌ی QP ارائه شده در حالت‌های مختلف در بخش قبل (LD) که از SVM بدست می‌آید، می‌توان فهمید که این مسئله به سادگی و با روش‌های حل بهینه‌سازی درجه‌ی دوم قابل حل نمی‌باشد. مسئله‌ی درجه‌ی دوم ارائه شده شامل ماتریسی می‌باشد که تعداد عناصر آن مربع تعداد عناصر مجموعه‌ی آموزشی می‌باشد..

برای مسائل بزرگ روش‌های و تکنیک‌های مختلفی وجود دارد که در آنها بر خلاف روش نیوتن که حل ماتریس Hessian مربوط به کل مسئله را در محاسبات دخالت می‌دهد، فقط قسمتی از ماتریس که در هر تکرار نیاز می‌باشد، مورد استفاده قرار می‌گیرد. [۱۵]

اساساً این راه‌حل‌ها به صورت تکراری عمل نموده و در هر تکرار مجموعه‌ای کوچک از متغیرها که به مجموعه‌ی کاری^۷ معروف است را انتخاب نموده و با توجه به این مجموعه مسئله‌ی QP را حل می‌نمایند.

ایده‌ی کلی به شکل مراحل اساسی زیر قابل بیان می‌باشند: [5]

۱- توجه به شرایط بهینه بودن (KKT) که راه حل باید آنرا دارا باشد.

۲- تعریف یک استراتژی برای رسیدن به حالت بهینه با کاهش یکنواخت دوگان تابع هدف با در نظر گرفتن قیود مسئله.

۳- تصمیم‌گیری برای یک الگوریتم تجزیه بطوریکه در هر زمان تنها بخش‌هایی از داده‌های آموزش نیاز به مدیریت داشته باشند.

برای کاهش پیچیدگی زمانی و فضایی، روش‌های مختلفی برای تخمین ماتریس Kernel کاهش یافته با استفاده از روش Nystrom (Williams and Seeger, 2001)، تخمین حریصانه (Smola and Scholkopf, 2000)، نمونه‌برداری (Achlioptas et al., 2002) و تجزیه‌ی ماتریس (Fine and Scheinberg, 2001) ابداع شده است. از میان کاره‌های انجام

شده، روش Decomposition Method [۱۴]، از همه مهمتر می باشد. به همین سبب در قسمت بعد، ابتدا این روش بررسی شده و سپس به تعدادی دیگر از روش ها به صورت گذرا اشاره خواهد شد:

۱-۴ - Decomposition Method :

در ۱۹۹۷ Osuna اثبات نمود که یک مسئله ی بزرگ QP می تواند به یک سری از زیر مسئله های QP کوچکتر شکسته شود. تا زمانی که حداقل یک نمونه که شرایط KKT را محقق نمی نماید (نقض می نماید)، به نمونه های زیر مسئله ی قبلی اضافه شده، هر مرحله کل تابع هدف را کاهش می دهد و نقطه ای را که قیود را رعایت می کند نگهداری می نماید. بنابراین دنباله ای از زیر مسئله های QP که همیشه حداقل یک عنصر را اضافه می کنند همگرایی را گارانتی می نماید. [4]

الگوریتم کلی که در این روش استفاده می شود، اساس بسیاری از روش های دیگر می باشد. هر روش مبتنی بر تجزیه^۸،^۳ ویژگی اصلی دارد: [۲۱]

- نحوه ی انتخاب مجموعه ی کاری

- نحوه ی حل زیر مسئله ی ایجاد شده

- شرایط پایان الگوریتم

این روش که توسط Osuna در [۱۴] مطرح شده است، مسئله ی بهینه سازی بدست آمده را به دو قسمت inactive و active (یا همان working set) تقسیم می نماید. [۱۶]

در [۱۴] بهینه سازی سراسری را تضمین نموده و می تواند در زمانی که حجم داده های آموزشی بسیار بالا می باشد، کارآمد باشد. ایده ی اصلی در این الگوریتم حل زیر مسئله ها و بررسی شرایط بهینه بودن در هر تکرار می باشد. همچنین در این مقاله این الگوریتم برای کاربرد تشخیص چهره بکار رفته که تعداد نمونه ها در آن ۵۰۰۰۰ نمونه بوده است.

حافظه ی مورد نیاز برای حل یک مسئله ی QP با ابعاد مثال ارائه شده در این مقاله (ماتریسی با ابعاد ۵۰۰۰۰*۵۰۰۰۰) و با فرض اینکه برای هر نمونه تعداد ۸ بایت اختصاص داده شده باشد، 20GB خواهد بود.

برای ارائه ی الگوریتم مطرح شده ابتدا ۲ شرط مطرح می شود:

۱- شرایط بهینه بودن^۹: این شرایط به ما اجازه می دهد که به صورت محاسباتی در مورد بهینه شدن راه حل در یک تکرار خاص تصمیم گیری نماییم.

^۸ Decomposition Method

^۹ Optimality Conditions

۲- استراتژی بهبود^{۱۰}: اگر یک راه حل خاص بهینه نباشد، این استراتژی راهی برای بهبود تابع هزینه ارائه می دهد. این استراتژی عمدتاً در ارتباط با متغیرهایی می باشد که شرایط بهینه بودن را نقض می نمایند.

الگوریتم Decomposition :

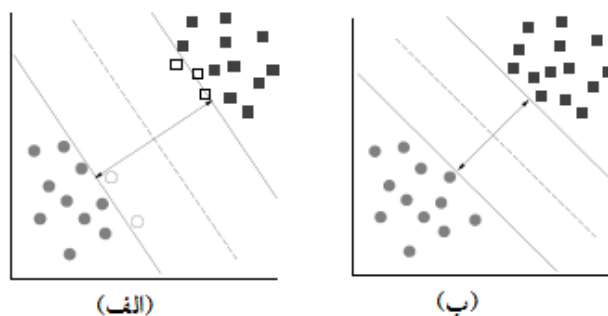
در هر تکرار α_i های مرتبط به مسئله، به دو دسته تقسیم می شوند (B, N). فرض کنید بتوان مجموعه ی کاری B را طوری انتخاب نمود که $|B| \leq 1$ و همچنین به اندازه ی کافی بزرگ باشد که بتواند شامل بردارهای پشتیبان باشد ($\alpha_i > 0$) و نیز به اندازه ای نیز کوچک باشد که انجام محاسبات روی آن مقدور باشد (توسط کامپیوتر). عناصر مجموعه ی B شرایط بهینه بودن را دارا هستند و سایر عناصر که این شرایط را نقض می کنند در مجموعه ی N دسته بندی می گردند. با این شرایط الگوریتم به شکل زیر خواهد بود :

۱- به صورت دلخواه $|B|$ نقطه را از مجموعه ی داده ها انتخاب کنید.

۲- زیر مسئله ای را که توسط متغیرهای B تعریف می شود حل کنید.

۳- تا زمانی که $j \in N$ وجود دارد بگونه ای که $g(x_j)y_j < 1$ ، $\alpha_i = 0, i \in B$ ، با $\alpha_j = 0$ جایگزین نموده و زیر مسئله ی جدید را حل کنید. عبارت دیگر یعنی تا زمانی که شرایط بهینه بودن نقض می شود ($g(x_j) = (\sum_{i=1}^l \alpha_i u_j K(x_i, x_j)) + b$)

باید توجه نمود که از آنجایی که این الگوریتم به صورت تکه ای تابع هدف را در هر تکرار بهبود می بخشد، لذا دارای



شکل ۱۱- الف) یک زیر مسئله ی بهینه سازی که نقاط شرایط بهینه بودن را نقض کرده و داخل ناحیه ی ± 1 قرار می گیرند. ب) ناحیه ی تصمیم گیری مجدداً تعریف می شود. از آنجایی که هیچ نقطه ای با $\alpha = 0$ در فاصله ی ± 1 قرار ندارد، لذا راه حل بهینه است. همانطور که دیده می شود حاشیه نیز کاهش یافته و شکل سطح تصمیم گیری متفاوت خواهد بود. [۱۴]

دور^{۱۱} نبوده و تابع هدف دارای کران می باشد (محدب بوده و توسط ناحیه ای امکان پذیر، محدود می باشد) و الگوریتم در تعداد محدودی تکرار به مقدار بهینه ی سراسری همگرا خواهد شد. شکل ۱۱، تفسیر هندسی از روشی را نشان می دهد که الگوریتم برای تعریف مجدد سطح جدا کننده با افزودن نقاطی که شروط بهینه را نقض می نمایند، بکار می برد. در [۱۶] روش فوق به شکل الگوریتم زیر آورده شده است :

While (the optimality conditions are violated)

Select q variables for the working set B . The remaining $l-q$ variables are fixed at their current value.

Decompose problem and **solve** QP-subproblem: **optimize** $W(\alpha)$ on B .

Terminate and **return** α .

از آنجایی که ماتریس Q ، positive-semidefinite می باشد و تمام قیود ذکر شده خطی می باشند، لذا مسئله ی مورد نظر یک مسئله ی بهینه سازی محدب خواهد بود. برای این دسته از مسائل شرایط KKT، شرایط لازم و کافی برای بهینه بودن می باشند. با توجه به قیود مسئله، یعنی :

$$\text{Subject to : } \alpha^T y = 0$$

$$0 \leq \alpha \leq C1$$

(۴۴)

ضرائب لاگرانژ را به این ترتیب علامت گذاری می نماییم: برای شرایط تساوی با λ^{eq} و برای کران بالا λ^{up} و پایین با λ^{lo} به این ترتیب α بهینه خواهد بود اگر λ^{eq} ، λ^{up} و λ^{lo} وجود داشته باشند، بنابراین با اعمال شرایط KKT خواهیم داشت:

$$g(\alpha) + (\lambda^{eq} y - \lambda^{lo} + \lambda^{up}) = 0$$

(۴۵)

$$\forall i \in [1..l]: \quad \lambda_i^{lo}(-\alpha_i) = 0$$

$$\forall i \in [1..l]: \quad \lambda_i^{up}(\alpha_i - C) = 0$$

$$\lambda^{lo} \geq 0$$

$$\lambda^{up} \geq 0$$

$$\alpha^T y = 0$$

$$0 \leq \alpha \leq C1$$

(۴۶)

^{۱۱} cycle

$g(\alpha)$ بردار مشتق های جزئی در α است. برای این مسئله داریم:

$$g(\alpha) = -1 + Q\alpha$$

(۴۷)

همانطور که در الگوریتم مشخص است، اگر شرایط بهینه بودن وجود نداشته باشد، الگوریتم مسئله را می شکند و به حل مسئله ی کوچکتر می پردازد. با توجه به تقسیم نمونه ها در دو دسته ی B و N، می توان ماتریس Q و متغیرهای α و y را نیز متناظر با B و N اندیس گذاری نمود:

$$\alpha = \begin{bmatrix} \alpha_B \\ \alpha_N \end{bmatrix} \quad y = \begin{bmatrix} y_B \\ y_N \end{bmatrix} \quad Q = \begin{bmatrix} Q_{BB} & Q_{BN} \\ Q_{NB} & Q_{NN} \end{bmatrix}$$

(۴۸)

با توجه به تقارن Q داریم که $Q_{NB} = Q_{BN}^T$ ، در نتیجه می توان مسئله را به شکل زیر باز نویسی نمود:

$$W(\alpha) = -\alpha_B^T(1 - Q_{BN}\alpha_N) + \frac{1}{2}\alpha_B^T Q_{BB}\alpha_B + \frac{1}{2}\alpha_N^T Q_{NN}\alpha_N - \alpha_N^T \mathbf{1}$$

(۴۹)

$$\text{subject to: } \alpha_B^T y_B + \alpha_N^T y_N = 0$$

$$0 \leq \alpha \leq C \mathbf{1}$$

(۵۰)

با توجه به اینکه متغیرهای مجموعه ی N به صورت ثابت فرض می شوند، لذا عبارت $\frac{1}{2}\alpha_N^T Q_{NN}\alpha_N - \alpha_N^T \mathbf{1}$ مقداری ثابت خواهد بود. مسئله ی به دست آمده یک مسئله ی بهینه سازی درجه دوم از نوع positive semidefinite می باشد که با توجه به کوچک بودن نسبی ابعاد مسئله، با استفاده از روش های مرسوم قابل حل می باشد. در این روال مهمترین فرایند انتخاب α برای تشکیل B می باشد، توجه داریم که B در هر تکرار به روز می شود. بنابراین سرعت روش حل می تواند وابسته به انتخاب صحیح مجموعه ی کاری باشد.

در خصوص همگرایی این روش در [۲۳] یک حالت ساده از شرایط بهینه بودن برای مسئله فرض شده و سپس نشان داده شده است که این روش می تواند با رعایت شرایط بهینه بودن مفروض، پس از تعداد محدودی تکرار به جواب بهینه همگرا شود.

۴-۲- SVM^{light}:

در [۱۶] روشی ارائه شده با عنوان SVM^{light} که در آن تعداد ثابتی، که می تواند هر عدد زوج کمتر از تعداد نمونه های آموزشی باشد، از متغیرها با استفاده از روش Steepest Descent به صورت متقارن، برای مجموعه ی کاری انتخاب می شوند و زیر مسئله ی ایجاد شده با روش های سنتی قابل حل می باشد.

پس از انتخاب مجموعه ی کاری، باید از بین متغیرهای به گونه ای انتخاب انجام شود که در تکرار فعلی بیشترین پیشرفت در بهبود تابع هدف ایجاد گردد. در [۱۶] از روش Steepest Descent برای یافتن بهترین جهت نزولی با q عنصر غیر صفر استفاده می نماید. برای این منظور مسئله ی بهینه سازی زیر حل می شود:

$$\begin{aligned} \text{minimize} \quad & v(\mathbf{d}) = \mathbf{g}(\alpha^t)^T \mathbf{d} \\ (51) \end{aligned}$$

$$\text{subject to: } \mathbf{y}^T \mathbf{d} = 0$$

$$d_i \geq 0 \quad \text{for } i: \alpha_i = 0$$

$$d_i \leq 0 \quad \text{for } i: \alpha_i = C$$

$$(52)$$

$$-1 \leq \mathbf{d} \leq 1$$

$$|\{d_i: d_i \neq 0\}| = q$$

$$(53)$$

همانطور که قبلاً اشاره شد، $\mathbf{g}(\alpha)$ بردار مشتق های جزئی در α است. مسئله ی بیان شده نشان می دهد که جهت نزولی در نقطه ی α^t مورد نظر می باشد. این روش مجموعه ی کاری را با بیشترین شیب نزولی انتخاب می نماید.

در [۲۲] در خصوص همگرایی این روش بحث شده و آمده است که تعداد اعضای مجموعه ی کاری هر عدد زوجی می تواند باشد.

در این میان برخی روش های دیگر نیز وجود دارند که با ایجاد تغییراتی در روش های ذکر شده ایجاد شده اند. بعنوان مثال، در [۲۱] روشی برای انتخاب مجموعه ی کاری ارائه شده است که مبتنی بر روش Conjugate Gradient می باشد. همچنین در [۲۱] نشان داده شده است که روش تجزیه با توجه به مجموعه ی کاری انتخاب شده، همیشه در پس از تعداد محدودی تکرار پایان می یابد و جواب بهینه را پیدا خواهد نمود. تفاوت این روش با روش SVM^{light} در روش انتخاب متغیرهای مجموعه ی کاری می باشد. در SVM^{light} از با استفاده از روش Steepest Decent عناصر مجموعه ی کاری انتخاب می شوند ولی در این روش از Conjugate Gradient استفاده می نماید. روش GC روشی کارا در حل مسائل

درجه‌ی دوم نامعید است. در این روش فرض می‌شود که $h(\alpha)$ تابع هدف باشد، بناب این جهت $d(k)$ در تکرار k -ام خواهد بود:

$$d(k) = -\nabla h(\alpha(k)) + \beta(k)d(k-1), \quad k = 0, 1, 2, \dots$$

(۵۴)

$$\beta(0) = 0, \quad \beta(k) = \frac{\|\nabla h(\alpha(k))\|^2}{\|\nabla h(\alpha(k-1))\|^2}, \quad k = 1, 2, \dots$$

(۵۵)

در [۲۱] نشان داده شده است که این الگوریتم در تعداد محدودی تکرار متوقف خواهد شد.

۳-۴- Chunking

Vapnik روشی را برای حل QP SVM مطرح نمود که به الگوریتم *chunking* معروف می‌باشد و اولین الگوریتم ارائه شده از این دسته می‌باشد. این روش براساس این واقعیت استوار است که مقدار در درجه‌ی دوم، با حالتی که از حذف سطرها و ستون‌های ماتریس متناظر با ضرائب لاگرانژ برابر صفر می‌باشند، حاصل می‌شود، معادل است. [4] بنابراین مسئله‌ی QP اصلی می‌تواند به یک سری از مسائل QP با اندازه‌ی کوچکتر شکسته شود، که هدف نهایی آن شناسایی تمام ضرائب لاگرانژ غیر صفر و حذف ضرائب لاگرانژ صفر می‌باشد. این روش سبب می‌شود که سائز ماتریس از توان دوم تعداد نمونه‌های آموزشی به حدود توان دوم تعداد ضرائب لاگرانژ غیر صفر، کاهش یابد، ولی باز هم این روش نمی‌تواند برای مقادیر داده‌ی بسیار زیاد مناسب باشد، زیرا گنجاندن این ماتریس کاهش یافته نیز در حافظه کار مشکلی می‌باشد. این روش نیاز به بهینه کردن تمام ضرائب لاگرانژ غیر صفر دارد، در عین حال ممکن است ماتریس Kernel ساخته شده، به اندازه‌ای بزرگ باشد که در حافظه گنجانده نشود.

۴-۴- SMO(Sequential Minimal Optimization)

Sequential minimal optimization (SMO) الگوریتمی برای حل برنامه‌ریزی درجه‌ی دوم منتج شده در ضمن فرآیند آموزش ماشین بردارهای پشتیبان (SVM)، می‌باشد. این روش توسط John Platt (از Microsoft Research) در ۱۹۹۸ ابداع گردید. SMO به طور وسیعی برای تعلیم ماشین بردارهای پشتیبان، بکار گرفته شده و توسط بسته‌ی بسیار محبوب LIBSVM پیاده‌سازی شده است. انتشار SMO در ۱۹۹۸ سبب موج استقبال وسیعی توسط انجمن SVM گردید، زیرا روش‌های موجود برای حل SVM بسیار پرهزینه و پیچیده بوده‌اند. [13]

در حقیقت SMO، الگوریتم ساده ایست که به سرعت مسئله ی QP SVM را بدون هیچگونه ذخیره سازی ماتریس و بدون استفاده از راه حل عددی بهینه سازی برنامه ریزی درجه ی دوم، حل می کند. SMO با استفاده از تئوری Osuna اقدام به تجزیه ی مسئله ی اصلی به تعدادی زیر مسئله می نماید و به این ترتیب همگرایی روش را نیز تضمین می نماید. برخلاف سایر روش های پیشین، SMO در هر مرحله کوچکترین مسئله را برای بهینه سازی انتخاب می نماید. برای مسئله ی QP، SVM استاندارد کوچکترین مسئله ی بهینه سازی ممکن شامل ۲ ضریب لاگرانژ می باشد، زیرا ضرائب لاگرانژ باید از قیود خطی تساوی تبعیت نمایند. در هر مرحله، SMO دو ضریب لاگرانژ را برای بهینه کردن انتخاب می نماید و مقدار بهینه را برای این ضرائب پیدا نموده و SVM را برای یافتن مقدار بهینه به روز می نماید.

مزیت روش SMO در این است که راه حل مسئله برای دو ضریب لاگرانژ می تواند به صورت تحلیلی حاصل شود. بنابراین در سراسر راه حل از بهینه سازی QP به صورت عددی پرهیز می گردد. بعلاوه همانطور که گفته شد این روش نیازی به حافظه جهت ذخیره سازی ماتریس ندارد، بنابراین مجموعه های آموزشی بسیار بزرگ را نیز می توان توسط این روش مدیریت نمود.

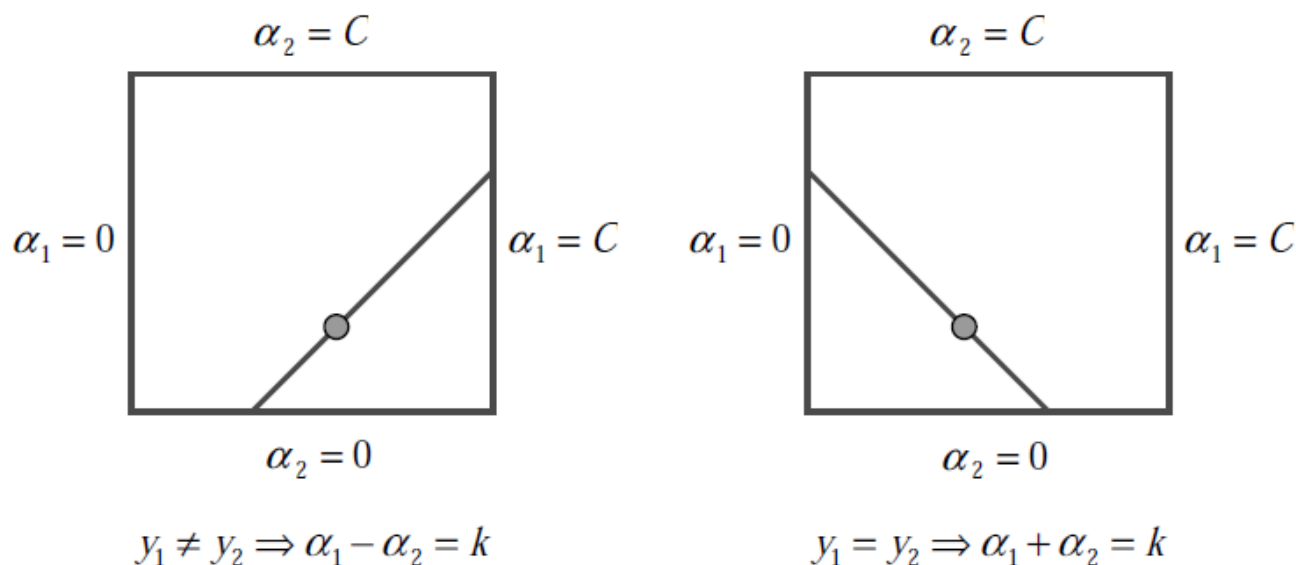
الگوریتم SMO در [4] بری آموزش ماشین های بردار پشتیبان، به طور کامل ارائه شده است. آموزش یک ماشین بردار پشتیبان مستلزم حل یک معادله ی بهینه سازی از نوع برنامه ریزی درجه ی دوم بسیار بزرگ می باشد. میزان فضایی مورد نیاز برای SMO نسبت به اندازه ی مجموعه ی آموزش خطی می باشد که این خود به SMO اجازه ی مدیریت مجموعه های آموزشی بزرگ را می دهد. با توجه به اینکه در این روش از محاسبات ماتریسی اجتناب شده است، لذا SMO بین مرتبه های خطی و دوم نسبت به اندازه ی مجموعه ی آموزشی قرار می گیرد، در حالیکه در SVM استاندارد در میان مرتبه های دوم و سوم قرار می گیرد. در دنیای واقعی SMO می تواند بیش از ۱۰۰۰ بار از الگوریتم قطعه بندی استاندارد برای SVM سریعتر باشد. [4]

در روش SMO [4]، در هر تکرار دو متغیر انتخاب شده و مسئله به صورت تحلیلی حل می شود، بنابراین دو عنصر در SMO وجود دارد [۱۷]، یک روش آنالیزی برای حل مسئله با دو ضریب لاگرانژ، و یک روال اکتشافی برای انتخاب اینکه کدام ضرائب برای بهینه سازی انتخاب شوند.

۴-۱- حل مسئله برای دو ضریب لاگرانژ

برای محاسبه ی این مسئله به ازای دو ضریب لاگرانژ، ابتدا باید محدودیت های روی ضرائب را در نظر گرفت و سپس به محاسبه ی می نیمم مقید شده پرداخت. با توجه به اینکه در این حالت فقط دو متغیر وجود دارد، می توان آنها را در دو بعد نمایش داد. قیدها نامساوی در مسئله QP SVM سبب قرار گرفت ضرائب در یک چهارچوب و قید مساوی مسئله سبب واقع شدن آنها روی یک خط (قطری از چهارچوب) می گردد (شکل ۱۲). بنابراین می نیمم مقید، برای تابع هدف،

روی این قطر واقع خواهد شد. قیدهای مسئله علت این نکته را که چرا دو ضریب لاگرانژ، حداقل تعداد ضریب برای بهینه سازی می باشد: انتخاب تنها یک ضریب، سبب عدم برآورده شدن قیدهای مسئله می شود. انتهای قطعه خط قطری را برای دومین ضریب (بدون از دست دادن کلیت مسئله) برحسب ضریب α_2 ، نوشته می شود.



شکل ۱۲ - دو ضریب لاگرانژ باید تمامی قیدهای مطرح شده در مسئله را مرتفع نمایند. قید نامساوی در مسئله سبب می شود که ضرایب در چهارچوب محدود شده و شرط تساوی سبب واقع شدن آنها روی خط قطری می گردد. بنابراین SMO باید مقدار بهینه را روی این قطعه خط قطری بدست آورد.

برای نوشتن حدود خط ها خواهیم داشت:

$$y_1 \neq y_2: \quad L = \max(0, \alpha_2 - \alpha_1), \quad H = \min(C, C + \alpha_2 - \alpha_1) \quad (56)$$

$$y_1 = y_2: \quad L = \max(0, \alpha_2 + \alpha_1 - C), \quad H = \min(C, \alpha_2 + \alpha_1) \quad (57)$$

همانطور که در [25] نشان داده شده است پس از جایگذاری و گرفتن مشتق دوم از دوگان مسئله ی بهینه سازی خواهیم داشت:

$$\frac{d^2 L_D}{d\alpha^2} = K(x_1, x_1) + K(x_2, x_2) - 2K(x_1, x_2) \quad (58)$$

که تعریف می کنیم:

$$\eta = K(x_1, x_1) + K(x_2, x_2) - 2K(x_1, x_2) \quad (59)$$

که مثبت یا منفی بودن η بستگی به تابع Kernel مورد استفاده دارد. در اینجا می خواهیم با داشتن مقادیر قبلی برای α_1 و α_2 و سایر α_i ها، مقدار جدید α_2 را محاسبه نماییم. با تعریف u به شکل زیر :

$$u = \sum_{j=1}^N y_j \alpha_j K(x_j, x) - b \quad (60)$$

و تعریف E_i به عنوان خطا در نمونه ی i -ام آموزشی:

$$E_i = u_i - y_i \quad (61)$$

با جایگذاری در مشتق مرتبه اول خواهیم داشت:

$$\frac{dL_D}{d\alpha_2} = \eta \alpha_2 + (y_2(E_1 - E_2) - \eta \alpha_2) \quad (62)$$

از مساوی صفر قرار دادن رابطه ی فوق داریم:

$$\alpha_2^{new} = \alpha_2 + \frac{y_2(E_1 - E_2)}{\eta} \quad (63)$$

مقدار بهینه ی محاسبه شده برای α_2 باید توسط قیدهای مسئله محدود گردد، لذا خواهیم داشت:

$$\alpha_2^{new, clipped} = \begin{cases} H & \text{if } \alpha_2^{new} \geq H; \\ \alpha_2^{new} & \text{if } L < \alpha_2^{new} < H \\ L & \text{if } \alpha_2^{new} \leq L \end{cases} \quad (64)$$

با تعریف $\alpha_1, s=y_1y_2$ جدید با توجه به α_2 جدید محاسبه می شود:

$$\alpha_1^{new} = \alpha_1 + s(\alpha_2 - \alpha_2^{new,clipped})$$

(۶۵)

۴-۲-۴- روش اکتشافی برای انتخاب ضرایب لاگرانژ

برای افزایش سرعت همگرایی، در SMO از روش اکتشافی برای انتخاب اینکه کدام ضرائب لاگرانژ جهت بهینه سازی انتخاب شوند، استفاده می کند. برای این کار از دو انتخاب اکتشافی مجزا برای ضرائب استفاده می شود. در الگوریتم ارائه شده [4]، از دو حلقه ی متداخل استفاده شده است. در حلقه ی خارجی در اولین اجرا، تمام مجموعه ی داده های آموزشی را برای یافتن نمونه ای که شرایط KKT را نقض نماید، پویش می نماید. اگر نمونه ای شرایط KKT را نقض نماید، مناسب برای بهینه سازی خواهد بود. در [26] این موضوع به شکل واضح تری اشاره شده است و بیان نموده که، حلقه خارجی به دنبال ضربی می گردد که شرایط زیر را نقض نماید:

$$\alpha_i = 0 \Rightarrow y_i f(x_i) \geq 1$$

(۶۶)

$$0 < \alpha_i < C \Rightarrow y_i f(x_i) = 1$$

(۶۷)

$$\alpha_i = C \Rightarrow y_i f(x) \leq 1$$

(۶۸)

در واقع این شرایط به همان شرایط KKT اشاره دارند. بعد از اولین تکرار حلقه ی خارجی تکرار را بر روی نمونه هایی از ضرایب لاگرانژ که بین 0 و C قرار دارند ($0 < \alpha_i < C$)، نباشند، مثلاً $0 < \alpha_i < C$)، ادامه می دهد. مجدداً، هر نمونه با شرایط KKT مقایسه شده و نمونه ای که آنرا نقض نماید، مناسب برای بهینه سازی خواهد بود. این بهینه سازی در حلقه ی خارجی تا زمانی اجرا می شود که تمام این نمونه های مذکور شرایط KKT را با خطای ε برآورده سازند. در ادامه حلقه خارجی بازگشت نموده و تکرار را روی تمام مجموعه ی آموزشی انجام می دهد. حلقه خارجی به طور مداوم تا زمانی که تمامی نمونه های آزمایشی شرایط KKT را با خطای ε برآورده سازند، بین این دو انتخاب این جابجا می شود.

بعد از انتخاب اولین ضریب توسط حلقه ی خارجی، دومین ضریب توسط حلقه داخلی و با تکرار بر روی ضرائب باقیمانده، برای انجام بهینه سازی انتخاب می شود. ارزیابی تابع $Kernel$ زمانبر می باشد، لذا SMO از قدرمطلق اندازه ی

گام $(E_1 - E_2)$ برای تخمین استفاده می‌نماید. SMO از یک $cache$ برای نگهداری مقدار خطای نمونه‌های غیرمرزی ($0 < \alpha_i < C$) استفاده نموده و سپس از یک خطا برای تخمین اندازه‌ی گام ماکزیمم، استفاده می‌نماید. در صورتی که E_1 مثبت باشد، نمونه‌ای با حداقل خطا E_2 انتخاب می‌گردد. در صورتی که E_1 منفی باشد، نمونه‌ای با حداکثر خطا E_2 انتخاب می‌گردد.

در الگوریتم ارائه شده در [4]، یک مرحله‌ی به‌روزرسانی برای b در نظر گرفته شده‌است که در [27] روشی بهینه برای الگوریتم SMO آمده است که طی آن نیازی با به‌روزرسانی b نمی‌باشد.

الگوریتم ارائه شده برای SMO در [4] به صورت شبه کد آمده است. در [26] نیز این الگوریتم به صورت ساده‌تری ذکر شده است.

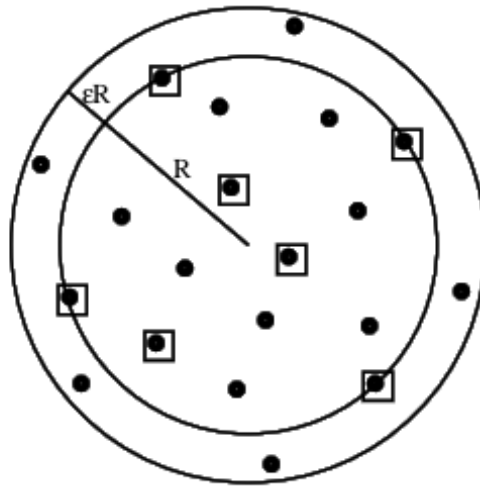
Core Vector Machines-۵-۴

در [2] روش CVM به عنوان روشی سریع برای آموزش SVM با داده‌های با حجم بسیار بالا مطرح شده است. SVM استاندارد دارای زمان آموزشی از مرتبه‌ی زمانی $O(m^3)$ و مرتبه پیچیدگی فضایی $O(m^2)$ می‌باشد. (m اندازه‌ی مجموعه‌ی آموزش می‌باشد). با مشاهده‌ی اینکه پیاده‌سازی‌های تجربی SVM فقط تخمینی از راه حل بهینه را با روش‌های تکرار ارائه می‌دهند، در [2] روش‌های kernel با توجه به چنین تخمین‌هایی، گسترش داده شده است. ابتدا نشان داده شده است که هر متد Kernel به طور مشابه می‌تواند با مسائلی MEB در هندسه‌ی محاسباتی، فرموله شود. سپس با انطباق یک الگوریتم موثر در تخمین MEB، مقدار تخمینی بهینه با توجه به ایده‌ی core set به دست می‌آید. روش Core Vector Machine ارائه شده در [2]، می‌تواند با kernel‌های غیر خطی بکار رود و دارای پیچیدگی زمانی خطی از m و پیچیدگی فضایی مستقل از m می‌باشد. نتایج تجربی نشان داده است که در مقادیر زیاد داده‌ها، دقت CVM به خوبی SVM می‌باشد و در عین حال می‌تواند با سرعت بالاتر حجم زیادی از داده‌ها را پاسخ‌گو باشد.

نکته کلیدی در پیاده‌سازی‌های تجربی SVM این است که مانند بسیاری از روش‌های عددی، فقط تقریبی از راه حل بهینه با استفاده از روش‌های تکراری حاصل می‌شود. این روش‌های تقریبی در بسیاری از مسائل محاسباتی مشکل، در کامپیوتر کاربرد داشته و به طور وسیعی مورد استفاده قرار می‌گیرند. در این الگوریتم از مفهوم minimum enclosing ball (MEB) برای تخمین استفاده می‌شود. مسئله‌ی MEB کوچکترین شعاعی را که می‌تواند نقاط یک مجموعه را پوشش دهد، محاسبه می‌کند. Clarkson و Bladoiu نشان دادند که تقریبی از MEB می‌تواند با استفاده از Core Sets به دست آید. Core Set زیر مجموعه‌ای از نقاط ورودی است که می‌تواند تقریب خوبی از جواب بهینه را از حل مسئله‌ی بهینه‌سازی روی این مجموعه حاصل نماید در واقع Core Set در CVM نقشی مشابه با Working Set در الگوریتم Decomposition را بازی می‌کند.

برای محاسبه ی تقریبی از MEB به صورت زیر عمل می شود

برای مجموعه ی مفروض $S = \{x_1, x_2, \dots, x_m\}$ که x_i متعلق به R^d می باشد، minimum enclosing ball MEB(S) نشان داده می شود، عبارتست از کوچکترین حجمی که بتواند تمام نقاط S را شامل شود. اگر فرض شود که $B(c, R)$ حجمی به مرکز c و شعاع R باشد، آنگاه برای $\epsilon > 0$ ، $B(c, (1 + \epsilon)R)$ تخمینی از MEB(S) خواهد بود اگر $R \leq r_{MEB(S)}$ بوده و $S \subset B(c, (1 + \epsilon)R)$ باشد. (شکل)



شکل ۱۳- مسئله ی MEB

الگوریتم ارائه شده توسط CVM برای بهینه سازی و یافتن بهترین جواب تقریبی برای مسئله ی بهینه سازی، با بدست آوردن مجموعه ای از بردارها با نام Core Vectorها پیاده سازی می شود. در این الگوریتم مرکز حجم با در تکرار t -ام برای مجموعه ی S_t با c_t و شعاع آن با R_t نمایش داده می شود. برای $\epsilon > 0$ ، CVM به شکل زیر عمل می کند:

1. **Initialize** S_0 , c_0 and R_0 .
2. **Terminate** if there is no training point z such that $\phi(z)$ falls outside the $(1 + \epsilon)$ -ball $B(c_t, (1 + \epsilon)R_t)$.
3. **Find** z such that $\phi(z)$ is furthest away from c_t . Set $S_{t+1} = S_t \cup \{z\}$.
4. **Find** the new MEB(S_{t+1}) from (5) and set $c_{t+1} = c_{MEB(S_{t+1})}$ and $R_{t+1} = r_{MEB(S_{t+1})}$ using (3).
5. **Increment** t by 1 and go back to Step 2.

$\phi(z)$ تابع نگاشت می باشد. نقاطی که به مجموعه ی Core Set اضافه می شوند، Core Vector گفته می شود. جزئیات مراحل فوق در [۱۵] آمده است. الگوریتم ارائه شده برای CVM برخلاف سادگی اش، ضمن گارانتی جواب تقریبی، دارای زمان کوتاه و پیچیدگی فضایی کم نیز می باشد.

۶-۴-۶ (Lagrangian SVM) LSVM

الگوریتم LSVM (Fung and Mangasarian, 2003; Mangasarian and Musicant 2001)، راه حل را با تکرارهای سریع محاسبه می نماید، ولی این روش نیز در مواردی با Kernel غیرخطی، نیاز به محاسبه ی ماتریس معکوس با ابعاد $m \times m$ می باشد. [2]

۷-۴-۶ Reduced SVM

Lee و Mangasarian در ۲۰۰۱ اقدام به ابداع روش Reduced SVM یا RSVM نمودند. این روش از یک زیرمجموعه ی مستطیل تصادفی از ماتریس Kernel استفاده می نماید. [2] RSVM یکی از جایگزین های SVM استاندارد می باشد. انگیزه ی پیدایش این روش مشکلات SVM استاندارد در مدیریت مجموعه های بزرگ داده ها با SVM با Kernel غیرخطی می باشد. این روش زیر مجموعه ای از داده ها را بعنوان بردارهای پشتیبان انتخاب نموده و مسئله ی بهینه سازی کوچکتری را حل می کند. در [1] نشان داده شده است که RSVM نیز به شکلی فرمول SVM خطی می باشد و همچنین ۴ پیاده سازی از RSVM مورد بحث قرار گرفته است. نتایج تجربی نشان می دهد که دقت RSVM اندکی از SVM استاندارد کمتر است ولی در مسائلی با حجم ۱۰ هزار داده و یا با تعداد زیاد بردار پشتیبان، RSVM می تواند مفید باشد. در [1] دو مقایسه صورت گرفته است:

۱- پیاده سازی های مختلف برای SVM خطی

۲- SVM استاندارد با استفاده از تابع هزینه ی درجه ی دوم و خطی.

در مواجهه شدن با مجموعه ای بزرگ از داده ها، RSVM بعنوان راه حلی برای کاهش پیچیدگی محاسباتی و کاهش پیچیدگی های مدل، ارائه شده است. در [3] RSVM از دیدگاه نمونه برداری و آماری مورد مطالعه قرار گرفته است.

در پایان این بخش به Relevance Vector Machine به عنوان توسعه ای از SVM و نه روشی برای حل معادله ی QP و سپس به مسئله ی انتخاب مجموعه ی کاری، اشاره می گردد.

۸-۴-۶ Relevance Vector Machine

اساساً روش SVM جهت بدست آوردن کارایی بهتر در کاربردهای عملی طراحی شده است. ولی این روش مشکلاتی نیز دارد. در [6] از Non-Bayesian بودن این روش بعنوان بزرگترین مشکل آن اشاره شده است. در SVM به ازای

ورودی جدید خروجی مشخصی وجود دارد. مزایای زیادی در دانستن احتمال وقوع اعضای یک کلاس وجود دارد. روش RVM (Relevance Vector Machine) که توسط Tipping (از Microsoft) ارائه شده، به صورت احتمالاتی پیش بینی را انجام می دهد و همچنان کارائی بالا و خلوت بودن SVM را حفظ می نماید. [6]

RVM روشی است در یادگیری ماشین که با استفاده از استنباط Bayesian راه حل های کم هزینه ای را برای رگرسیون و طبقه بندی احتمالاتی (Probabilistic Classification)، ارائه می دهد. RVM فرم تابعی مشابه با SVM دارد ولی طبقه بندی احتمالاتی را ارائه می کند. [11]

۹-۴- انتخاب مجموعه ی کاری^{۱۲}

در [10] به مسئله ی انتخاب مجموعه ی کاری به عنوان یک مسئله ی مهم در روش های مبتنی بر تجزیه^{۱۳}، پرداخته شده است. در این مقاله روشی برای انتخاب مجموعه ی کاری در روش تجزیه ی مبتنی بر SMO ارائه شده است. در بسیاری از مسائل برای سادگی نمایش، مسئله ی بهینه سازی مطرح شده در آموزش SVM به شکل زیر بیان می شود:

$$\min f(\alpha) = \frac{1}{2} \alpha^T Q \alpha - e^T \alpha$$

(۶۹)

$$\text{subject to: } 0 \leq \alpha_i \leq C, i = 1, \dots, l, y^T \alpha = 0$$

(۷۰)

که در آن e آرائه ای از ۱ها، C کران بالای متغیرها، Q یک ماتریس $l \times l$ متقارن که در آن $Q_{ij} = y_i y_j K(x_i, x_j)$ و $K(x_i, x_j)$ عبارتست از تابع هسته ی مورد استفاده. ماتریس Q عموماً بسیار پر (خلوت نمی باشد) بوده و بسیار بزرگ می باشد، به طوری که در حافظه قرار گرفتن آن مشکل می باشد. روش های تجزیه، مانند روش های ذکر شده در این بخش اساساً در پی مدیریت این پیچیدگی می باشند. برخلاف سایر روش های بهینه سازی که در آنها تمام بردار α به صورت یکجا و در هر تکرار، به روز می شود، روش های مبتنی بر تجزیه فقط زیر مجموعه ای از α را در هر تکرار به روز رسانی می نمایند. این زیر مجموعه به مجموعه ی کاری معروف می باشد و با B نمایش داده می شود. به این ترتیب به روز رسانی B در هر تکرار سبب می شود که یک زیر مسئله ی کوچک در هر تکرار می نمایم شود. همانطور که پیشتر اشاره شد، اوج این روش در SMO می باشد که در هر مرحله B ، تنها دارای ۲ عنصر می باشد. بنابراین در هر مرحله این مسئله ی دو متغیره به سادگی حل می شود. الگوریتم زیر روش تجزیه در SMO را نمایش می دهد:

^{۱۲} Working Set

^{۱۳} Decomposition Methods

۱- α^1 را به عنوان راه حل ممکن اولیه پیدا کن. $K=1$

۲- اگر α^k راه حلی بهینه برای (۴۱) می باشد متوقف شو، در غیر اینصورت مجموعه ی دو عنصری کاری $B=\{i,j\} \subset \{1,...,l\}$ را انتخاب کن. $N \equiv \{1,...,l\}/B$ و α^k_B و α^k_N را به عنوان زیربرداري از α^k ، مناظر با B و N تعريف كن.

۳- زیر مسئله ی زیر را با متغیر α_B حل کن:

$$\begin{aligned} \min_{\alpha_B} \frac{1}{2} [\alpha_B^T (\alpha_N^k)^T] \begin{bmatrix} Q_{BB} & Q_{BN} \\ Q_{NB} & Q_{NN} \end{bmatrix} \begin{bmatrix} \alpha_B \\ \alpha_N^k \end{bmatrix} - [e_B^T e_N^T] \begin{bmatrix} \alpha_B \\ \alpha_N^k \end{bmatrix} \\ = \frac{1}{2} \alpha_B^T Q_{BB} \alpha_B + (-e_B + Q_{BN} \alpha_N^k)^T \alpha_B + constant \\ = \frac{1}{2} [\alpha_i \ \alpha_j] \begin{bmatrix} Q_{ii} & Q_{ij} \\ Q_{ij} & Q_{jj} \end{bmatrix} \begin{bmatrix} \alpha_i \\ \alpha_j \end{bmatrix} + (-e_B + Q_{BN} \alpha_N^k)^T + constant \end{aligned} \quad (71)$$

$$subject\ to: 0 \leq \alpha_i, \alpha_j \leq C, y_i \alpha_i + y_j \alpha_j = -y_N^T \alpha_N^T$$

(72)

۴- α^{k+1}_B را به عنوان راه حل بهینه برای (71) قرار بده، $\alpha^{k+1}_N = \alpha^k_N$. $k=k+1$. برو به ۲.

باید توجه داشت که B در هر تکرار تغییر نمی نماید و در اینجا برای سادگی بجای B^k از B استفاده شده است. از آنجایی که در هر تکرار تنها تعداد کمی از عناصر به روز می شوند، لذا در مسائل بزرگ، روش تجزیه به کندی همگرا می شود. راه حل های مناسب تر می توانند از طریق کاهش تعداد مراحل تکرار، سریعتر عمل نمایند.

روش های موجود عموماً بر اساس نقض شرایط بهینه عمل می کنند که وابسته به دانستن اطلاعات مرتبه ی اول (مانند گرادیان)، از تابع هدف می باشند. مطالعات بر روی بهینه سازی نشان داده است که استفاده از اطلاعات مرتبه ی دوم، عموماً سبب سرعت بیشتر در همگرایی می شوند. از آنجایی که معادله ی (۶۹) یک مسئله ی درجه ی دوم می باشد، لذا اطلاعات مرتبه ی دوم مستقیماً به کاهش تابع هدف مربوط می شود. تلاش های بسیاری برای یافتن مجموعه ی کاری براساس کاهش تابع هدف انجام شده است. اما این روش ها اکتشافی بوده و تضمینی بر همگرایی آنها نمی باشد. در [10] روش ساده ای برای انتخاب مجموعه ی کاری با استفاده از اطلاعات مرتبه ی دوم ارائه شده است که سبب زمان تعلیم کوتاه تری می گردد.

روش های مختلفی برای انتخاب این مجموعه در [10] معرفی شده و مقایسه شده اند و نشان داده شده است که روش معرفی شده، از کارایی بیشتری نسبت به تمام آنها برخوردار می باشد.

۵- کاربردهای SVM

در این بخش مروری خواهیم داشت به کاربردهای تشخیص الگو^{۱۴} با استفاده از SVM. با توجه به هدف کاربردها در [8] هفت دسته کاربرد مختلف برای SVM دسته بندی شده است که در اینجا به آنها اشاره می نماییم. تعدادی از کاربردها که در دسته بندی خاصی گنجانده نشده است در گروه سایر کاربردها، ذکر شده اند.

۵-۱- تشخیص و شناسایی چهره^{۱۵}

این دسته یکی از محبوب ترین زمینه های در بایومتریک، تایید هویت، کنترل دسترسی و نظارت ویدئویی می باشد. زمینه های تحقیقاتی فعالی در این دسته وجود دارند که برای این کاربردها از روش های مختلفی استفاده می نمایند. به هرحال رسیدن به کرائی مطلوب در این زمینه بسیار مشکل خواهد بود. شناسایی چهره های افرادی که ویژگی های ظاهری نزدیک به هم دارند بسیار مشکل خواهد بود. همچنین حالات مختلف برای یک چهره ی خاص نظیر آرایش های متفاوت و ... تشخیص را بسیار مشکل می نماید. همچنین استفاده از عینک و یا ریش و سیبیل می تواند تشخیص را مشکل و پیچیده نماید. در سال های اخیر تحقیقات زیادی برای استفاده از SVM در کاربردهای تشخیص چهره، شناسایی و بازشناسی چهره و تشخیص حالات چهره صورت گرفته است. هر کدام از روش های فوق از ورودی های، پایگاه داده ها و هسته های مختلفی برای طبقه بندی کننده ی SVM استفاده کرده اند. اولین بار Osuna از SVM برای تشخیص چهره استفاده نمود. در این کاربرد وی از یک تصویر چهره ی ۱۹*۱۹ و تابع هسته ی چند جمله ای درجه ی دوم استفاده نمود. شناسایی ماشین در مورد چهره به دو دسته ی شناسایی چهره و تایید چهره تقسیم می شود. Guo از SVM چندکلاسه و یک درخت باینری برای شناسایی چهره استفاده نمود. ویژگی های نرمال شده که توسط روش PCA استخراج شده بودند، به عنوان ورودی SVM بکار گرفته شدند.

۵-۲- تشخیص و شناسایی اشیاء

هدف از تشخیص یا شناسایی اشیا یافتن و تعقیب افراد در حال حرکت و یا بررسی ترافیک برای کاربردهای کنترل ترافیک می باشد. شناسایی اشیا ۳-بعدی نیز یکی از این زمینه ها می باشد. COIL یک پایگاه داده ی معروف می باشد که شامل ۷۲۰۰ تصویر از ۱۰۰ شیء و از ۷۲ زاویه دید مختلف می باشد. این پایگاه داده در بسیاری از کاربردهای تحقیقاتی در این زمینه مورد استفاده قرار گرفته است. Roobaert تشخیص اشیا ۳ بعدی را توسط SVM انجام داد و توانایی

^{۱۴} Pattern Recognition

^{۱۵} Face Detection and Recognition

SVM را در این شناسایی از زاویه های مختلف، بررسی کرد. M.Pontil و A.Verri با آمیختن تصاویر پایگاه داده ی COIL با نویز و شیف دادن تصاویر، کارائی بالایی را نتیجه گرفتند.

۳-۵- تشخیص دست نوشته و ارقام

در میان کاربردهای مختلفی که بر پایه ی SVM انجام شده است، تشخیص ارقام دست نویس توسط SVM از تمامی الگوریتم های یادگیری دیگر، کارائی بالاتری داشته است. مشکل عمده در مسئله ی تشخیص دست نوشته، وجود الگوهای مختلف و متنوع نوشتاری بسیار زیاد می باشد. مدل های Elastic که برپایه ی مشاهدات محلی و برنامه نویسی پویا شکل گرفته اند مانند HMM، برای تحلیل این تنوع کارا و مناسب هستند. Choisy برای ترکیب این توانایی محلی با ویژگی های سراسری، از NSPH-HMM برای مشاهدات محلی و نرمال سازی و از SVM برای مشاهدات سراسری بر روی خروجی نرمال شده توسط NSPH-HMM، استفاده نمود. کارهای مختلفی در این زمینه انجام شده است که در مجموع برتری این روش را به سایر روش ها در شناسایی دست نوشته ها نشان می دهد.

۴-۵- تشخیص صحبت و گوینده

در مسئله ی تشخیص صحبت و یا شناسایی گوینده، دو روش بسیار محبوب discriminative classifier و generative model classifier می باشند. روش هایی که از discriminative classifier استفاده می کنند، شامل درخت های تصمیم گیری، شبکه های عصبی و SVM می باشند. معروفترین روش generative model classifier شامل مدل Hidden Markov (HMM) و مدل ترکیبی Gaussian (GMM) می باشد. Bengio و Wan از SVM برای شناسایی گوینده با پایگاه داده های مختلف، استفاده نمودند. آنها بر روی داده های وابسته به متن و غیر وابسته به متن کار کرده و روش SVM را با آستانه گذاری کلاسیک برای تصمیم گیری در رد یا پذیرش تایید گوینده ی، جایگزین نمودند. استفاده های گوناگونی از SVM در این زمینه نیز صورت گرفته است.

۵-۵- بازیابی تصاویر و اطلاعات^{۱۶}

Content-based image retrieval نتیجه ی زمینه ی تحقیقات مهمی در زمینه ی کتابخانه های دیجیتال و پایگاه داده های چند رسانه ای می باشد. Guo معیاری جدیدی را تحت عنوان فاصله از مرز^{۱۷}، را برای بازیابی بافت تصاویر ابداع نمود که در آن مرز بین دسته ها توسط SVM بدست آمده است. برای بازیابی تصاویر مرتبط بیشتر با تصویر ورودی، طبقه بندی کننده ی SVM جدا کردن تصاویر به دودسته ی مرتبط و غیر مرتبط، استفاده شده است. Zhang و Drucker, Tian

^{۱۶} Information and Image Retrieval

^{۱۷} Distance-from-boundary

روش اتوماتیکی را با استفاده از SVM برای بازیابی تصاویر ابداع نمودند که در آن وزن‌ها با فاصله از ابر صفحه تعیین می‌شوند و داده‌ها به دو دسته‌ی داده‌های مثبت و منفی تقسیم می‌شوند.

۶-۵- پیش‌بینی

هدف اصلی در بسیاری از روش‌های پیش‌گویی غیرخطی، پیش‌بینی نقطه‌ی بعدی در یک سری زمانی^{۱۸} می‌باشد [7]. Tay و Cao روش C-ascending SVM را با کاهش C ارائه دادند. این ایده بر این فرض استوار است که بهتر است که وزن بیشتری را به داده‌های اخیر بدهیم تا به داده‌های دورتر. نتایج آنها نشان داد که روش C-ascending SVM کارایی بالاتری را نسبت به SVM استاندارد در پیش‌گویی سری‌های زمانی اقتصادی، ارائه می‌دهد. Fan روش SVM را بر مسئله‌ی پیش‌بینی مشکلات شرکت‌ها با توجه به وضعیت اقتصادی آنها، منطبق نمود. برای این مسئله انتخاب متغیرهای ورودی (شاخص‌های اقتصادی)، روی کارایی سیستم تاثیرگذار می‌باشد. در مقاله‌ی ارائه شده، پیشنهاد داده که از داده‌های مناسب استفاده شود که فاصله‌ی بین دسته‌های بیشترین مقدار و فاصله‌ی بین داده‌های مشابه را کمترین مقدار قرار می‌دهند.

۷-۵- سایر کاربردها

کاربردهای بسیار زیاد دیگری برای SVM در زمینه‌ی تشخیص الگو وجود دارد. بعنوان مثال Yang از SVM برای دسته‌بندی بصری بر اساس جنسیت، با استفاده از تصاویر کوچک (۱۲*۲۱) و با کیفیت پایین استفاده نمود. او برای این کار از پایگاه داده‌ی FERET استفاده نمود و از ۱۷۵۵ تصویر استفاده نمود. پس تعلیم نشان داده شد که کارایی این روش در مقابل روش‌های قبلی مانند FLD, RBF و ... بسیار بالاتر می‌باشد. (۳.۴٪ خطا)

Gutta از SVM برای دسته‌بندی حالات چهره بروی پایگاه داده‌ی FERET استفاده نمود و به ۱۰۰٪ دقت رسید. Yao پنج دسته اثر انگشت را فرض نمود و از SVM برای این دسته‌بندی استفاده نمود و کارایی خوبی را نتیجه گرفت. علاوه بر کاربردهای ذکر شده، بسیاری از کاربردهای دیگر مانند خلاصه کردن داده‌ها^{۱۹}، شناسایی هدف و... وجود دارد که SVM به خوبی از عهده‌ی آنها برمی‌آید.

در کاربرد خلاصه‌سازی داده‌ها، زیرمجموعه‌ای کوچک از پایگاه داده‌ای بزرگ انتخاب می‌شود و دقت تفکیک کننده که بر روی این داده‌های کاهش یافته عمل می‌کند، با نتیجه‌ی آموزش روی تمام داده‌ها، مقایسه می‌شود.

^{۱۸} Time Series

^{۱۹} Data Consentation

۶- ابزارهای پیاده سازی SVM

برای استفاده از طبقه بندی کننده ی SVM ابزارها و زیرساخت های مختلفی ایجاد شده است که هر کدام ویژگی ها و امکانات خاصی را ارائه می دهند. در این بخش به تعدادی از معروفترین این ابزارها پرداخته خواهد شد.

۶-۱-MATLAB

یکی از ابزارهای مهم در این زمینه دستورات تعبیه شده در نرم افزار MATLAB می باشد. در MATLAB علاوه بر اینکه راهنمای خوبی برای آشنایی با مفهوم ماشین های بردار پشتیبان، ارائه شده است دستوراتی نیز برای کار با SVM طراحی شده است.

مانند هر مدل یادگیری نظارت شده ای در MATLAB نیز ابتدا ماشین بردار پشتیبان توسط داده های تعلیمی آموزش داده شده، سپس از ماشین تعلیم دیده برای دسته بندی (پیش بینی) داده های جدید استفاده می گردد. بعلاوه برای رسیدن به دقت پیش بینی مناسب، می توان از توابع هسته ی مختلف استفاده نمود که پارامترهای این توابع هسته، قابل تنظیم می باشند. آموزش دسته بندی کننده ی SVM در MATLAB با استفاده از دستور svmtrain انجام می شود. الگوی این دستور به شکل کلی به این ترتیب می باشد:

```
SVMstruct=svmtrain(data,groups,'Kernel_Function','rbf');
```

که در آن:

data : ماتریسی از نقاط داده است. هر سطر یک مشاهده و هر ستون یک ویژگی می باشد.

groups : بردارهای ستونی که هر سطر متناظر با مقداری متناظر در سطری در data می باشد. groups می تواند تنها دو نوع از ورودی ها را شامل شود.

Kernel_Function : مقدار پیش فرض آن 'linear' داده ها را توسط ابرصفحه جدا می نماید. مقدار 'rbf' از تابع Gaussian radial base استفاده می نماید.

نتیجه ی این تابع در ساختار SVMstruct قرار می گیرد که شامل مقادیر بهینه برای الگوریتم SVM می باشد. یعنی در این مرحله ماشین SVM تعلیم دیده است و می توان از آن در پیش بینی مقادیر جدید استفاده نمود.

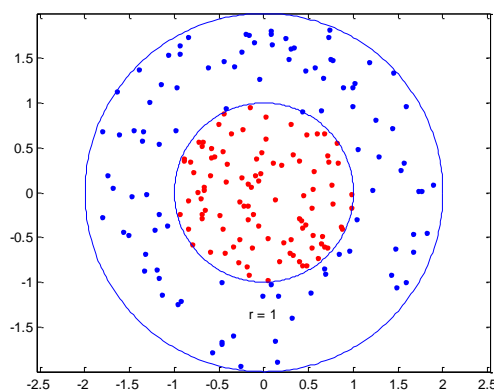
برای دسته بندی داده های جدید از تابع svmclassify استفاده می شود. از این تابع از ساختار SVMstruct که در دستور قبلی تکمیل شده است استفاده می نماید:

```
newClasses=svmclassify(SVMstruct,newData)
```

بردار خروجی در newClasses شامل نمایش دسته بندی برای داده ی جدید newData می باشد.

مثال:

در اینجا تعداد ۱۰۰ نقطه با توزیع نرمال به طور دایره ای تولید می شوند، سپس تعداد ۱۰۰ نقطه ی دیگر با توزیع نرمال به طور محاط بر داده های قبلی تولید می گردد. (شکل ۱۴)



شکل ۱۴- داده های آزمایشی تولید شده

در پیوست ۱ کد مربوط به تولید این داده ها آمده است.

با فرض داشتن داده های فوق، می خواهیم با استفاده از دستورات MATLAB ضمن دسته بندی داده ها بردارهای پشتیبان را نیز مشخص نماییم. برای این کار نیاز است که علاوه بر داده ها (x_i) خروجی متناظر با آنها نیز (y_i) نیز مشخص باشد، تا ماشین بتواند با خواندن داده های ورودی و خروجی متناظر با آنها یادگیری خود را تکمیل نماید.

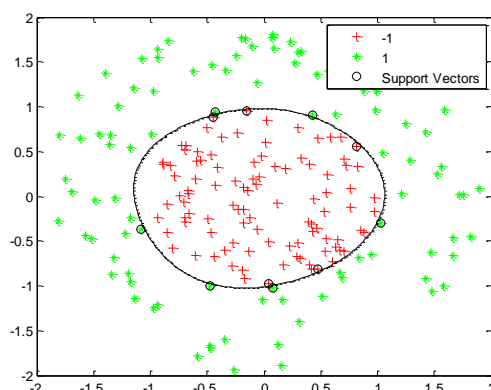
اگر فرض کنیم که داده ها دیسک داخلی در متغیر `data1` و دیسک خارجی در `data2` قرار گرفته باشند، آنگاه با دستورات زیر متغیر `data3` شامل هر دو مجموعه از داده ها شده و خروجی متناظر آنها نیز تولید می گردد که در دو دسته ی ۱- و ۱+ قرار گرفته اند:

```
data3 = [data1;data2];
theclass = ones(200,1);
theclass(1:100) = -1;
```

در مرحله ی بعد نیاز به آموزش ماشین می باشد:

```
cl = svmtrain(data3,theclass,'Kernel_Function','rbf',...
    'boxconstraint',Inf,'showplot',true);
```

که این دستور سبب می شود ضمن آنکه یادگیری انجام می شود، دسته بندی و بردارهای پشتیبان نیز نمایش داده شوند (شکل ۱۵):



شکل ۱۵- داده های دسته بندی شده با svm

در شکل ۱۵، بردارهای پشتیبان با دایره های کوچکی بر روی منحنی جداکننده مشخص شده اند.

سپس نوبت به تست بردار ماشین برای کاربرد دسته بندی می رسد. در این قسمت MATLAB از تابع `svmclassify` استفاده کرده و با ورود داده ی تست گروه متناظر با آنرا مشخص می نماید. بعنوان مثال اگر چند نقطه داخل دیسک مرکز و چند نقطه داخل دیسک محاطی انتخاب نماییم (که جز داده های ورودی نباشند)، خروجی های زیر را خواهیم داشت: نقاط داخل دیسک مرکزی:

```
s1=[-0.9,0.4];
svmclassify(cl,s1)
-1
s2=[0.3,0.9];
svmclassify(cl,s2)
-1
s3=[0.8,-0.2];
svmclassify(cl,s3)
-1
```

نقاط داخل دیسک محاطی:

```
sn1=[1.8,-1.2];
svmclassify(cl,sn1)
```

```

1
sn2=[-2.2,-2.2];
svmclassify(cl,sn2)

1
sn3=[-1.7,2.2];
svmclassify(cl,sn3)

1

```

لازم به ذکر است ۱- و ۱+ اشاره به گروه بندی تعریف شده در قسمت آموزش دارند.

LIBSVM - ۲-۶

بسته ی دیگری که برای کار با SVM ارائه شده است LIBSVM می باشد که به صورت مجموعه ای از توابع کتابخانه ای توسط Chih-Chung Chang and Chih-Jen Lin ارائه شده است. این مجموعه دارای امکانات فراوانی برای استفاده از SVM می باشد که برخی از آنها در زیر آمده اند:

- برای کاربردهای کلاسه بندی، رگرسیون و تخمین توزیع شده قابل استفاده می باشد.

- از دسته بندی چند کلاسه پشتیبانی می نماید.

- حمایت از تخمین احتمالاتی

- از توابع هسته ی گوناگون استفاده می نماید، ضمن آنکه ماتریس هسته در آنها از قبل محاسبه شده است.

- سورس آن به دو زبان C++ و Java موجود می باشد.

- در زبان ها و محیط های مختلف قابل استفاده می باشد.

(Python, R, MATLAB, Perl, Ruby, Weka, Common LISP, CLISP, Haskell, OCaml, LabVIEW, and PHP interfaces. C# .NET code and CUDA)

- از نسخه ی 2.8 به بعد از روش SMO استفاده می نماید.

همچنین چندین دیتاست مختلف توسط ارائه دهندگان LIBSVM طراحی و آماده ی استفاده می باشد.

آخرین نسخه ی این بسته تا زمان نگارش این مقاله، نسخه ی ۳.۱۸ می باشد که در آوریل ۲۰۱۴ منتشر شده است.

۶-۳- svm-struct-matlab

svm-struct-matlab عنوان بسته‌ای است که توسط Andrea Vedaldi برای استفاده‌ی ساده‌تر و کارآمدتر از SVM در MATLAB به صورت مجموعه‌ای از توابع پیاده‌سازی شده است که در MATLAB قابل فراخوانی و استفاده می‌باشد.

۷- جمع‌بندی و نتیجه‌گیری

در این مقاله مروری بر ماشین‌های بردار پشتیبان، حالت‌های مختلف دسته‌بندی داده‌ها، روش‌های حل مسئله‌ی بهینه‌سازی مرتبه‌ی دوم حاصل از SVM و کاربردهای مختلف SVM داشتیم. SVM در کاربردهای گوناگونی از شناسایی چهره تا تشخیص گفتار و ... را شامل می‌شود.

روش بردارهای پشتیبان یک روش بسیار کارا در دسته‌بندی داده‌ها می‌باشد از این رو یکی از ابزارهای مهم در یادگیری ماشین از نوع نظارت شده می‌باشد. این روش علاوه بر تمامی نقاط قوت خود محدودیت‌هایی را نیز دارا می‌باشد.

بزرگ‌ترین محدودیت روش بردارهای پشتیبان در انتخاب Kernel می‌باشد. زمانی که Kernel ثابت باشد، طبقه‌بندی SVM تنها دارای یک پارامتر قابل تغییر توسط کاربر می‌باشد (پارامتر جرمیه‌ی خطا). انتخاب بهترین Kernel برای یک مسئله‌ی خاص از مسائل مهم و مطرح می‌باشد. [5]

دومین محدودیت در سرعت و اندازه در آموزش و تست می‌باشد. آموزش پایگاه داده‌های بسیار بزرگ، مسئله‌ای غیرقابل حل به روش عددی می‌باشد. راه‌حلی برای این قبیل از مسائل ارائه گردیده که به برخی از آنها در مقاله اشاره گردید.

داده‌های گسسته مشکل دیگری را سبب می‌شوند، هرچند با مقیاس‌بندی‌های مجدد می‌توان به نتایج خوبی رسید. [5] یکی دیگر از مسائل، مشکل در طراحی تفکیک‌کننده برای SVM چند کلاسه می‌باشد که در اینجا نیز راه‌حلی ابداع گردیده است.

این روش به نسبت سایر روش‌های کلاسه‌بندی داده‌ها مانند Neural Network در صورتی که انتخاب‌ها مناسب انجام شوند (انتخاب حل‌کننده‌ی درجه‌ی دوم، انتخاب Kernel و ...) از دقت بسیار خوب و بالایی برخوردار می‌باشد و در کارهایی نظیر شناسایی الگو و شناسایی دست‌خط کارایی بالایی را به نمایش گذارده است.

پیوست ۱

کد مربوط به تولید شکل ۱۰ در MATLAB :

```
r = sqrt(rand(100,1)); % radius
t = 2*pi*rand(100,1); % angle
data1 = [r.*cos(t), r.*sin(t)]; % points

r2 = sqrt(3*rand(100,1)+1); % radius
t2 = 2*pi*rand(100,1); % angle
data2 = [r2.*cos(t2), r2.*sin(t2)]; % points

plot(data1(:,1),data1(:,2),'r.')
hold on
plot(data2(:,1),data2(:,2),'b.')
ezpolar(@(x)1);
ezpolar(@(x)2);
axis equal
hold off
```

پیوست ۲:

روش ضرائب لاگرانژ و شرایط KKT

روش ضرائب لاگرانژ

هنگامیکه با قیود به شکل معادله (تساوی) سر و کار داریم، این روش بهینه سازی بسیار کارا خواهد بود. به علاوه، هنگامیکه مقادیر ثابت قیود تغییر می کند، اطلاعات با ارزشی از مقدار بهینه تابع هدف در اختیار می گذارد.

ابتدا، تابع هدف F را در نظر بگیرید که بر زیرمجموعه ای چون X از فضای اقلیدسی تعریف شده است (پس عناصر X می توانند بردار باشند)؛ که این زیرمجموعه، توسط قیود به شکل $g(x)=b$ تعریف می شود. تابع لاگرانژ در این حالت عبارت خواهد بود از

$$L(x,y)=F(x)+y.(b-g(x))$$

(۷۳)

شرایط لازم را می توان، از طریق یافتن نقاط بحرانی تابع لاگرانژ (ماکزیمم سازی بدون قید)، به دست آورد.

مثال: فرض کنید مساله ی بهینه سازی زیر را داریم:

$$\max \{u(x_1, x_2): px_1 + px_2 = m\}$$

(۷۴)

لاگرانژین این مسئله به شکل زیر خواهد بود:

$$L(x_1, x_2, \lambda) := u(x_1, x_2) + \lambda(m - p_1 x_1 - p_2 x_2)$$

(۷۵)

شرایط کاروش-کان-تاکر (KKT)

در ریاضیات، شرایط کاروش-کان-تاکر، (Karush-Kuhn-Tucker) و یا کان-تاکر (Kuhn-Tucker)، شرایط لازم برای جواب بهینه در برنامه ریزی غیرخطی می باشد، مشروط بر آنکه برخی شرایط Regularity برقرار شود. کان-تاکر در واقع تعمیمی است بر روش ضرایب لاگرانژ که در آن قیود نامساوی هم در نظر گرفته شده است. مسئله بهینه سازی غیرخطی زیر را در نظر بگیرید:

$$\text{Minimize } f(x)$$

$$\text{subject to: } g_i(x) \leq 0, h_l(x) \leq 0$$

$$(۷۶)$$

تابع هدف $f(\cdot)$ تابعی است که قصد حداقل سازی آن را داریم؛ $g_i(\cdot)$ ($i=1, \dots, m$) توابعی هستند که قیود نامساوی را تشکیل می دهند و $h_l(\cdot)$ ($l=1, \dots, m$) توابعی اند که قیود تساوی را شکل می دهند، m و l به ترتیب تعداد قیود نامساوی و تساوی هستند. شرایط کان-تاکر، به افتخار هارولد و. کان (Harold W. Kuhn) و نیز آلبرت و. تاکر (Albert W. Tucker) نام گذاری شده است که اول بار، این شرایط را منتشر کرده اند. پژوهشگران بعدی، شرایط لازم این مسئله را یافتند. ویلیام کاروش (William Karush) در پایان نامه کارشناسی ارشد خود، این شرایط را تشریح کرده است.

- شرایط لازم

فرض که تابع هدف $f: \mathbb{R}^n \rightarrow \mathbb{R}$ باشد و توابع قیود نامساوی به شکل $g_i: \mathbb{R}^n \rightarrow \mathbb{R}$ و توابع قیود مساوی به شکل $h_l: \mathbb{R}^n \rightarrow \mathbb{R}$ باشند. اگر x^* ، مینیمم موضعی باشد که برخی شرایط Regularity را برآورده کند، آنگاه ثابتی چون u_i ($i=1, \dots, m$) و λ_j ($j=1, \dots, m$) وجود خواهد داشت بطوریکه:

$$\nabla f(x^*) + \sum_{i=1}^m \mu_i \nabla g_i(x^*) + \sum_{j=1}^l \lambda_j \nabla h_j(x^*) = 0$$

$$(۷۷)$$

$$g_i(x^*) \leq 0, \text{ for all } i = 1, \dots, m$$

$$(۷۸)$$

$$h_j(x^*) = 0, \text{ for all } j = 1, \dots, l$$

(۷۹)

$$\mu_i \geq 0, \text{ for all } i = 1, \dots, m$$

(۸۰)

$$\mu_i g_i(x^*) = 0, \text{ for all } i = 1, \dots, m$$

(۸۱)

- شرایط کافی

اگر مجموعه جواب های ممکن، مجموعه ای ناتهی و فشرده و ضمناً محدب باشد، و تابع هدف بر این مجموعه، پیوسته و مقعر باشد؛ در این صورت، ماکزیمم موضعی، ماکزیمم مطلق است. بعلاوه، اگر تابع هدف اکیدا مقعر باشد، در این صورت جواب یکتاست، یعنی یک ماکزیمم مطلق یکتا وجود دارد. پس یک راه اطمینان حاصل کردن از اینکه، جواب های حاصل از شرایط مرتبه اول کان-تاکر، ماکزیمم هستند، بررسی شرایط فوق الذکر است.

هرگاه تابع هدف f و قیود نامساوی g_i بطور پیوسته مشتق پذیر و محدب باشند و قیود تساوی h_i توابع آفین باشند؛ شرایط لازم، کافی نیز خواهند بود.

مارتین ۱۹۸۵ (Martin 1985)، دسته ای از توابع را مشخص کرد که در آن ها شرایط کان-تاکر، وجود بهینه سرتاسری را تضمین می کند، که توابع $invex$ خوانده می شوند. اگر قیود تساوی، توابع آفین باشند، قیود نامساوی و تابع هدف بطور پیوسته مشتق پذیر و $invex$ باشند، شرایط کان-تاکر برای بهینه سرتاسری، کافی نیز خواهد بود.

مراجع

- [1] K.M.Lin,C.J.Lin, " A Study on Reduced Support Vector Machines," IEEE TRANSACTIONS ON NEURAL NETWORKS, VOL. 14, NO. 6,pp.1449-1459,NOVEMBER 2003.
- [2] I.W. Tsang, J.T. Kwok, P.M.Cheung, "Core Vector Machines:Fast SVM Training on Very Large Data Sets," Journal of Machine Learning Research 6 (2005) 363–392,2005.
- [3] Y.J.Lee, S.Y.Huang, "Reduced Support Vector Machines: A Statistical Theory, " IEEE TRANSACTIONS ON NEURAL NETWORKS, 2006.
- [4] J.Platt, "Sequential minimal optimization: A fast algorithm for training support vector machines," 1998.
- [5] C.J.C.Burges, "A tutorial on support vector machines for pattern recognition," *Data mining and knowledge discovery* 2.2, 121-167, 1998.
- [6] C.M.Bishop, M.E. Tipping, "Variational relevance vector machines," *Proceedings of the Sixteenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 2000.
- [7]J. Quinonero-Candela, L.K.Hansen, "Time series prediction based on the Relevance Vector Machine with adaptive kernels," *Acoustics, Speech, and Signal Processing (ICASSP), 2002 IEEE International Conference on* , vol.1, no., pp.I-985,I-988, 13-17 May 2002.
- [8] H.Byun, S.W.Lee, "Applications of support vector machines for pattern recognition: A survey," *Pattern recognition with support vector machines*. Springer Berlin Heidelberg, pp. 213-236, 2002.
- [9] C.Campbell, "Kernel methods: a survey of current techniques," Elsevier, *Neurocomputing* 48.1, pp. 63-84, 2002.
- [10] R.E.Fan, C.Pai-Hsuen, C.J.Lin, "Working set selection using second order information for training support vector machines," *The Journal of Machine Learning Research* 6, pp.1889-1918, 2005.
- [11] M.Tipping, "Sparse Bayesian learning and the relevance vector machine," *The journal of machine learning research* 1, pp.211-244,2001.
- [12] "Hilbert space," http://en.wikipedia.org/wiki/Hilbert_space, [Online Access:Jol. 11, 2014]
- [13] " Sequential minimal optimization," http://en.wikipedia.org/wiki/Sequential_minimal_optimization, [Online Access: Jol. 11, 2014]

- [14] E. Osuna, R. Freund, and F. Girosi, "Training support vector machines: An application to face detection," in Proc. CVPR '97, 1997.
- [15] The Analysis of Decomposition Methods for Support Vector Machines
- [16] T.Joachims, "Making large-scale SVM learning practical," in Advances in Kernel Methods—Support Vector Learning, B. Schölkopf, C. J. C.Burges, and A. J. Smola, Eds. Cambridge, MA: MIT Press, 1998.
- [17] L.Kaufman, "Solving the quadratic programming problem arising in support vector classification," in Advances in Kernel Methods—Support Vector Learning, B. Schölkopf, C. J. C. Burges, and A. J. Smola, Eds. Cambridge, MA: MIT Press, 1999.
- [18] O.L.Mangasarian, D.R.Musicant, "Successive overrelaxation for support vector machines," IEEE Trans. Neural Networks, vol. 10, pp. 1032–1037, Sept. 1999.
- [19] J.C.Platt, "Fast training of support vector machines using sequential minimal optimization," in Advances in Kernel Methods—Support Vector Learning, B. Schölkopf, Ed. Cambridge, MA: MIT Press, 1998.
- [20] C.Saunders, M.O.Stitson, J.Weston, L.Bottou, B.Schölkopf, and A.Smola, "Support vector machine reference manual," Royal Holloway Coll., Univ. London, Egham, U.K., Tech. Rep. CSD-TR-98-03, 1998.
- [21] N.Takahashi, M.Kuranoshita, Y.Kawazoe, J.Guo, J.I.Takeuchi, "A New Working Set Selection for Decomposition-Type SVM Learning Algorithms."
- [22] Lin, Chih-Jen. "On the convergence of the decomposition method for support vector machines," Neural Networks, IEEE Transactions on 12.6,p.p.1288-1298, 2001.
- [23] Takahashi, Norikazu, N.Tetsuo, "Global convergence of decomposition learning methods for support vector machines," Neural Networks, IEEE Transactions on 17.6 ,p.p.1362-1369, 2006.
- [24] Joachims, Thorsten, "Introduction to support vector machines," 2003.
- [25] Vasant Honavar, "Sequential Minimal Optimization for svm," Iowa State University.
- [26] R.R.Padron, "A Roadmap to SVM Sequential Minimal Optimization for Classification, " School of Electrical Engineering and Computer Science University of Central Florida Orlando.
- [27] S.Keerthi, Sathiya,"Improvements to Platt's SMO algorithm for SVM classifier design," Neural Computation 13.3, p.p. 637-649, 2001.