



دانشگاه صنعتی امیرکبیر

(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد

گرایش مهندسی کامپیوتر-هوش مصنوعی

عنوان

دسته بندی کارا مبتنی بر رگرسیون تنک

نگارش

پردیس نورزاد

استاد راهنما

محمد رحمتی

اساتید مشاور

نصرالله مقدم چارکاری

محمد مهدی عبادزاده

تیرماه ۱۳۹۱

صفحه فرم ارزیابی و تصویب پایان نامه - فرم تأیید اعضاء کمیته دفاع

در این صفحه فرم دفاع یا تایید و تصویب پایان نامه موسوم به فرم کمیته دفاع - موجود در پرونده آموزشی - را قرار دهید.

اینجانب پردیس نورزاد متعهد می‌شوم که مطالب مندرج در این پایان نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آنها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است.

در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است. نقل مطالب با ذکر مآخذ بلامانع است.

نام و نام خانوادگی دانشجو

امضا

تقدیر و تشکر

از استاد راهنمای عزیزم، پروفسور محمد رحمتی، برای پذیرش بنده به عنوان دانشجوی خود در آزمایشگاه پروازش سیکنال و تشخیص الگو و فراهم سازی مکانی مناسب و مطلوب برای کارهای اینجانب، کمال تشکر و قدردانی را دارم. همچنین از وی بدلیل اختصاص زمان و بردباری در پرداختن به ایده های فراوانم در پایان نامه و ارائه جهت گیری های مناسب و اطمینان بخشی از قرار داشتن در مسیر صحیح در طول مطالعاتم در دانشگاه صنعتی امیرکبیر، پاسکزارم.

خود را حقیقتاً مدیون پروفسور باب. ال. استورم از دانشگاه آلبرک کپنهاگ می دانم. وی مرا با الگوریتم های تقریبی تنگ آشنا نمود، که انگیزه ای برای این پایان نامه شد. از ایشان به پاس الهام بخشی قابل توجهشان برای نگاه به رسیدن به اطمینان و مهارت لازم جهت دستیابی به آینده در تحقیقات، نهایت تشکر را دارم.

و اما قدردان برادرم پرللم، بستم، به جهت مشاوره فنی در فضایی، بیلبرت توابع، ضرب های داخلی و نرم ها، و برای همه ساعات مشاوره ای که به من اختصاص داد.

بزرگ ترین قدردانی من از پدر و مادرم، معلمان فوق العاده ام، برای حمایت، آرایش و مراقبت در سراسر خطرات زندگی من است.

چکیده

اخیراً مساله کمینه‌سازی زیان مربعی ℓ_1 -منظم، توجه بسیاری را به خود جلب کرده است. بویژه، صحبت زیادی از لسو می‌شود که اصل بهینه‌سازی به کار رفته برای حل مساله معکوس خطی $y = Xa$ تحت شرط تنک بودن بردار ضریب a است. این اصل بهینه‌سازی دارای دو کاربرد اصلی در یادگیری ماشین می‌باشد. کاربرد اول در رگرسیون است که در آن y بردار خروجی یا برچسب داده‌ی آموزش می‌باشد. زمانی که لسو برای کدگذاری تنک و یادگیری ویژگی، یا در زمینه دسته‌بندی نمایش تنک به کار گرفته شود، y بردار ویژگی یا خود سیگنال است.

قسمت اول این پایان‌نامه به نحوه‌ی انتخاب الگوریتم مناسب برای مسائل دسته‌بندی عظیم می‌پردازد. در این مسائل نه تنها سرعت آموزش مدل مهم می‌باشد، بلکه سرعت ارزیابی دسته‌بندی‌کننده و حافظه‌ی اشغال شده توسط آن نیز دارای اهمیت می‌باشد. در این قسمت به اهمیت استفاده از دسته‌بندی‌کننده‌ی خطی با زیان مربعی و تنظیم ℓ_1 وزن‌ها، که همان اصل بهینه‌سازی لسو می‌باشد، برای مسائل دسته‌بندی عظیم می‌پردازیم. لازم به ذکر است که استفاده از لسو برای رگرسیون رایج است ولی ما فواید استفاده از این اصل بهینه‌سازی را در مسائل دسته‌بندی ارائه می‌دهیم. برای این کار دلایل مبتنی بر نظریه‌های آمار ارائه داده و همچنین با چند آزمایش صحت این ادعا را به صورت شهودی نیز اثبات می‌کنیم.

در قسمت دوم پایان‌نامه اهمیت استفاده از تنظیم ℓ_1 بردار وزن را برای بازسازی سیگنال (برداری ویژگی) می‌پردازیم و دلایل آماری برای موفقیت این روش در این کاربرد ارائه می‌دهیم. به دنبال موفقیت تنظیم ℓ_1 برای مسئله‌ی بازسازی سیگنال، روشی برای دسته‌بندی مبتنی بر این کاربرد تنظیم ℓ_1 توسط محققان ارائه شده (بخش ۳،۲). ما این

روش را به مسئله‌ی رگرسیون تعمیم می‌دهیم. علاوه بر این، با ارائه آزمایش، کارکرد این روش را (هم برای دسته‌بندی و هم برای رگرسیون) با الگوریتم متداول k نزدیک‌ترین همسایه مقایسه می‌کنیم. نشان می‌دهیم که این روش با آنکه گران‌تر از روش ساده و متداول نزدیک‌ترین همسایه می‌باشد، کارکرد بیشتری فراهم نمی‌کند.

در قسمت سوم بر اساس ساده‌سازی اثبات‌های قبلی انجام شده‌ی قبلی نشان می‌دهیم که مسئله‌ی رگرسیون بردار پشتیبان (با زیان محوری) هم ارز مسئله‌ی نویززدایی تعاقب پایه (با زیان مربعی) است. دلیل اهمیت این اثبات این است که راهی برای توضیح موفقیت این دو روش و عملکرد مشابهشان در مسئله‌ی رگرسیون فراهم می‌کند.

واژه‌های کلیدی:

کمینه‌سازی زیان مربعی، ℓ_1 -منظم، دسته‌بندی دودویی، بازنمایی تنک

۱ فصل اول مقدمه	۱
۱,۱ دسته‌بندی خطی و SVM	۲
۱,۱,۱ رگرسیون لجستیکی	۷
۲,۱ تقریب تنک	۹
۱,۲,۱ نمایش تنک سیگنال‌ها و ویژگی‌ها	۱۰
۲ فصل دوم کمینه‌سازی زیان مربعی ℓ_1 -منظم برای دسته‌بندی	۱۱
۱,۱,۲ کارهای مرتبط در کمینه‌سازی زیان مربعی برای دسته‌بندی	۱۳
۲,۱,۲ کارهای مرتبط در ℓ_1 -منظم‌سازی برای دسته‌بندی‌کننده تنک	۱۴
۲,۲ دسته‌بندی و کمینه‌سازی زیان محدب	۱۷
۱,۲,۲ احتمال پسین و دسته‌بندی‌کننده‌ی جانشین	۱۹
۲,۲,۲ کمینه‌سازی ریسک آزمایشی	۲۲
۳,۲,۲ ارزیابی احتمال پسین	۲۶
۳,۲ چرا منظم‌سازی ℓ_1 بردار وزن تنک ایجاد می‌کند؟	۲۷
۴,۲ ارزیابی تجربی لسو برای دسته‌بندی	۲۸
۳ فصل سوم کمینه‌سازی زیان مربعی ℓ_1 -منظم برای بازسازی	۳۰
۱,۳ کدگذاری تنک و یادگیری فرهنگ لغات برای یادگیری ویژگی	۳۱
۲,۳ دسته‌بندی بازنمایی تنک	۳۶
۳,۳ SPARROW: رگرسیون وزن‌دار مبتنی بر تقریب تنک	۳۷
۱,۳,۳ رگرسیون وزن‌دار مبتنی بر تقریب تنک	۴۰
۲,۳,۳ تعریف وزن‌های موثر	۴۰
۳,۳,۳ تعریف وزن‌های شهودی	۴۲
۴,۳ ارزیابی آزمایشی SPARROW	۴۴
۱,۴,۳ نتیجه‌گیری	۴۸
۵,۳ مقایسه kNN با SRC	۴۹
۴ فصل چهارم یک هم‌ارزی بین ϵ -SVR و BPDN	۵۰
۱,۴ فضای هیلبرت هسته تکثیری	۵۱
۲,۴ ماشین‌های بردار پشتیبان برای رگرسیون	۵۴
۳,۴ ϵ -SVR و تنک بودن	۵۹
۴,۴ ارتباط با تقریب تنک	۵۹

۶۴.....	۵ فصل پنجم نتیجه گیری و کارهای آینده.....
۶۵.....	۱,۵ رگرسیون غیرخطی
۶۶.....	۲,۵ منظم سازی α : متغیر دوگان SVM
۶۶.....	۳,۵ ناپایداری و غیریکتا بودن جواب های لسو
۶۷.....	۴,۵ kNN یا تعاقب تطابقی برای یادگیری ویژگی
۷۲.....	منابع و مراجع.....

- شکل ۱.۱ نمودار، تعداد بردارهای پشتیبان را نشان می‌دهد، که با افزایش تعداد نمونه‌های آزمایشی رشد می‌کند. ۱۶
- شکل ۲.۱ نمودارها نشان می‌دهد که تعداد بردارهای پشتیبان ارتباط معنی‌داری با فرایارامتر C در بهینه‌سازی SVM ندارند. ۱۷
- شکل ۱.۲ در این اشکال، می‌بینیم که، تعداد عناصر ناصفر جواب، زمانی افزایش می‌یابد که، پارامتر منظم‌سازی λ افزایش پیدا کند، بنابراین، ابزاری را برای کنترل پراکندگی جواب در اختیار ما قرار می‌دهد. ۱۹
- شکل ۲.۲ مقایسه‌ای از توابع زبان محدب. زبان دسته‌بندی نادرست نیز نشان داده شده است. ۲۵
- شکل ۳.۲ همان‌طور که p به صفر می‌رود، $|x|^p$ به تابع شاخص تبدیل می‌شود و تعداد درآیه‌های ناصفر در x را می‌شمارد [۶]. ۲۷
- شکل ۱.۳ خطوط لوله دسته‌بندی تصویر [۱۲]. ۳۴
- شکل ۲.۳ این اشکال، توانایی روش‌های رگرسیون محلی را برای مدل‌سازی داده‌ها با یک توزیع ناشناخته، نشان می‌دهند. تابع مولد داده‌ها عبارت است از: $y_i = f(x_i) + \epsilon_i$ $y_i = f(x_i) + \epsilon_i$ که در آن $f(x) = (x^3 + x^2)I(x) + \sin(x)I(-x)$ ۳۹
- شکل ۳.۳ ترسیم جعبه‌ای برای ارزیابی اعتبار گذری ۱۰ تایی از خطای مربع میانگین (۱۰۰ بار اجرای مستقل) برای ۴ مجموعه داده گوناگون. هر جعبه ۲۵ تا ۷۵ درصد را معین می‌کند و خط قرمز، میانه را نشان می‌دهد. اکستریم با علامت + و برون هشته‌ها با ضرب‌آنها مشخص شده است. ۴۷

جدول ۱.۱ اطلاعات مجموعه داده ها: n نشان دهنده تعداد مشاهدات و p نشان دهنده تعداد گرایش های هر مشاهده است.	۱۵
جدول ۱.۲ توابع زیان محدب شناخته شده و تابع کمینه سازی متناظرشان.	۲۴
جدول ۲.۲ مقایسه ای از سه دسته بندی کننده دیگر با دسته بندی براساس لسو روی ۷ مجموعه داده. فرایارامتر با C یا λ بسته به الگوریتم نشان داده می شود. تعداد درآیه های ناصفر در بردار جواب نیز با $\#nz$ و درصد نمونه های آزمون بدرستی دسته بندی شده با Acc نمایش داده می شود.	۲۶
جدول ۳.۲ نتایج برای دسته بندی براساس رگرسیون مرزی. توجه داشته باشید که، تعداد درآیه های ناصفر جواب برابر با تعداد صفت های مشاهدات است.	۲۹
جدول ۱.۳ اطلاعات مجموعه داده ها. آخرین ستون پارامتر k تزار شده در آزمایش های مربوط به k -NNR و Wk -NNR را نشان می دهد.	۴۵
جدول ۲.۳ مقایسه ای از ارزیابی های MSE براساس ۴ مجموعه داده با ۱۰ توالی اعتبار گذری ۱۰ تایی از C -SPARROW و L -SPARROW با و بدون منظم سازی. آخرین ستون، نشان دهنده پارامتر مرزی بکار رفته برای دستیابی به ارزیابی L -SPARROW است.	۴۸
جدول ۳.۳ مقایسه دقت بدست آمده توسط kNN و SRC روی ۵ مجموعه داده دسته بندی چندکلاسه.	۴۹

فهرست علائم

علائم لاتین

بردار ضریب یا وزن	\mathbf{w}, \mathbf{a}
عرض از مبدا	b
فراپارامتر ضریب هزینه	C
کرنل	K
اندازه‌ی مجموعه‌ی داده‌ی آموزشی	n
بردار ویژگی	\mathbf{x}
فرهنگ لغات یا ماتریس داده یا ماتریس ویژگی	\mathbf{X}, \mathbf{D}
برجسب متناظر با یک بردار ویژگی	y

علائم یونانی

حاشیه خطا در دسته‌بندی‌کننده‌ی ماشین بردار پشتیبان	ξ
فراپارامتر ضریب زیان	λ

۱

فصل اول مقدمه

این پایان نامه دو کاربرد اصلی کمینه سازی زیان مربعی ℓ_1 -منظم را بررسی می کند. نخستین کاربرد، که در فصل ۲ پوشش داده شده است، بکارگیری مساله بهینه سازی مطرح شده برای دسته بندی دودویی خطی است. کاربرد دوم، که در فصل ۳ پوشش داده شده است، اعمال کمینه سازی زیان مربعی ℓ_1 -منظم برای نمایش سیگنال و ویژگی ها است. در فصل ۴، یک هم ارزی بین کمینه سازی زیان مربعی ℓ_1 -منظم و رگرسیون برداری ϵ -پشتیبان، تحت تعدادی از شرایط محدود سازی توضیح داده می شود. در فصل آخر، نتایج و بعضی زمینه ها برای کارهای آینده ارائه می شوند.

مسئله ی بهینه سازی لسو که در ۱,۲ معرفی می کنیم دارای سه خصیصه ی مهم است: اولاً خطی است، ثانیاً دارای زیان مربعی است و نهایتاً دارای تنظیم ℓ_1 روی وزن ها می باشد. در ادامه این فصل این سه موضوع را بررسی کرده و با اهمیت هر یک آشنا می شویم. در بخش ۱,۱ راجع به دسته بندی با خطای محوری (که یک حالت آن دسته بندی کننده ی بردار پشتیبان می باشد) صحبت می کنیم. بخش ۱,۲ به توضیح مسائل معکوس خطی، رگرسیون خطی (با خطای مربعی) و مفهوم تنظیم ℓ_1 و ℓ_2 می پردازیم. در بخش ۱,۲,۱ رگرسیون لجیستیکی را توضیح داده و با نحوه ی به دست آوردن احتمال واپسین توسط این روش آشنا می شویم. در نهایت در بخش ۱,۳ راجع به قسمت دوم پایان نامه که تقریب تنک سیگنال می باشد توضیح می دهیم.

۱.۱ دسته بندی خطی و SVM

دسته بندی ماشین بردار پشتیبان (SVM)^۱ [۴]، [۵۳] دارای دو ویژگی شناخته شده است: کارایی تعمیمی بسیار خوب نسبت به یک جریمه ℓ_2 روی وزن ها؛ و یک بردار ضریب پراکندگی نسبت به زیان محوری به عنوان جمله برازش داده ها. پس از آن، نشان می دهیم که به محض محاسبه خروجی برای یک نمونه آزمون مفروض، دسته بندی SVM، تنها زیرمجموعه ای از نمونه های آزمایشی، شناخته شده به عنوان بردارهای پشتیبان است که منجر به ارزیابی سریع تر می شود. در عمل، با این حال، تعداد بردارهای پشتیبان با تعداد نمونه های آزمایشی قابل مقایسه است: بردار ضریب پراکنده نیست. بنابراین، دسته بندی

¹ Support Vector Machine

حاصل هزینه بر است و سرعت با افزایش نمونه‌های آزمایشی، رشد می‌کند. مساله بهینه‌سازی SVM بشکل زیر است

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{Subject to } \begin{cases} y_i(w \cdot x_i - b) & \geq 1 - \xi_i \\ \xi_i & \geq 0 \end{cases} \quad \text{for } i = 1, \dots, n, \end{aligned} \quad (1.1)$$

این می‌تواند به صورت زیر جایگزین شود

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{Subject to } \xi_i \geq \max(0, 1 - y_i(w \cdot x_i - b)) \quad \text{for } i = 1, \dots, n. \end{aligned} \quad (2.1)$$

که در آن

$$\xi_i \geq [1 - y_i(w \cdot x_i - b)]_+, \text{ where } [x]_+ = \max(0, x) \quad (3.1)$$

و منجر به فرمول زیر برای مساله بهینه‌سازی SVM می‌شود

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n [1 - y_i(w \cdot x_i - b)]_+ \quad (4.1)$$

که شامل دو جمله به نام های زیان محوری و ℓ_2 - منظم از وزن ها است.

تنک بودن دسته‌بندی حاصل برای کاربردهای زمان واقعی نظیر تشخیص گفتار، مهم است. همچنین، برای کاربردهایی با تعداد زیادی مشاهده برای ارزیابی؛ برای مثال، در سیستم‌هایی برای خواندن چک‌ها و کدهای پستی، بسیار حیاتی است. یک دسته‌بندی کننده تنک، حافظه کمتری می‌گیرد. مسایل معکوس خطی، رگرسیون و منظم‌سازی

با T_n داده شده، هدف ما، ارزیابی ضرایب رگرسیون حقیقی a_0, a_1, \dots, a_p از مدل خطی

$$Y = \sum_{j=1}^p a_j X_j + a_0 \quad (9.1)$$

بگونه‌ای است که، یک ارزیابی تجربی از مقدار مورد انتظار یک معیار زیان خاص، کمینه شود. توجه داشته باشید که، متغیر تصادفی Y متغیر پاسخ (واکنش) و متغیر تصادفی X شامل متغیرهای پیشگو یا متغیرهای توصیفی X_1, \dots, X_p است. در غیاب نویز و با استفاده از ماتریس برداری، مساله رگرسیون به یافتن جواب برای دستگاه معادلات زیر تقلیل پیدا می‌کند

$$y = Xa \quad (10.1)$$

که در آن $y = [y_1, \dots, y_n]^T$ بردار پاسخ‌های مشاهده شده، X یک ماتریس $n \times p$ تعریف شده به صورت $X = [x_1 | x_2 | \dots | x_n]$ است که در آن، $x_i = [x_{i1}, \dots, x_{ip}]^T$ و $a = [a_1, \dots, a_p]^T$ بردارهای هدف ضرایب رگرسیون هستند. در ادامه بحث، فرض می‌کنیم که جمله انحراف a_0 در نمایش ما صریح است، یعنی

$$x_i^T = [x_i^T, 1] \text{ and } a^T = [a^T, a_0] \quad (11.1)$$

می‌خواهیم مستقیماً معادله (10,1) را حل کنیم، اما X در مساله ما، تقریباً معکوس پذیر نیست. اگر $n = p$ یا $n < p$ باشد، دستگاه معادلات مشخص شده با (10,1) فرامعین² است. در این حالت، جوابی موجود نیست. اگر $n < p$ باشد، دستگاه فرو معین³ است و تعداد نامتناهی جواب وجود دارد. بجای حل (10,1)، می‌توان یک اندازه مشخص از باقی‌مانده بین طرفین معادله را کمینه کرد. یک معیار زیان مشترک برای کاربردهای رگرسیون، تابع زیان کمترین مربعات است. خطای کمترین مربعات، مجموعه باقی‌مانده‌های مربعی است، آن هم در جایی که، باقی‌مانده‌ها به صورت

$$r = y - \sum_{j=1}^p a_j x_j - a_0. \quad (12.1)$$

تعریف می‌شوند، با قرار دادن همه آنها با هم، خطای کمترین مربعات به صورت

$$\|y - Xa\|_2^2. \quad (13.1)$$

در می‌آید، به عبارت دیگر، ما مجاز به مقداری نویز در مدل هستیم، و هدف کمینه‌سازی یک تابع مشخص از آن نویز است. راهی برای ارزیابی ضرایب رگرسیون a ، کمینه‌سازی (13,1) است. کمینه‌سازی

² Overdetermined

³ Underdetermined

(۱,۱۳)، یک مساله بهینه‌سازی محدب نامشروط با یک هدف دیفرانسیل پذیر است، که دارای یک کمینه کلی یکتاست. این هدف در (۱,۱۳)، زمانی کمینه می‌شود، که گرادیان آن نسبت به a صفر باشد. این مطلب منجر به دستگاه معادلات زیر به نام معادلات نرمال می‌شود

$$(X^T X) \hat{a} = X^T y. \quad (۱۴.۱)$$

معادلات نرمال می‌توانند به گونه‌ای موثر حل شوند، اگر $X^T X$ رتبه کامل باشد، در این حالت، X باید دارای رتبه ستونی کاملی باشد. با این حال، این شرط در بسیاری از کاربردهای منجر به ناپایداری عددی، برقرار نیست. این امر می‌تواند منجر به ارزیابی کمترین مربعات با بالاترین بی دقتی شود. جوابی بر این مساله، منظم‌سازی تیخونف^۴ است. در منظم‌سازی تیخونف، مربع ℓ_2 نرم ضرایب به هدف افزوده می‌شود، و در نتیجه اولویت را به جواب‌های با ℓ_2 نرم‌های کوچک‌تر می‌دهد. ایده، قرار دادن فرض‌های پیشین در جواب است. منظم‌سازی تیخونف منجر به مساله بهینه‌سازی

$$\|y - Xa\|_2^2 + \lambda \|a\|_2^2 \quad (۱۵.۱)$$

می‌شود، که در آن، λ پارامتر منظم‌سازی است که، تعامل بین کمینه‌سازی جمله کمترین مربعات و کمینه‌سازی جمله جریمه ℓ_2 را کنترل می‌نماید. هدف در (۱۵,۱)، با یک کمینه کلی یکتای داده شده به عنوان جواب دستگاه معادلات زیر، محدب و دیفرانسیل پذیر باقی می‌ماند

$$(X^T X + \lambda I) \hat{a} = X^T y. \quad (۱۶.۱)$$

توجه داشته باشید که، $X^T X + \lambda I$ حتی اگر $X^T X$ تکیه‌ناپذیر باشد، تکیه‌پذیر است. کمینه‌سازی (۱۵,۱) یک رگرسیون مرزی^۵ نامیده می‌شود.

جریمه ℓ_2 بخشی از خانواده توان جریمه‌ها^۶ با اندیس‌های $\gamma \geq 0$ است

$$\sum_{j=1}^p |a_j|^\gamma. \quad (۱۷.۱)$$

^۴ Tikhonov regularization

^۵ Ridge regression

^۶ Power family of penalties

این یک ℓ_2 نرم از پارامتر بر خاسته از توان γ است [۲۰]. برای $\gamma = 2$ منظم سازی تیخونف یا رگرسیون مرزی را انتخاب می کنیم. مشخص شده است زمانی که $\gamma = 2$ باشد، منظم سازی منجر به یک بردار ضریب چگالی می شود، یعنی تقریباً همه مقادیر ناصفر هستند. منظم سازی با ℓ_2 تنها دارای ویژگی انقباضی^۷ است - انقباض مقادیر مطلق ضریب. در سوی دیگر طیف، $\gamma = 0$ یک بردار ضریب تنک را تولید می کند. بنابراین، می گوئیم دارای ویژگی انتخاب متغیر است. رگرسیون کمترین مربعات با $\gamma = 0$ ، رگرسیون همه زیرمجموعه ها^۸ نامیده می شود. جریمه با ℓ_0 نرم بسیاری از مقادیر بردار ضریب را صفر می سازد، اما هیچ مقدار ناصفیری منقبض نمی شود. در این بین $\gamma = 1$ وجود دارد. بهینه سازی با جریمه ℓ_1 بهترین نتایج انقباض و انتخاب را به دنبال دارد. بهینه سازی کمترین مربعات با منظم سازی ℓ_1 توسط [۱۰] به عنوان نویز زدایی تعاقبی پایه ای^۹ (BPDN) و توسط Tibshirani [۴۹] به عنوان کمترین انتخاب مطلق و عملگر انقباض^{۱۰} (LASSO) معرفی شد و به شهرت رسید. نگاهی دقیق تر به مساله کمترین مربعات ℓ_1 منظم شده، می اندازیم

$$\min_a \frac{1}{2} \|y - Xa\|_2^2 + \lambda \|a\|_1. \quad (18.1)$$

هدف در (۱۸،۱)، محدب است، اما به دلیل ℓ_1 نرم، دیگر دیفرانسیل پذیر نیست. پس، هیچ جواب بسته ای مشابه با (۱۳،۱) و (۱۵،۱) موجود نیست. در فصل ۲، به سراغ چندین روش برای حل (۱۸،۱) می رویم.

همچنین به جریمه شبکه کشسانی^{۱۱} علاقه مند هستیم. علاوه بر این، بسراغ ۳ شکل هم ارز منظم سازی تیخونف نیز می رویم.

⁷ Shrinkage

⁸ All-subset regression

⁹ Basis Pursuit Denoising

¹⁰ Least Absolute Selection and Shrinkage Operator

¹¹ Elastic-net penalty

۱.۱.۱ رگرسیون لجستیکی

در رگرسیون لجستیکی، تابع رگرسیون دارای یک رابطه غیرخطی با ترکیب خطی متغیرهای توصیفی است. این رابطه با تابع پروبیت^{۱۲} مدل سازی می شود.

در تعاریف دسته بندی، پاسخ یک متغیر دودویی یعنی $y_i \in \{-1, 1\}$ است. داده پاسخ برای منظم سازی یک متغیر تصادفی برنولی Y با احتمال موفقیت $\eta = P\{Y=1\}$ است. احتمال موفقیت وابسته به پیشگو، یعنی $\eta = \eta(x)$ ؛ در نظر گرفته می شود. برای توزیع برنولی، مشخص شده است که $E\{Y\} = \eta$. اگر پیشگو برای منظم سازی یک متغیر تصادفی X در نظر گرفته شود، سپس، $\eta(x)$ انتظار شرطی Y است

$$E\{Y | X\} = \eta(x). \quad (۱۹.۱)$$

در رگرسیون خطی، انتظار شرطی Y مقدار x از X را به عنوان تابعی آفین از x ارائه می کند

$$E\{Y | X\} = \eta(x). \quad (۲۰.۱)$$

با این حال، در دسته بندی، یک تابع اتصال^{۱۳} یکنوای g انتظار را به ترکیب خطی پیشگوها، منتقل می کند

$$g(E\{Y | X\}) = a^T x. \quad (۲۱.۱)$$

مدلهایی به این شکل، مدل های خطی تعمیمی نامیده می شوند. در حالت رگرسیون لجستیکی، تابع اتصال به صورت لگاریتمی زیر انتخاب می شود

$$g(a) = \ln \frac{a}{1-a}. \quad (۲۲.۱)$$

معکوس تابع لگاریتمی، تابع لجستیکی است که با $\sigma(z)$ نشان داده می شود

$$g^{-1}(z) = \sigma(z) = \frac{1}{1 + \exp(-z)} \quad (۲۳.۱)$$

^{۱۲} Probit function

^{۱۳} Link function

بنابراین، انتظار شرطی به شکل زیر در می‌آید

$$E\{Y | X\} = g^{-1}(a^T x) = \sigma(a^T x). \quad (۲۴.۱)$$

به منظور ارزیابی پارامتر ناشناخته a از مدل، با تشکیل تابع درست نمایی شروع می‌کنیم. فرض می‌کنیم که، مشاهدات بطور مستقل تولید می‌شوند. تابع احتمال از نظر جبری، با تابع چگالی احتمال مشترک مشاهدات، یکسان است، پس

$$\begin{aligned} L(a) &= p(y | X; a) = \prod_{i=1}^n p(y_i | x_i; a) \\ &= \prod_{i=1}^n \sigma(a^T x_i)^{y_i} (1 - \sigma(a^T x_i))^{1-y_i}. \end{aligned} \quad (۲۵.۱)$$

به دلایل محاسباتی، بجای لگاریتم تابع درست نمایی، تابع

$$\ell(a) = \log L(a) = \sum_{i=1}^n \log p(y_i | x_i; a). \quad (۲۶.۱)$$

را در نظر می‌گیریم، بیشینه‌سازی احتمال معادل با کمینه‌سازی منفی لگاریتم درست نمایی است. مساله بهینه‌سازی به حل مساله زیر تبدیل می‌شود

$$\min_a \sum_{i=1}^n -\log p(y_i | x_i; a) = \min_a \sum_{i=1}^n -\log \sigma(y_i a^T x_i) \quad (۲۷.۱)$$

$$\min_a \sum_{i=1}^n (1 + \exp(-y_i a^T x_i)). \quad (۲۸.۱)$$

که در آن (۲۷،۱) را دنبال می‌کنیم، زیرا $y_i \in \{-1, 1\}$ است. تابع

$$\log(1 + \exp(-y_i a^T x_i)) \quad (۲۹.۱)$$

در (۲۸،۱)، اغلب به عنوان زیان لجستیکی^{۱۴} مورد اشاره قرار می‌گیرد. همان‌طور که دیده می‌شود، بهینه‌سازی زیان لجستیکی برای تولید یک ارزیابی درست نمایی بیشینه، رخ می‌دهد. بنابراین، رگرسیون

^{۱۴} Logistic loss

لجستیکی، برای ارزیابی درست نمایی بیشینه بودن، بررسی می‌شود، البته اگر احتمال پسین $\eta(x)$ بتواند به صورت $1/(1 + \exp(-f(x)))$ برای $f(x) \in F$ بیان شود.

در عوض، بیایید فرض کنیم که علاقه‌مند به ارزیابی احتمال پسین بیشینه (MAP) از یک پارامتر a هستیم. یک لاپلاسین پیشین روی پارامترها را در نظر می‌گیریم. توزیع احتمال لاپلاسین چند متغیره به شکل

$$p(a) = \left(\frac{\lambda}{2}\right)^n \exp(-\lambda \|a\|_1). \quad (30.1)$$

است، مساله بهینه‌سازی، نیازمند حل

$$\max_a \sum_{i=1}^n -\log p(y_i | x_i; a) p(a) = \max_a \sum_{i=1}^n -\log p(y_i | x_i; a) + \log p(a) \quad (31.1)$$

$$\min_a \sum_{i=1}^n -\log p(y_i | x_i; a) + \lambda \|a\|_1 \quad (32.1)$$

$$\min_a \sum_{i=1}^n \log(1 + \exp(-y_i a^T x_i)) + \lambda \|a\|_1. \quad (33.1)$$

است. بهینه‌سازی در $(33,1)$ ، به عنوان مساله رگرسیون لجستیکی ℓ_1 نامگذاری می‌شود.

۲.۱ تقریب تنک

به تازگی، روش‌هایی برای تقریب و کدگذاری تنک، مورد توجه قرار گرفته است [۶]. بنابراین، نیاز به بهره‌گیری از نتایج جدید و تجدید شده برای ارائه وجود دارد. در ادامه، به صورت خلاصه، به مساله تقریب تنک می‌پردازیم.

ماتریس رتبه کامل $A \in R^{n \times m}$ را با $n < m$ در نظر بگیرید، می‌خواهیم دستگاه معادلات $Ax = b$ را حل کنیم. بوضوح، دستگاه معادلات فرو معین و دارای تعدادی نامتناهی جواب است. مساله را با نیاز به پراکنده بودن x محدود می‌کنیم، یعنی x دارای تعدادی درآیه ناصفر باشد. ℓ_0 شبه نرم را به صورت

$$\|x\|_0 = \#\{i : x_i \neq 0\}$$

تعریف می‌کنیم، براساس این تعریف، x پراکنده است، اگر $\|x\|_0 < m$ باشد. مساله بهینه‌سازی که می‌خواهیم حل کنیم، به صورت زیر است

$$\text{minimize } \|x\|_0 \quad \text{subject to } Ax = b. \quad (34.1)$$

با این حال، بدلیل جنبه ترکیبیاتی ℓ_0 نرم، مساله بهینه‌سازی NP-سخت است. در طول سال‌ها، تلاش‌های بسیاری برای حل مساله تقریب تنک (34,1) صورت گرفته است، برای مثال [36] را ببینید. همچنین، تلاش‌هایی برای برداشتن محدودیت‌ها و مهارسازی تابع هدف (34,1) صورت پذیرفته است. بسیاری از این روش‌ها، یک نسخه بی قید از این مساله بهینه‌سازی را حل می‌کنند، برای مثال [10]؛ [49] را ببینید. در طول پایان‌نامه در موارد بسیاری به این موضوع برخورد خواهیم گشت.

۱.۲.۱ نمایش تنک سیگنال‌ها و ویژگی‌ها

در سال‌های اخیر، تقریب تنک به عنوان یک روش یادگیری ویژگی بدون نظارت موفق، به ثبوت رسیده است [۱۲]. می‌توان از تقریب تنک (بخش ۳,۱) برای دستیابی به یک نمایش تنک (بعد بالاتر) از یک بردار ویژگی داده شده، استفاده نمود. این روش به عنوان کدگذاری تنک در منابع، شناخته می‌شود. Yang و همکاران [۶۱] نشان دادند که، کدگذاری تنک، زمانی که ویژگی‌های یادگیری در یک کار دسته‌بندی تصویری بکار رود، کارایی دسته‌بندی SVM خطی را بهبود می‌بخشد.

اگر فرهنگ لغات بکار رفته در کدگذاری تنک نمونه‌های آزمایشی با برچسب‌های شناخته شده، ساخته شده باشد، بازنمایی تنک حاصل می‌تواند برای دسته‌بندی بکار رود [۵۵]، نشان می‌دهیم که این آماده سازی می‌تواند برای رگرسیون نیز بکار گرفته شود.

۲

فصل دوم

کمینه‌سازی زیان مربعی ℓ_1 -منظم برای دسته‌بندی

همان‌طور که در فصل ۱ توضیح داده شد، ما علاقه‌مند به مساله بهینه‌سازی کمترین مربعات ℓ_2 -منظم شده (۱۸،۱)، هستیم. در این فصل، به سراغ روال‌هایی موثر برای اعمال این روش‌ها در مسایل بهینه‌سازی و مهم‌تر از آن، بکارگیری آنها در تنظیمات دسته‌بندی می‌رویم.

دلایل بسیاری برای استفاده از اصل بهینه‌سازی لسو برای آموزش دادن دسته‌بندی‌کننده‌ها وجود دارد. در این فصل اهمیت استفاده از قاعده‌ی بهینه‌سازی لسو را برای مسائل دسته‌بندی هم از دیدگاه نظریه‌های آمار و احتمال و هم از دیدگاه آزمایشی بررسی می‌کنیم. همانطور که در ابتدای فصل قبل گفته شد لسو دارای سه خصیصه‌ی مهم است: ۱. خطی بودن با ۲. زیان مربعی، و ۳. منظم‌سازی ℓ_2 روی وزن‌ها. در ادامه اهمیت این خصایص را بررسی می‌کنیم.

در این پایان‌نامه، بر دسته‌بندی خطی تمرکز داریم. دو دلیل اصلی برای این تصمیم وجود دارد. دسته‌بندی غیرخطی، مانند SVM غیرخطی، چندان برای آموزش و ارزیابی بکار نمی‌رود. این امر، بویژه، شکستی در حالت چند کلاسه است که، دسته‌بندی‌های بسیاری برای یک کار آموزش می‌بینند، و به‌صورت یک در مقابل یک یا یک در مقابل همه، تنظیم می‌شوند. علاوه بر این، برای مسایل با بعد بالا، دسته‌بندی غیرخطی، مزیت خاصی را ارائه نمی‌نماید. فضاهای داده با ابعاد بالا -مگر در کاربردهایی خاص- معمولاً پراکنده هستند، و بنابراین، به احتمال زیاد جدایی‌پذیر خطی می‌شوند.

با حرکت به‌سوی ساخت یک دسته‌بندی‌کننده^۱ سریع‌تر، تنک بودن را از طریق منظم‌سازی توسط نرم القاکننده‌ی تنک بودن، تضمین می‌کنیم. این کار، امکان کنترل صریح را بر تنک بودن بردار ضریب از طریق پارامتر منظم‌سازی فراهم می‌سازد- برخلاف بهینه‌سازی SVM که چنین کنترلی در آن چندان مشخص نیست (شکل ۱،۲ را ببینید). بویژه، ما از الگوریتم‌های کمینه‌سازی زیان مربعی ℓ_2 -منظم شده از جمله، SPGL1 [۵۲]، برای حل مساله دسته‌بندی استفاده می‌کنیم. این مطلب کار ما را از مطالعه اخیر صورت گرفته توسط Yuan و همکاران [۶۲]، [۶۳] متمایز می‌سازد. آنالیزهای آنها تنها بر الگوریتم‌های بهینه‌سازی لجستیکی و محوری مربعی- که L_2 زیان هم نامیده می‌شود- تاکید دارد، این مسایل به‌صورت زیر تعریف می‌شود

$$\max(0, 1 - y a^T x)^2. \quad (۵.۲)$$

^۱ Classifier

زیان لجستیکی برای حالت زیان مربعی، دوبار دیفرانسیل پذیر است. با این حال، زیان محوری مربعی این گونه نیست.

پس از مرور بر کارهای مرتبط در این زمینه به بیان نتایج آماری پرداخته و سپس نتایج آزمایش‌های انجام شده برای نشان دادن قدرت لسو در دسته‌بندی را بیان می‌کنیم. نشان می‌دهیم که لسو در عین دقت بالا، دسته‌بنده‌کننده‌ی تنک و لذا سریع‌تر به دست می‌دهد.

۱.۱.۲ کارهای مرتبط در کمینه‌سازی زیان مربعی برای دسته‌بندی

Rifkin [۴۲] در رساله دکترای خود، ادعا می‌کند که، زیان محوری راهی برای موفقیت SVM نیست. Rifkin دسته‌بندی کمترین مربعات منظم شده (RLSC)^۱

$$\min_w \|y - Xw\|^2 + \lambda \|w\|^2 \quad (۶.۲)$$

و حالت غیرخطی

$$\min_c \|y - Kc\|^2 + \lambda c^T Kc \quad (۷.۲)$$

را ارائه می‌نماید. که در آن، K ماتریس تعریف شده با ارزیابی هسته K روی زوج نمونه‌های آزمایشی است. با استفاده از چندین مطالعه آزمایشی، Rifkin نشان داد که، RLSC به اندازه SVM روی مجموعه داده‌های گوناگون، خوب است. با این حال، RLSC منجر به دسته‌بندی کننده تنک نمی‌شود. علاوه بر این، RLSC غیرخطی بیش از SVM برای آزمایش زمان می‌برد.

بهینه‌سازی کمترین مربعات می‌تواند برای دسته‌بندی دودویی بکار رود [۲۶]، [۴۴]، [۴۴]، [۶۴]، [۶۵]. آن را همچنین می‌توان در زمینه مدل‌های خطی تعمیمی-برای مثال رگرسیون لجستیکی با بهینه‌سازی کمترین مربعات باز وزنی تکراری [۳۰]، [۴۹] و برای m -ارزیابی [۱۷] بکار برد. ما یک دسته‌بندی کننده

^۱ Regularized Least Square Classification

سریع را بدست می آوریم، که هنوز بعضی از قابلیت های تعمیمی را بدلیل ℓ_1 منظم‌سازی، حفظ می کند. با این حال، بعدها آنالیزهای بیشتری صورت گرفت [۲۲]، [۴۰]، [۶۰].

۲.۱.۲ کارهای مرتبط در ℓ_1 -منظم‌سازی برای دسته‌بندی کننده تنک

Yuan و همکاران (۲۰۱۰) [۶۲] چندین دسته‌بندی کننده خطی تنک به فرم

$$\min_w \|w\|_1 + C \sum_{i=1}^n \xi(w; x_i, y_i) \quad (۸.۲)$$

را با زیان لجستیکی، محوری و محوری مربعی تعریف شده به صورت

- $\xi \log(w; x_i, y_i) = \log(1 + \exp(-yw^T x))$
- $\xi L_1(w; x_i, y_i) = \max(1 - yw^T x, 0)$
- $\xi L_2(w; x_i, y_i) = \max(1 - yw^T x, 0)^2$

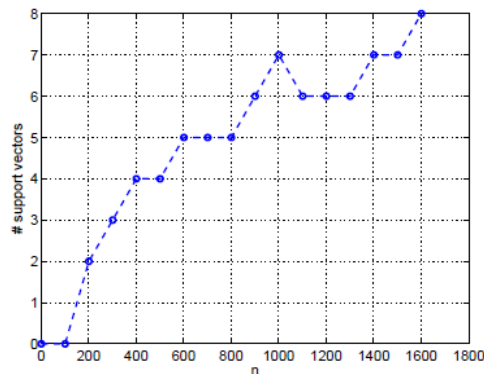
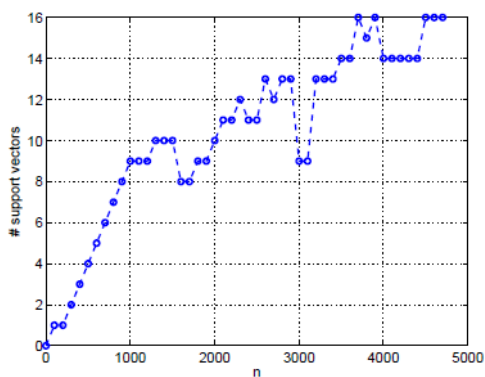
مقایسه می کنند، با این حال، آنها بهینه‌سازی زیان مربعی برای دسته‌بندی را بررسی نمی کنند. همان طور که در فصل بعدی نشان می دهیم، زیان مربعی دارای ویژگی های محاسباتی و آماری مناسبی است.

از طریق چندین آزمایش و اثبات ساده نشان می دهیم که، مزیت قابل توجهی در ترجیح دادن یک تابع زیان محدب بر دیگران برای کمینه‌سازی در تعاریف دسته‌بندی وجود ندارد. علاوه بر این، مزیت هایی برای بهینه‌سازی زیان مربعی بر زیان محوری و زیان لجستیکی موجود است [۲۶]، [۶۴]. زمانی که الگوریتم های کمینه‌سازی زیان مربعی نه چندان گرانی وجود داشته باشد، استفاده از آنها برای دسته‌بندی، بنظر قابل قبول و موثر است. جدول ۱، شامل اطلاعاتی درباره مجموعه داده هایی است که ما برای آزمایش خود در این فصل و فصل بعدی بکار می بریم. آنها از منزلگاه (وبسایت) LIBSVM برای مجموعه داده های رده دودویی، گرفته شده اند.^۱ این مجموعه داده ها برای نمایش های گوناگون امور عملی، انتخاب شده اند.

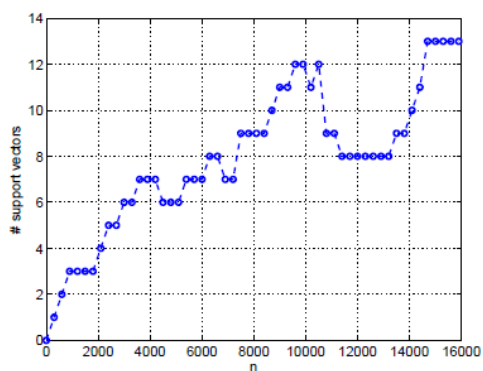
^۱ <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/dataset/binary.html>

جدول ۱.۱ اطلاعات مجموعه داده‌ها: n نشان‌دهنده تعداد مشاهدات و p نشان‌دهنده تعداد ویژگی‌های هر مشاهده است.

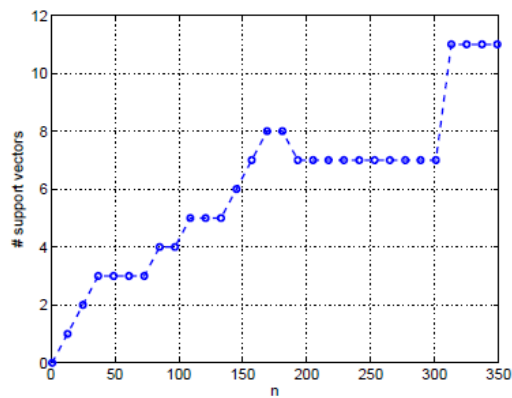
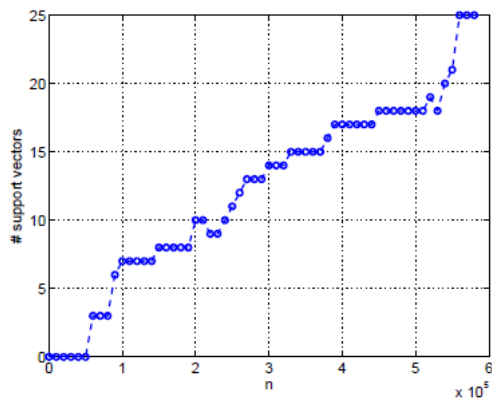
تعداد ویژگی‌ها (p)	تعداد مشاهدات (n)	مجموعه داده‌ها
۱۲۳	۱۶۰۵	بزرگسال ۱
۱۲۳	۴۷۸۱	بزرگسال ۴
۱۲۳	۱۶۱۰۰	بزرگسال ۷
۱۲۳	۳۲۵۶۱	بزرگسال ۹
۱۴	۶۹۰	استرالیا
۲۰۰۰	۶۲	سرطان کولون
۵۴	۵۸۱۰۱۲	نوع پوشش
۸	۷۶۸	دیابت
۱۳	۲۷۰	قلب
۳۴	۳۵۱	یونکره
۷۱۲۹	۳۸	ال.ای.یو
۶	۳۴۵	تجویزات کبدی
۱۱۲	۸۱۲۴	قارچ‌ها



(الف) مجموعه داده بزرگسال ۱



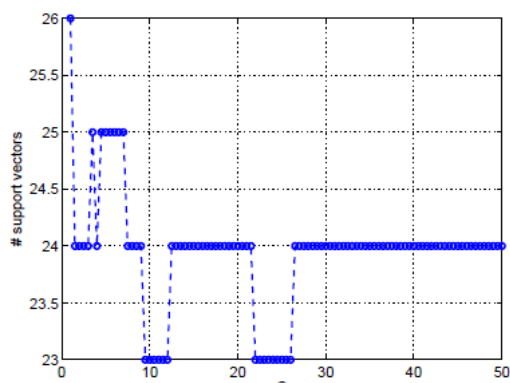
(ج) مجموعه داده بزرگسال ۷



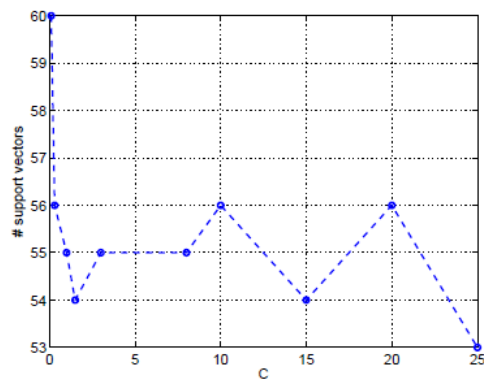
(ه) مجموعه داده یونکره

(د) مجموعه داده نوع پوشش

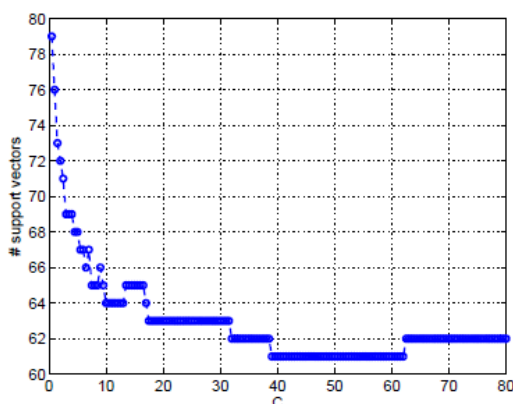
شکل ۱.۱ نمودار، تعداد بردارهای پشتیبان را نشان می‌دهد که با افزایش تعداد نمونه‌های آزمایشی رشد می‌کند.



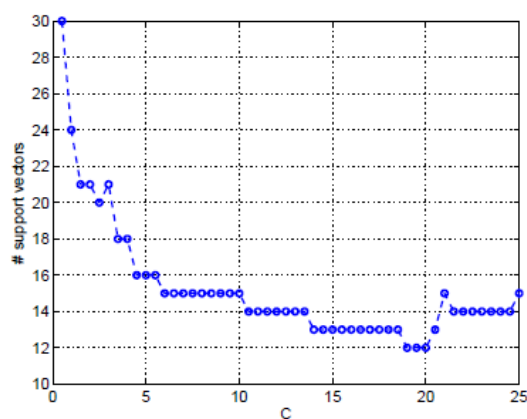
(ب) مجموعه داده‌های قلب



(الف) مجموعه داده‌های استرالیا



(د) مجموعه داده‌های تجزیه‌های کبدی



(ج) مجموعه داده‌های یونکره

شکل ۲.۲ نمودارها نشان می‌دهد که تعداد بردارهای پشتیبان ارتباط معنی‌داری با فرایارامتر C در بهینه‌سازی SVM ندارند.

۲.۲ دسته‌بندی و کمینه‌سازی زیان محدب

تشخیص الگو، که دسته‌بندی نیز نامیده می‌شود، فرآیند انتساب یک برچسب گسسته به یک مشاهده ناشناخته است [۱۵]. در تشخیص الگو، کار اصلی، یافتن یک تابع $g: R^p \rightarrow \{1, \dots, M\}$ است، که یک

مشاهده نشان داده شده توسط $x \in R^p$ را گرفته و آن را به $y \in \{1, \dots, M\}$ نسبت می‌دهد: یکی از کلاس‌های^۱ در دسترس M . تابع g یک دسته‌بندی کننده^۲ نامیده می‌شود.

برای آنالیز فعلی، فرض کنید X و Y متغیرهایی تصادفی باشند که مقادیرشان را به ترتیب از R^p و $\{1, \dots, M\}$ می‌گیرند. احتمال خطا را برای دسته‌بندی کننده g به صورت

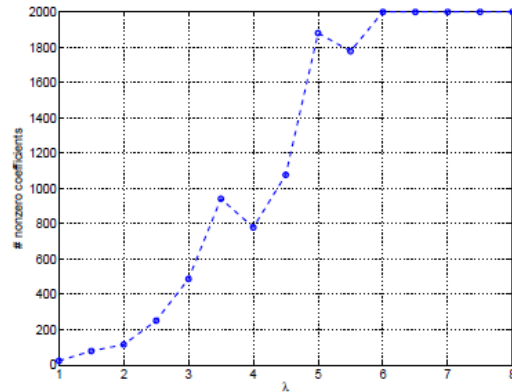
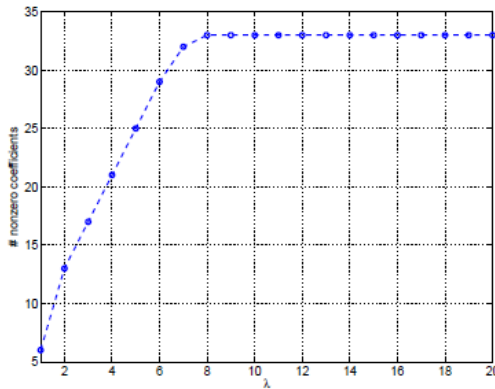
$$L(g) = P\{g(X) \neq Y\}. \quad (۲.۲)$$

تعریف می‌کنیم، در نتیجه، دسته‌بندی کننده بهینه g^* عبارت است از

$$g^* = \underset{g: R^p \rightarrow \{1, \dots, M\}}{\operatorname{argmin}} P\{g(X) \neq Y\} \quad (۳.۲)$$

و دسته‌بندی کننده بیزی نامیده می‌شود. احتمال متناظر از خطا- احتمال کمینه خطا- خطای بیزی نامیده شده و با $L^* = L(g^*)$ نشان داده می‌شود.

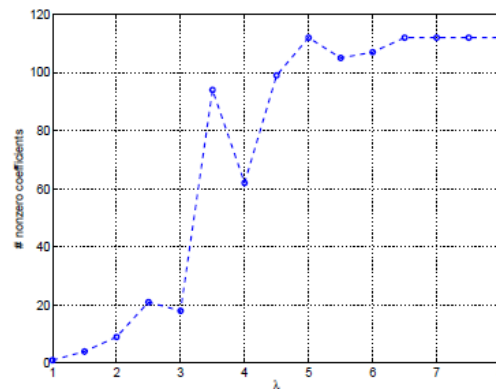
برای محاسبه g^* ، نیازمند دانش از توزیع (X, Y) ای هستیم که ناشناخته است. با این حال، با داده‌های کافی، یک دسته‌بندی کننده g را با $L(g)$ کم، می‌یابیم. برای دسته‌بندی، به مجموعه داده‌های $T_n = \{(x_i, y_i) : i = 1, \dots, n\}$ از n مشاهده دسترسی داریم، و فرض می‌کنیم که $i.i.d.$ های نمونه‌گیری شده از توزیع (X, Y) باشند.



¹ Class

² Classifier

(الف) مجموعه داده‌های سرطان کولون (ب) مجموعه داده‌های یونکره



(ج) مجموعه داده‌های قارچ‌ها

شکل ۱.۲ در این اشکال، می‌بینیم که، تعداد عناصر ناصفر جواب، زمانی افزایش می‌یابد که، پارامتر منظم‌سازی λ افزایش پیدا کند، بنابراین، ابزاری را برای کنترل تنک بودن جواب در اختیار ما قرار می‌دهد.

در باقی‌مانده این پایان‌نامه، حالت دسته‌بندی دودویی یعنی $M = 2$ را در نظر می‌گیریم. در این حالت، Y مقادیرش را از مجموعه دودویی $\{-1, 1\}$ می‌گیرد. دسته‌بندی دودویی می‌تواند به دسته‌بندی چندکلاسی ($M > 2$) در یک تنظیم یک در مقابل همه (OVA)^۱ گسترش یابد [۴۳]. قبل از کاوش دقیق موضوع دسته‌بندی، مفهوم مهمی بنام احتمال پسین را در این زمینه معرفی می‌کنیم.

۱.۲.۲ احتمال پسین و دسته‌بندی کننده‌ی جانشین

در آنالیز رگرسیون، هدف ارزیابی Y برای X مفروض با $r(X)$ است، که در آن $r: R^p \rightarrow R$ یک تابع می‌باشد. می‌توان نشان داد که، تابعی که خطای مربع میانگین را در این چارچوب کمینه می‌سازد، احتمال پسین η است

$$\eta(x) = P\{Y = 1 | X = x\} = E\{Y | X = x\}. \quad (۴.۲)$$

برای همه، $r: R^p \rightarrow R$ داریم

^۱ One-Versus-All

$$E\{(\eta(X) - Y)^2\} \leq E\{(r(X) - Y)^2\} \quad (5.2)$$

اثبات. برای هر $x \in R^p$ داریم

$$\begin{aligned} E\{(r(X) - Y)^2 \mid X = x\} \\ &= E\{(r(X) - \eta(X) + \eta(X) - Y)^2 \mid X = x\} \\ &= (r(x) - \eta(x))^2 + 2(r(x) - \eta(x))E\{\eta(X) - Y \mid X = x\} + E\{\eta(X) - Y\}^2 \mid X = x\} \\ &= (r(x) - \eta(x))^2 + E\{\eta(X) - Y\}^2 \mid X = x. \end{aligned} \quad (6.2)$$

توجه داشته باشید که تساوی در (5,2) برقرار است، اگر و تنها اگر برای همه $x \in R^p$ داشته باشیم

$$r(x) = \eta(x)$$

اهمیت احتمال پسین در این است که با یک η داده شده، می‌توان یک دسته‌بندی کننده را با احتمال کمینه خطا ساخت. دسته‌بندی کننده $g^*: R^p \rightarrow \{-1, 1\}$ را با استفاده از تابع رگرسیون η به صورت

$$g^*(x) = \begin{cases} -1 & \text{if } \eta(x) \leq \frac{1}{2} \\ 1 & \text{otherwise.} \end{cases} \quad (7.2)$$

تعریف می‌کنیم، ادعا این است که، g^* یک دسته‌بندی کننده بیزی است، یعنی احتمال خطا را کمینه می‌سازد. برای اثبات این مطلب، باید نشان دهیم که برای هر دسته‌بندی کننده

$$E\{g^*(\eta(X) - Y)^2\} \leq P\{g(X) \neq Y\}. \quad (8.2)$$

اثبات. برای $X = x$ داده شده، احتمال شرطی خطای دسته‌بندی کننده g عبارت است از

$$\begin{aligned} P\{g(X) \neq Y \mid X = x\} \\ &= 1 - P\{g(X) = Y \mid X = x\} \\ &= 1 - (P\{g(X) = 1, Y = 1 \mid X = x\} + P\{g(X) = -1, Y = -1 \mid X = x\}) \\ &= 1 - (\chi_{\{g(x)=1\}} P\{Y = 1 \mid X = x\} + \chi_{\{g(x)=-1\}} P\{Y = -1 \mid X = x\}) \\ &= 1 - (\chi_{\{g(x)=1\}} \eta(x) + \chi_{\{g(x)=-1\}} (1 - \eta(x))) \end{aligned} \quad (9.2)$$

که در آن χ_A تابع شاخص مجموعه A است. برای هر $x \in R^p$ می‌توانیم بنویسیم

$$\begin{aligned} & P\{g(X) \neq Y \mid X = x\} - P\{g^*(X) \neq Y \mid X = x\} \\ &= \eta(x) \left(\chi_{\{g^*(x)=1\}} - \chi_{\{g(x)=1\}} \right) + (1 - \eta(x)) \left(\chi_{\{g^*(x)=-1\}} - \chi_{\{g(x)=-1\}} \right) \quad (10.2) \\ &= (2\eta(x) - 1) \left(\chi_{\{g^*(x)=1\}} - \chi_{\{g(x)=1\}} \right) \geq 0 \end{aligned}$$

که آخرین تساوی برقرار است، زیرا

$$\begin{aligned} \chi_{\{g^*(x)=-1\}} &= 1 - \chi_{\{g^*(x)=1\}} \\ \chi_{\{g(x)=-1\}} &= 1 - \chi_{\{g(x)=1\}}. \end{aligned} \quad (11.2)$$

با انتگرال گیری نسبت به x می توانیم به $(8,2)$ برسیم.

تابع $\eta(x)$ ناشناخته است. برای تقریب دسته بندی کننده بیزی، ما تابع نامنفی $\tilde{\eta}(x)$ را به عنوان تقریبی از $\eta(x)$ بکار گرفته و آن را در شکل دسته بندی کننده در معادله $(18,2)$ جایگذاری می کنیم

$$g(x) = \begin{cases} -1 & \text{if } \tilde{\eta}(x) \leq \frac{1}{2} \\ 1 & \text{otherwise.} \end{cases} \quad (12.2)$$

دسته بندی کننده g ، دسته بندی کننده ی جانشین نامیده می شود. دسته بندی کننده نشانده بخوبی عمل می کند. اگر $\tilde{\eta}(x)$ نزدیک به $\eta(x)$ باشد - که با L_1 نرم مورد انتظار، اندازه گیری می شود - سپس احتمال خطای دسته بندی کننده نشانده نزدیک به خطای بیزی است، یعنی

$$P\{g(X) \neq Y\} - L^* \leq 2E\{|\eta(X) - \tilde{\eta}(X)|\}. \quad (13.2)$$

اثبات. توجه داشته باشید که تمایز بین احتمالات خطای شرطی g و g^* زمانی که برای هر $x \in R^p$ ، $g(x) = g^*(x)$ باشد، صفر است. زمانی که $g(x) \neq g^*(x)$ باشد، براساس $(10,2)$ ، می توانیم تمایز را به صورت

$$\begin{aligned} & P\{g(X) \neq Y \mid X = x\} - P\{g^*(X) \neq Y \mid X = x\} \\ &= (2\eta(x) - 1) \left(\chi_{\{g^*(x)=1\}} - \chi_{\{g(x)=1\}} \right) \\ &= |2\eta(x) - 1| \chi_{\{g(x) \neq g^*(x)\}} \\ &= 2|\eta(x) - 1/2| \chi_{\{g(x) \neq g^*(x)\}} \\ &\leq 2|\eta(x) - \tilde{\eta}(x)|. \end{aligned} \quad (14.2)$$

بنویسیم، نامساوی اخیر برقرار است، زیرا $g(x) \neq g^*(x)$ ایجاب می کند

$$|\eta(x) - \tilde{\eta}(x)| \geq |\eta(x) - 1/2|. \quad (۱۵.۲)$$

برای بررسی درستی این مطلب، هر حالت را به ترتیب، بررسی می‌کنیم

if $g^*(x) = -1 \Rightarrow \tilde{\eta}(x) \leq 1/2$ and $g(x) = 1$ so $\tilde{\eta}(x) > 1/2$

$$|\eta(x) - \tilde{\eta}(x)| = \tilde{\eta}(x) - \eta(x) \leq 1/2 - \eta(x) = |1/2 - \eta(x)|$$

if $g^*(x) = 1 \Rightarrow \eta(x) > 1/2$ and $g(x) = -1$ so $\tilde{\eta}(x) \leq 1/2$

$$|\eta(x) - \tilde{\eta}(x)| = \eta(x) - \tilde{\eta}(x) \leq \eta(x) - 1/2 = |1/2 - \eta(x)|.$$

بر می‌گردیم به اثبات، می‌توانیم با انتگرال گیری طرفین نامساوی در (۱۴,۲) نسبت به X به

$$P\{g(X) \neq Y\} - L^* \leq 2E\{|\eta(X) - \tilde{\eta}(X)|\}$$

برسیم.

این نتیجه نشان می‌دهد که، یک ارزیاب خوب از η می‌تواند یک دسته‌بندی کننده نشانده خوب را تولید کند. آنچه که روشن است این است که، $\tilde{\eta}(x)$ می‌تواند دور از $\eta(x)$ باشد و دسته‌بندی کننده یکسانی را بدست آورد، همراه با این که هر دو روی طرف یکسانی از $1/2$ باشد. ما اکنون، حاضریم روی آموزش بر راه‌هایی جهت دست یابی به دسته‌بندی کننده خوب، تمرکز کنیم.

۲.۲.۲ کمینه‌سازی ریسک آزمایشی

کمینه‌سازی احتمال خطا (احتمال دسته‌بندی نادرست)

$$L(g) = E\{\chi_{\{g(X) \neq Y\}}\} = P\{g(X) \neq Y\} \quad (۱۶.۲)$$

تنها با دانستن توزیع توامان X و Y ممکن است. یک ارزیاب احتمال خطا از یک دسته‌بندی کننده g با T_n داده شده، شمارش خطای متوسط زیر است

$$L_n(g) = \frac{1}{n} \sum_{i=1}^n \chi_{\{g(x_i) \neq y_i\}}. \quad (۱۷.۲)$$

ارزیاب $L_n(g)$ خطای تجربی g است. برای یک بررسی کلی از پیشرفت‌های اخیر [۵] را ببینید.

می‌توانیم به رهیافت کنونی دسته‌بندی با در نظر گرفتن یک کلاس C از دسته‌بندی کننده‌های $\{-1, 1\} : R^p \rightarrow g^*$ برسیم. با T_n داده شده، دسته‌بندی کننده در C را انتخاب می‌کنیم که، منجر به کمترین خطای تجربی $L_n(g)$ می‌شود. با این حال، مساله کمینه‌سازی خطای تجربی، از نظر محاسباتی، بسیار دشوار است. در مواجهه با این مساله، نیازمند اصلاح عملکرد برای کمینه‌سازی هستیم. دسته‌بندی کننده‌هایی به شکل زیر را در نظر می‌گیریم

$$g_f(x) = \begin{cases} -1 & \text{if } f(x) < 0 \\ 1 & \text{otherwise} \end{cases} \quad (18.2)$$

که در آن $f : R^p \rightarrow R$ یک تابع حقیقی مقدار در F است. احتمال خطای g_f به صورت

$$\begin{aligned} Lg_f &= L(f) = P\{\text{sgn}(f(X)) \neq Y\} \\ &= P\{Y f(X) \leq 0\} \\ &= E\{\chi_{Y f(X) \leq 0}\}. \end{aligned} \quad (19.2)$$

است، کمیت $Yf(x)$ حاشیه^۱ نامیده می‌شود و زمینه تکراری در باقی این بخش دارد. با T_n داده شده، می‌توان $L(f)$ را با $L_n(f)$ ارزیابی نمود

$$L_n(f) = \frac{1}{n} \sum_{i=1}^n \chi_{\{Y_i f(x_i) \leq 0\}} \quad (20.2)$$

که در آن $\chi_{\{Yf(x) \leq 0\}}$ تابع زیان ۰ تا ۱ است. همان‌طور که قبلاً توضیح داده شد، کمینه‌سازی خطای تجربی از نظر محاسباتی، دشوار است [۲]. در عوض، به دنبال کمینه‌سازی یک کران بالاتر محدب هموار از زیان ۰ تا ۱، یعنی $\chi_{\{Yf(x) \leq 0\}}$ هستیم. تابع محدب هموار Φ از حاشیه $v = Yf(x)$ تابع هزینه نامیده می‌شود.

^۱ Margin

جدول ۱.۲ توابع زیان محدب شناخته شده و تابع کمینه‌سازی متناظرشان.

نام تابع زیان	$\phi(v)$	$f_{\phi}^*(\eta)$
زیان مربعی	$(1-v)^2$	$2\eta-1$
زیان محوری	$\max(0, 1-v)$	$\text{sign}(2\eta-1)$
زیان محوری مربعی	$\max(0, 1-v)^2$	$2\eta-1$
زیان لجستیکی	$\ln(1+\exp(-v))$	$\ln \frac{\eta}{1-\eta}$

مثال‌ها شامل تابع زیان نمایی $\phi(v) = \exp(-v)$ بکار رفته در AdaBoost [۱۸] و تابع زیان محوری $\phi(v) = \max(0, 1-v)$ استفاده شده در SVM [۴] است. تابعک هزینه عبارت است از

$$A(f) = E\{\phi(Yf(X))\} \quad (۲۱.۲)$$

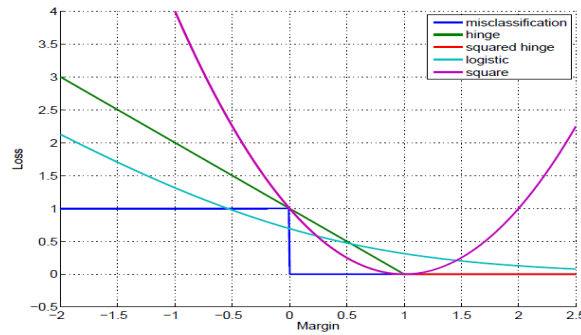
و فرم آزمایشی متناظر نیز به صورت

$$A_n(f) = \frac{1}{n} \sum_{i=1}^n \phi(y_i f(x_i)). \quad (۲۲.۲)$$

است، توجه داشته باشید که، ϕ یک کران بالاتر روی زیان \cdot تا ۱ است و بنابراین، $L_n(f) \leq A_n(f)$ و $L(f) \leq A(f)$.

اگر F شامل توابعی باشد که در پارامترهایشان خطی هستند- بنابراین، آنها را محدب می‌سازیم و سپس کمینه‌سازی هزینه تجربی $A_n(f)$ ، مساله بهینه‌سازی محدب است. پس، الگوریتم‌هایی موثر برای دستیابی به کمینه‌سازی $A_n(f)$ بر F وجود دارد.

توابع زیان مورد توجه ما در جدول ۱،۲ فهرست شده و در شکل ۲،۲ نشان داده شده‌اند. توابع زیان حقیقی مقدار، به عنوان جانشین‌هایی برای زیان \cdot عمل می‌کنند. این منجر به یک مساله رگرسیون می‌شود. نتیجه برای دسته‌بندی با آستانه خروجی تابع حاصل از رگرسیون، بدست می‌آید. در ادامه، به سراغ این واقعیت می‌رویم که، بهینه‌سازی یک جانشین محدب منجر به دسته‌بندی کننده بیزی می‌شود.



شکل ۲.۲ مقایسه‌ای از توابع زیان محدب. زیان دسته‌بندی نادرست نیز نشان داده شده است.

کمینه‌سازی یک کران بالای محدب از زیان $0-1$ ، نه تنها منجر به دسته‌بندی قابل قبولی می‌شود، بلکه کار موفق است. موفقیت‌های آزمایشی SVM، و تقویت‌ها، این نقطه را ایجاب می‌نمایند. می‌توانیم نشان دهیم که کمینه ساز f^* از $A_n(f)$ به گونه‌ای است که، دسته‌بندی کننده القایی g_f^* یک دسته‌بندی کننده بیزی است. برای نیل به این هدف، نگاه نزدیک‌تری به تابع هزینه می‌اندازیم

$$\begin{aligned} A(f) &= E\{\phi(Yf(X))\} \\ &= E\{\eta(X)\phi(f(X)) + (1-\eta(X))\phi(-f(X))\} \end{aligned} \quad (23.2)$$

که در آن، $\eta(x)$ نشان‌دهنده احتمال پسین $P\{Y=1 | X=x\}$ است. تعریف زیر را در نظر بگیرید

$$A(f, \eta) = \eta\phi(f) + (1-\eta)\phi(-f) \quad (24.2)$$

که در آن، $\eta \in [0,1]$ ، $f \in R$ است. فرض کنید که تابع $f_\phi^* : [0,1] \rightarrow R$ کمینه ساز $A(f, \eta)$ باشد

$$f_\phi^*(\eta) = \arg \min_{f \in R} A(f, \eta). \quad (25.2)$$

با تعریف f_ϕ^* از بالا، واضح است که $f_\phi^*(\eta(x))$ مقدار $A(f(x))$ را در میان همه توابع $f(x)$ در $(21,2)$ کمینه می‌سازد. با یک تابع زیان محدب مفروض ϕ ، کمینه ساز بهینه f_ϕ^* به آسانی قابل محاسبه است. در ستون سوم از جدول ۲.۱، f_ϕ^* بهینه را برای توابع زیان مطلوب، لیست کرده‌ایم. در این مثال‌ها، به آسانی می‌توان دید که، $f_\phi^*(\eta(x)) > 0$ تنها زمانی که $\eta > 1/2$ باشد، برقرار است. اگر فرض کنیم $f(x) = f_\phi^*(\eta(x))$ باشد، با کمینه‌سازی $(21,2)$ ، دسته‌بندی کننده القایی

$$g^*(x) = \begin{cases} -1 & \text{if } f(x)^*(x) < 0 \\ 1 & \text{Otherwise} \end{cases}$$

دارای علامتی یکسان با دسته‌بندی کننده بیزی است. این به ما اجازه می‌دهد تا نتیجه بگیریم که، کمینه ساز تابع هزینه $A(f)$ ، کران بالایی روی خطای دسته‌بندی درست $L(f)$ است، به گونه‌ای که دسته‌بندی کننده القایی g_f یک دسته‌بندی کننده بیزی است، و این سازگاری فشر توابع هزینه محذب را اثبات می‌کند [۳، ۴، ۵، ۳۲، ۶۵]. توجه داشته باشید که، سازگاری فشر در این زمینه، به عنوان دسته‌بندی کننده بیزی یعنی $\text{sign}(\eta - 1/2)$ است.

۳.۲.۲ ارزیابی احتمال پسین

یک مشاهده مهم براساس جدول ۲،۱ این است که، SVM بطور مستقیم دسته‌بندی کننده دودویی را ارزیابی می‌نماید. یعنی، SVM مستقیماً $\text{sign}(2\eta(x) - 1)$ را بجای احتمال پسین $\eta(x)$ ارزیابی می‌کند. این بدان معناست که، SVM نمی‌تواند اطلاعاتی درباره اطمینان از پیش‌بینی به ما بدهد. چنین اطلاعاتی بویژه، برای دسته‌بندی چندکلاسه با استفاده از دسته‌بندی کننده دودویی در تنظیمات یکی در مقابل همه، مفید است. در شرایطی که، احتمالات پسین برای همه کلاس‌ها زیر $1/2$ است، SVM نمی‌تواند تصمیمات صحیحی بگیرد.

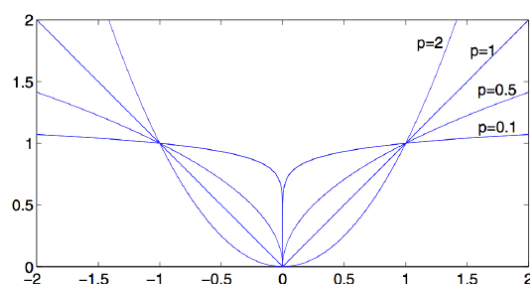
از سوی دیگر، دسته‌بندی کننده کمترین مربعات، احتمال پسین $\eta(x)$ را ارزیابی نموده و بنابراین می‌تواند اطلاعاتی مطمئن را ارائه نماید. برای دسته‌بندی کننده کمترین مربعات، احتمال پسین برابر با $(f(x) + 1)/2$ محدود شده به $[0, 1]$ است. برای آنالیز عمیق‌تر این موضوع [۶۴] را ببینید.

جدول ۲.۲ مقایسه‌ای از سه دسته‌بندی کننده دیگر با دسته‌بندی لسویی روی ۷ مجموعه داده. فراپارامتر C یا λ بسته به الگوریتم نشان داده می‌شود. تعداد درآیه‌های ناصفر در بردار جواب نیز با $\#nz$ و درصد نمونه‌های آزمون به‌درستی دسته‌بندی شده با Acc نمایش داده می‌شود.

ℓ_1 -reg logreg			ℓ_1 -reg L2SVM			SVM			lasso			مجموعه داده‌ها
#nz	دقت	C	#nz	دقت	C	#nz	دقت	C	#nz	دقت	λ	
۱۴	۸۷	۵	۱۴	۸۶	۲۰	۶۸	۸۶	۲	۱۴	۸۶	۱۰	استرالیا
۹۱	۸۳	۱۰	۱۱۲	۷۵	۱۰	۱۱	۸۷	۱	۱۶	۷۷	۱	سرطان کولون
۸	۷۶	۱۰	۸	۷۷	۱۰	۱۰۵	۷۵	۱	۸	۸۰	۱۰	دیابت
۱۳	۸۳	۵	۱۳	۸۰	۱۰	۳۰	۸۴	۱/۵	۱۳	۸۷	۶	قلب
۳۱	۸۲	۲	۲۸	۷۵	۱۰	۱۱	۸۷	۱	۱۶	۷۷	۱	یونکره
۶	۶۷	۵	۶	۶۶	۵	۷۰	۶۲	۲	۶	۴۶	۲	تجویزات کبدی
۹۵	۱۰۰	۲۰	۹۶	۱۰۰	۱۰	۹۰	۱۰۰	۲	۱۳	۴۸	۲	قارچ‌ها

۳.۲ چرا منظم‌سازی ℓ_1 بردار وزن تنک ایجاد می‌کند؟

شکل ۳.۲، شهود مناسبی را برای این که چرا منظم‌سازی ℓ_1 تنک بودن را در بردار وزن القا می‌کند، ارائه می‌دهد.



شکل ۳.۲ همان‌طور که p به صفر می‌رود، $|x|^p$ به تابع شاخص تبدیل می‌شود و تعداد درآیه‌های ناصفر در x را می‌شمارد [۶].

۴.۲ ارزیابی تجربی لسو برای دسته‌بندی

انواع گوناگونی از الگوریتم‌ها برای حل مساله بهینه‌سازی کمترین مربعات ℓ_1 -منظم شده وجود دارد [۳۳]. در این پایان‌نامه، یکی از موفق‌ترین این الگوریتم‌ها را بکار می‌گیریم که، قادر به حل مساله کمینه‌سازی زیان مربعات ℓ_1 -منظم شده، به‌صورت بهینه است. این الگوریتم، روش تصویری طیفی برای ℓ_1 -کمینه‌سازی (SPGL1)^۱ [۵۲] است و مساله بهینه‌سازی زیر (که لسو نام دارد) را حل می‌کند

$$\min_a \|y - Xa\|_2^2 \text{ s.t. } \|a\|_1 \leq \lambda. \quad (۱.۲)$$

همان‌طور که پیشتر اشاره کردیم، یکی از این دلایل، کنترل صریح روی تنک بودن جواب است. در شکل ۱،۲، ما نشان می‌دهیم که، می‌توان تنک بودن جواب را با استفاده از پارامتر منظم‌سازی λ در فرمول لسو بالا، کنترل کرد. دلیل مهم دیگر، این است که، زیان مربعی دارای مزایای محاسباتی و آماری فراوانی است. همچنین برتری آماری زیان مربعی بر زیان محوری را بیان می‌کند.

در جدول ۲،۲، مقایسه‌ای از روش‌های گوناگون برای دسته‌بندی دودویی با لسو را می‌بینیم. می‌بینیم که لسو قادر به دستیابی به تنک‌ترین جواب‌ها در مسایل بعد بالاست، درحالی‌که دقت قابل قبولی را حفظ می‌نماید. در جدول ۳،۲، نتایج را برای رگرسیون مرزی بکار رفته به عنوان دسته‌بندی کننده می‌بینیم.

^۱ Spectral Projected Gradient Method for ℓ_1 -minimization

جدول ۳.۲ نتایج دقت دسته‌بندی کننده برای دسته‌بندی براساس رگرسیون با تنظیم‌سازی ℓ_2 . در اینجا λ فرایارامتر بهینه است که با اعتبارسنجی متقابل به دست آوردیم. توجه داشته باشید که، تعداد درآیه‌های ناصفر جواب (nz) برابر با تعداد ویژگی‌های مشاهدات است.

مجموعه داده‌ها			تعداد nz
فرایارامتر λ	دقت دسته‌بندی (%)	تعداد nz	
۳۰	۸۶	۱۴	استرالیایی
۶	۸۷	۲۰۰۰	سرطان کولون
۴۰	۷۶	۸	دیابت
۹	۸۶	۱۳	قلب
۱۰	۷۴	۳۴	یون کره
۸	۳۵	۶	تجویزات کبدی
۲۰	۴۹	۱۱۲	قارچ‌ها

۳

فصل سوم

کمینه‌سازی زیان مربعی ℓ_1 -منظم برای بازسازی

طرح‌های تقریب تنک برای نمایش سیگنال و ویژگی (یعنی بازسازی)، بسیار مفید شناخته شده‌اند. برای مثال، الگوریتم دسته‌بندی چند کلاسه موفق بر مبنای نمایش‌های پراکنده وجود دارد [۵۵]. در این فصل، به کارایی بازسازی (در تعریف غیر نظارتی) الگوریتم‌های کمینه‌سازی زیان مربعی ℓ_1 -منظم شده نگاهی می‌اندازیم.

همان گونه که در فصل ۲ دیدیم، ℓ_1 -منظم‌سازی برای تعمیم به داده‌های دیده نشده، بخوبی ℓ_2 -منظم‌سازی نیست. با این حال، همان‌طور که آزمایش‌های در تقریب تنک نشان می‌دهد [۱۲]، [۵۵]، [۶۱]، ℓ_1 -منظم‌سازی برای نمایش ویژگی، کاملاً مناسب است. چه چیزی روش‌های تقریب تنک را در یادگیری و نمایش ویژگی، موفق می‌سازد؟ آیا موفقیت برخاسته از تنک نمایش است یا این واقعیت که، آن داده‌ها را بسیار خوب برازش می‌نماید؟ آنالیزهایی روی پایداری الگوریتم‌های ℓ_1 -منظم شده صورت پذیرفته است [۴۰]، [۶۰].

در ادامه این فصل، به سراغ کاربرد کمینه‌سازی مربعی ℓ_1 -منظم شده برای یادگیری ویژگی در یک کار دسته‌بندی تصویر می‌رویم. بویژه به کار [۱۲] می‌پردازیم. سپس، الگوریتم دسته‌بندی را با روش جدیدمان برای رگرسیون که توسیعی از SRC به رگرسیون است، بررسی می‌کنیم. درباره این بحث می‌کنیم که، این روش‌ها بهبودی را نسبت به روش‌های موثری نظیر k -میانگین ها و k NN ارائه نمی‌کنند.

۱.۳ کدگذاری تنک و یادگیری فرهنگ لغات برای یادگیری ویژگی

برای بیان اهمیت کمینه‌سازی زیان مربعی ℓ_1 -منظم شده برای کدگذاری تنک، به سراغ مقاله [۱۲] می‌رویم. در این مقاله، اهمیت کدگذاری در مقابل یادگیری با کدگذاری تنک و کمی‌سازی برداری، با آموزش، مولفان یادگیری فرهنگ لغت D را معنی کرده و با کدگذاری، آنها منظور از تصویر کردن (نگاشتن) ورودی x به ویژگی f با فرهنگ لغت مفروض D را روشن می‌سازند. مولفان بحث می‌کنند که، اگر ما قادر به شکستن هر روش یادگیری ویژگی به مسیرهایی شامل آموزش و کدگذاری باشیم، می‌توانیم مسیرهای فرعی را از روش‌های یادگیری ویژگی ترکیب کرده و الگوریتم‌های موثرتری را بدون فداکردن کارایی در دسته‌بندی، بدست آوریم. آنها نتایج دسته‌بندی را با استفاده از اعتبار بخشی گذری ۵ بخشی با SVM خطی روی CIFAR-10، NOTB و Caltech 101 (۸۱/۵٪، ۹۵٪ و ۷۲/۶٪،

دوتای اولی نتایج مدرن هستند) گزارش می‌کنند [۱۲]. مولفان، با ۶ روش برای بکارگیری فرهنگ لغت و چندین روش برای کدگذاری با استفاده از فرهنگ لغت، آزمایش می‌کنند. در این جا، روش‌هایی برای یادگیری/ساکن کردن فرهنگ لغت وجود دارند (ما در جستجوی فرهنگ لغت $D \in R^{n \times d}$ هستیم، بگونه‌ای که هر اتم (ستون) دارای یک ℓ_2 -نرم واحد باشد):

۱. کدگذاری تنک^۱ (SC) با کاهش مختصاتی:

$$\min_{D, s^{(i)}} \sum \|Ds^{(i)} - x^{(i)}\|_2^2 + \lambda \|s^{(i)}\|_1$$

راهی برای حل این مساله بهینه‌سازی، جایگزین کردن کمینه‌سازی بین فرهنگ لغت D و کدهای تنک $\{s^{(i)}\}$ است (یکی ثابت نگه داشته می‌شود، درحالی‌که، تابع هدف کمینه می‌گردد، و به همین ترتیب). مولفان پارامتر λ را با کمینه‌سازی خطای اعتبار گذریش بر یک شبکه مقادیر کاندید، بدست می‌آورند.

۲. تعاقب تطابق متعامد^۲ (OMP-k): به‌صورت

$$\min_{D, s^{(i)}} \sum_i \|Ds^{(i)} - x^{(i)}\|_2^2 \quad (۱.۳)$$

و

$$\text{subject to } \|s^{(i)}\|_0 \leq k, \forall i, \quad (۲.۳)$$

است، که در آن، k کرانی پایین روی تعداد عناصر ناصفر در $s^{(i)}$ است. برای حل این مساله بهینه‌سازی، می‌توان بین D و $\{s^{(i)}\}$ دقیقاً مانند بالا، جایگذاری کرد.

- نزول مختصاتی و OMP الگوریتم‌هایی برای دستیابی به کدهای تنک یک فرهنگ لغت مفروض هستند، فرهنگ لغت از سوی دیگر می‌تواند با استفاده از کاهش گرادیان بدست آید.

- با بهینه‌سازی (۱.۳) و (۲.۳)، شما کدهای تنک را به عنوان نتیجه فرعی یادگیری فرهنگ لغت (فازهای یادگیری و کدگذاری در هم تنیده هستند) بدست می‌آورید. با این‌حال، این ما

^۱ Sparse coding

^۲ Orthogonal Matching Pursuit

را از برقرار ساختن تنها یک فرهنگ لغت بدست آمده در این گام و محاسبه کدها با دیگر ابزارها باز نمی‌دارد.

۳. ماشین بولتزمن تحدیدی تنک^۱ (RBM) و کدگذار خودکار تنک.

۴. زیر نمونه‌گیری تصادفی ماتریس داده X شامل $s^{(i)}$ نرمالیزه شده است.

۵. وزن‌های تصادفی: پرکردن فرهنگ لغت با نمونه‌گیری ستونی نرمال شده از توزیع نرمال استاندارد.

و در این جا، روش‌هایی برای کدگذاری وجود دارد؛ SC: مساله بهینه‌سازی در بالا، با D ثابت، λ متفاوت و تنظیم کردن عناصر ویژگی f

به صورت

$$f_j = \max\{0, s_j\} \quad (3.3)$$

و

$$f_{j+d} = \max\{0, -s_j\}. \quad (4.3)$$

است، توجه داشته باشید که بجای d بعد، ویژگی f دارای بعد $2d$ است. مولفان این را تقسیم قطبیت^۲ می‌نامند.

۱. OMP-k: تعاریف مانند ۱ است.

۲. آستانه‌سازی نرم: برای آستانه ثابت α ، آنها f را به صورت زیر نسبت می‌دهند،

$$f_j = \max\{0, D^{(j)T} x - \alpha\} \quad (5.3)$$

و

$$f_{j+d} = \max\{0, -D^{(j)T} x - \alpha\}. \quad (6.3)$$

¹ Sparse restricted Boltzmann machine

² Polarity splitting

۳. کدگذاری طبیعی: اگر فرهنگ لغت با SC فراگرفته شود، سپس، کدهای یادگیری بکار گرفته می‌شوند. این بخش برای OMP یکسان است. برای RBM و کدگذار خودکار، می‌توان فعال‌سازی در گره‌های مخفی را با استفاده از تابع حلقوی لجستیکی g محاسبه کرد:

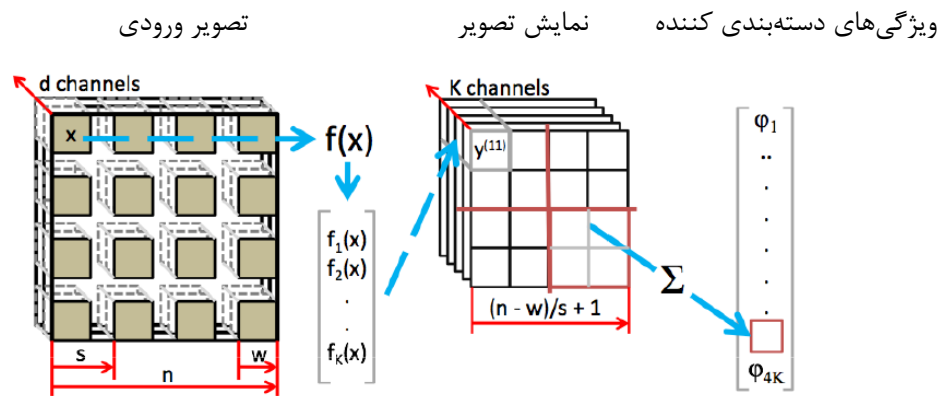
$$f_j = g(W^{(j)}x + b) \quad (۷.۳)$$

و

$$f_{j+d} = g(-W^{(j)}x + b), \quad (۸.۳)$$

که در آن $W = D^T$ و W^j ردیف j ام از W است. برای (۵.۳) و (۶.۳)، مولفان، فرهنگ لغت را به عنوان نگاشتی خطی بکار می‌برند، یعنی $f = D^T x$ (این مشابه تصویر تصادفی است، بجز این که بجای کاهش بعد، آن را افزایش می‌دهیم، با فرض $d > n$).

مولفان، بهترین نتایج را روی CIFAR-10 با استفاده از OMP-1 برای آموزش و آستانه‌سازی نرم برای کدگذاری بدست می‌آورند (که نشان‌دهنده فرهنگ لغت وسیع‌تر و بهتری است - $d=600$). آنها بهترین نتیجه را برای NORB با استفاده از وصله‌های تصادفی به عنوان فرهنگ لغت و SC برای کدگذاری، بدست آوردند. نتایج یکسانی برای Caltech 101 بدست می‌آید، هر چند این نتایج به میزان ۳/۱٪ حالت جدید را دنبال می‌کند.



شکل ۱.۳ مراحل دسته‌بندی تصویر [۱۲].

در اینجا، استفاده مولفان از مجموعه داده‌ها در فاز یادگیری بدون نظارت را بررسی می‌کنیم. در حالت CIFAR و NORB، آنها تنظیم کردند که $x^{(i)} \in R^n$ که به صورت تصادفی انتخاب شده، نرمال شود و در وصله‌های $3 \times (6 \times 6)$ بردارسازی گردد. بمانند Caltech 101، $s^{(i)}$ یک توصیف گر SIFT با ۱۲۸ بعد استخراج شده از هر وصله 16×16 تصادفی است. قبل از ارسال به الگوریتم آموزش فرهنگ لغت، آنها ZCA-سفیدسازی را روی تمام مجموعه داده $X = [x^{(1)}, \dots, x^{(1600)}]$ ارسال می‌کنند.

با فرض نگاشت ویژگی پارامتری شده با D، به بررسی خطوط لوله‌ای می‌پردازیم که، مولفان برای اجرای دسته‌بندی راه‌اندازی می‌کنند. نخست، آنها وصله‌های $\{s^{(i)}\}$ را (با اندازه مشخص شده در بالا) با یک جابجایی یک پیکسلی برای CIFAR-10 و NORB و ۸ پیکسل برای Caltech 101 استخراج می‌کنند تا تمام تصویر را پوشش دهند. برای CIFAR-10 و NORB، $s^{(i)}$ ها مقادیر پیکسل خام برای وصله هستند، درحالی‌که برای Caltech 101، آنها مقادیر توصیفگر SIFT تکی استخراج شده از وصله می‌باشند. برای هر زوج از روش‌های آزمایشی/کدگذاری، مولفان از فرهنگ لغت D برای رسیدن به ویژگی $f^{(i)}$ برای هر $s^{(i)}$ استفاده می‌کنند. برای مثال، برای هر تصویر 32×32 پیکسلی در مجموعه داده CIFAR-10، ما (در تنظیمات مشخص داده شده) بعد $2 \times 1600 \times 27 \times 27 = 2332800$ را دریافت می‌کنیم! برای کاهش بعد فضای ویژگی، یک گام تجمیع اجرا می‌شود (که برای هر مجموعه داده متفاوت است):

- CIFAR-10: مولفان مقادیر ویژگی را بر چهار ربع تصویر میانگین می‌گیرند، تا به بردار ویژگی نهایی نشان‌دهنده آن تصویر برسند (با این کار ما به بعد $2 \times 1600 \times 4$ می‌رسیم).
 - NORB: مولفان دو مرحله از نمونه‌گیری پایین را روی تصاویر 108×108 اصلی قبل از استخراج وصله‌ها، اجرا می‌کنند. همچنین به این نکته اشاره نمی‌کنند که، استراتژی تجمیع آنها پس از نگاشت ویژگی اجرا می‌شود.
 - Caltech 101: در اینجا، مولفان تجمیع هر می‌مکانی را اجرا می‌کنند، یعنی آنها تجمیع بیشینه را روی ویژگی‌هایی بر شبکه‌های 4×4 ، سپس 2×2 و بعد 1×1 در یک نظام سلسله مراتبی اجرا می‌نمایند. آنها نتایج را برای تشکیل بردار ویژگی نهایی نشان‌دهنده تصویر، الحاق می‌کنند.
- اکنون، برای در اختیار داشتن بردار ویژگی تکی برای هر مجموعه آزمون و آموزش، مولفان یک SVM خطی را برای دسته‌بندی تمرین می‌دهند.

به نظر ما نتایج این مقاله جالب است، اما چندان شگفت انگیز نیست. طرح PCA و تصادفی را در نظر بگیرید (هر چند مانند روش‌های بکار رفته در این مقاله، منجر به فرهنگ لغات وسیعی نمی‌شوند). می‌دانیم که (در پاره ای امور) وزن‌های متعامد یک تصادفی با بخش‌های آگاه از داده و یادگرفته شده یعنی مولفه‌های اصلی، قابل مقایسه هستند. نتایج همچنین، توضیح می‌دهند که چرا، الگوریتم k - میانگین در مقاله پیشین آنها [۱۳] -بسیار دورتر از طرح یادگیری اتم‌های فرهنگ لغت- بکار گرفته شده است.

۲.۳ دسته‌بندی بازنمایی تنک

برای ارائه پیش زمینه روشمان برای رگرسیون در بخش بعدی، گذری بر دسته‌بندی بازنمایی تنک خواهیم داشت [۵۵]. دسته‌بندی بازنمایی تنک^۱ (SRC) یک روش دسته‌بندی چندکلاسه است. SRC نیازمند این است که هر نمونه آزمونی با تنها چندین نمونه آزمایشی بازسازی شود. این با جستجوی بازنمایی تنک هر نمونه آزمون نسبت به فرهنگ لغت نمونه‌های آزمایشی، میسر می‌شود. در SRC، هر نمونه به یک وزن منتسب می‌شود که، برای درجه مشارکتش در بازسازی یک نمونه آزمون بحساب می‌آید. این اطلاعات امکان تصمیمی مطلع تر روی کلاس یک نمونه آزمون را فراهم می‌سازند.

مجموعه آزمایشی $\{(x_i, y_i) : i=1, \dots, n\}$ را در نظر بگیرید که در آن، $x_i \in R^p$ یک بردار ویژگی آموزشی و $y_i \in \{1, \dots, C\}$ برچسب کلاس متناظر آن است. بردارهای ویژگی $x_i, i=1, \dots, n$ ستون‌های فرهنگ لغت $D \in R^{p \times n}$ را می‌سازند. فرض می‌کنیم که، فرهنگ لغت بیش از اندازه کامل باشد ($p < n$)، در غیر این صورت، نخست کاهش بعد را روی ویژگی‌ها اجرا می‌کنیم. علاوه بر این، فرض می‌کنیم که هر ستون به نرم واحد، نرمال شده است. برای هر نقطه آزمون x ، یک نمایش پراکنده نسبت به فرهنگ لغت نقاط آزمایشی را جستجو می‌کنیم. حل مساله بهینه‌سازی زیر، منجر به نمایش a برای نقطه آزمون x می‌شود،

¹ Sparse Representation Classification

$$\min_{a \in \mathbb{R}^n} \frac{1}{2} \|x - Da\|_2^2 + \lambda \|a\|_1. \quad (9.3)$$

براساس مقادیر در بردار a ، یک تصمیم می‌تواند درباره برچسب کلاس x ساخته شود. فرض کنید $I_c \in \{1, \dots, n\}$ اندیس‌های ستون‌های D باشد که، متناظر با نقاط متعلق به کلاس c هستند. علاوه بر این، فرض کنید بردار $\delta_c(a)$ برابر با a در اندیس‌های I_c و در غیر این صورت صفر باشد. یک نمونه آزمون x به کلاسی یکسان به عنوان نمونه‌هایی که ترکیب خطی آنها بهترین بازسازی از x در معنای کمترین مربعات است، نسبت داده می‌شود، یعنی

$$\hat{c} = \arg \min_{a \in \{1, \dots, C\}} \|x - D\delta_c(a)\|_2. \quad (10.3)$$

در تعریف رگرسیون، خروجی‌های $y_i, i=1, \dots, n$ اعدادی حقیقی هستند. فرض کنید $\{y_1, \dots, y_n\}$ برداری شامل خروجی‌های متناظر برای همه نمونه‌های آزمایشی در D باشد. برای یک نمونه آزمون مفروض x ، پس از حل مساله بهینه‌سازی در (9.3)، خروجی \hat{y} را براساس فرمول زیر، پیش‌بینی می‌کنیم

$$\hat{y} = y^T a. \quad (11.3)$$

این مشابه k -نزدیک‌ترین همسایگی رگرسیون با مزیت برخورداری از وزن مشارکت هر نمونه ارائه شده در a است.

۳.۳ SPARROW: رگرسیون وزن‌دار مبتنی بر تقریب تنک

این بخش بر اساس کار منتشر شده در مقاله P.Noorzad و B.L.Strum با عنوان رگرسیون با تقریب تنک داده‌ها، در مجموعه مقالات کنفرانس پردازش سیگنال اروپا، اوت ۲۰۱۲م. است (Noorzad & Sturm, 2012).

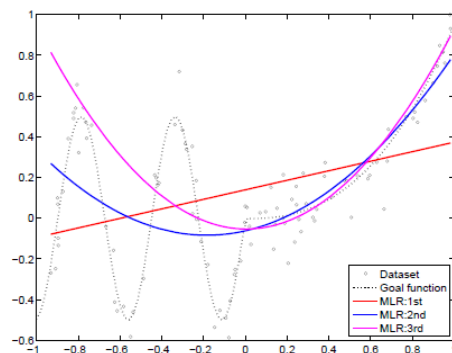
در این بخش، یک روش غیرپارامتری را برای رگرسیون چند متغیره محلی -رگرسیون وزنی تقریب تنک^۱- ارائه و مطالعه می‌کنیم که، تقریب تنک از یک نقطه برحسب متغیرهای پیش‌گوست. یک روش

¹ SPARse approximation Weighted regression

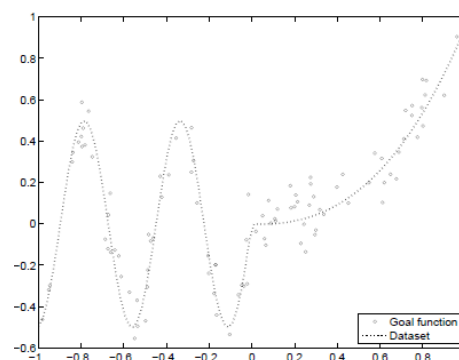
غیرپارامتری مشابه، رگرسیون k -نزدیک‌ترین همسایگی (k -NNR) است [۲۳]، که فرض می‌کند نزدیک‌ترین k متغیر پیشگو به یک نقطه آزمون، متغیرهای پاسخ یکسانی را تولید می‌کنند. هر دو روش می‌تواند به عنوان گونه‌هایی از رگرسیون هسته چندجمله‌ای محلی^۱ (LPKR) در نظر گرفته شود [۴۷] که تابع رگرسیون را در یک نقطه با برازش یک چند جمله‌ای در آن نقطه، ارزیابی می‌نماید.

علاوه بر روش‌های محلی نظیر k -NNR و LPKR، پژوهش‌های درخور توجهی با هدف روش‌های غیرپارامتری کلی، مانند مدل‌های جمعی (AMها) [۸] و مدل‌های جمعی تنک (ApAM) [۴۱] صورت پذیرفته است. در AMها، روش‌های تک متغیره برای ارزیابی یک تابع هموار از هر متغیر پیشگو-در یک مدل شامل مجموع چنین توابع مولفه‌ای تک متغیره‌ای- با اجتناب از نیاز به ارتباط مستقیم با ورودی‌های چند بعدی، بکار گرفته می‌شوند. در SpAM، هدف کاهش تعداد توابع مولفه‌ای یک مدل جمعی است [۴۱]. رگرسیون تعاقبی تصویری (PPR) [۱۹] توسیعی به AMهایی است که، قادر به مدل کردن یک کلاس کلی تر از توابع هستند.

هرچند روش‌ها برای رگرسیون پارامتری و غیرپارامتری کلی ممکن است خطای میانگین را در کل مجموعه داده‌ها کمینه سازند، اما امکان ارائه یک برازش محلی خوب را ندارند. شکل ۲.۳، توانایی روش‌های محلی در مدل‌سازی داده‌های تولید شده توسط یک توزیع نامعلوم را نشان می‌دهد.

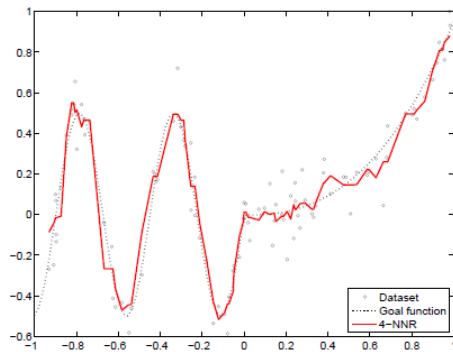


(ب) رگرسیون خطی چندگانه با جملات مرتبه ۱ تا ۳.

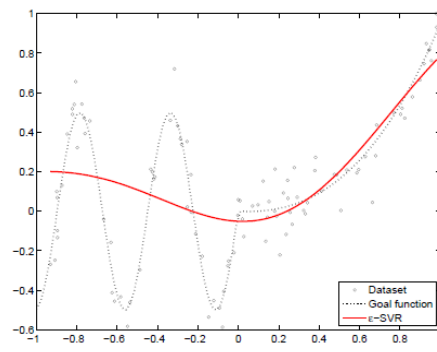


(الف) مجموعه داده تولید شده ما.

^۱ Local Polynomial Kernel Regression



(د) ۴ رگرسیون نزدیک‌ترین همسایگی.



(ج) رگرسیون برداری ϵ -پشتیبان با هسته RBF.

شکل ۲.۳ این اشکال، توانایی روش‌های رگرسیون محلی را برای مدل‌سازی داده‌ها با یک توزیع

ناشناخته، نشان می‌دهند. تابع مولد داده‌ها عبارت است از: $y_i = f(x_i) + \epsilon_i$ که

در آن $f(x) = (x^3 + x^2)I(x) + \sin(x)I(-x)$.

روش‌های محلی نظیر LPKR، بر فرض یک مدل پارامتری محلی برای داده‌ها، متکی هستند [۱۱]. در LPKP، می‌توان تابع رگرسیون را در هر نقطه با برازش یک چندجمله‌ای تیلور در آن نقطه، ارزیابی نمود. این می‌تواند مدل‌هایی را بر اساس مرتبه چندجمله‌ای تولید کند که، ثابت محلی، خطی محلی، مربعی محلی و غیره هستند. اصل این روش، کمینه‌سازی یک مجموع وزنی از خطای مربعی است. معمولاً، وزن‌ها با یک تابع نزولی از فاصله بین دو نقطه تعریف می‌شوند. SPARROW این وزن‌ها را با استفاده از تقریب تنک داده‌ای آزمون، تعریف می‌کند. این فرض برقرار است که یک نقطه آزمون با ترکیبی از پیشگوها بجای نزدیک به آنها، بهتر مدل‌سازی می‌شود.

مزایای مدل‌سازی داده‌ها با محدودیت‌های تنک بودن بخوبی مستندسازی شده است [۹]، [۱۶]، [۳۴]، به عنوان مثال در عدم پوشش کد فیزیولوژی قشر بصری اصلی پستانداران [۳۸] و در تولید کدهای تنک صداهای طبیعی [۳۱]، تصاویر [۶۱]، آواهای موسیقایی [۳۹]. در حوزه یادگیری نظارت شده، دسته‌بندی بازنمایی تنک [۵۵]، می‌تواند روش‌های استاندارد را در چارچوب‌های دشوار پیاده‌سازی نماید، برای مثال ترازبندی نادرست و تغییر درخشندگی [۵۵]. شرط تنک بودن همچنین، برای انتخاب متغیر بویژه در لسو بکار می‌رود [۴۹]. در بخش‌های بعدی، SPARROW را تعریف کرده و نشان

می‌دهیم که چگونه، گونه ای از k-NNR و LPKP است. سپس، نتایج چندین آزمایش را در مقایسه SPARROW با این‌ها و دیگر روش‌های شناخته شده، ارائه می‌کنیم.

۱.۳.۳ رگرسیون وزن‌دار مبتنی بر تقریب تنک

یک مجموعه داده (یا فرهنگ لغت) از N مشاهده را در نظر بگیرید؛ $D := \{(x_i, y_i)\}_{i \in \Omega}$ ، که ورودی $x_i = [x_{i1}, \dots, x_{Mi}]^T \in R^M$ متناظر با خروجی $y_i \in R$ است. فرض کنید $\Omega := \{1, 2, \dots, n\}$ اندیس فرهنگ لغت باشد. در رگرسیون غیرپارامتری، می‌توان فرض کرد که $y_i = f(x_i) + \varepsilon$ ، که در آن $f(x)$ نامعلوم است، اما تابعی است هموار و ε_i خطای مستقل از x_i است. با D و نقطه z مفروض، SPARROW تابع رگرسیون $f(z)$ را با ترکیب خطی خروجی‌های

$$\hat{f}(z) := \sum_{i \in \Omega} l_i(z, D) y_i \quad (12.3)$$

ارزیابی می‌کند که در آن، $l_i(z, D)$ وزن موثر نام است که، SPARROW را به عنوان تابعی از تقریب تنک z در D انتخاب می‌کند.

بجای برازش یک مدل تکی برای همه مجموعه داده‌ها، مانند رگرسیون پارامتری و غیرپارامتری کلی، SPARROW مدل‌های پارامتری درباره هر نقطه آزمون z را با استفاده از مثلاً یک بسط تیلور مرتبه صفر، یک یا دو، برازش می‌نماید. اکنون، بحث می‌کنیم که، چگونه SPARROW وزن‌های موثر را در (۱۲.۳) برای ارزیابی تابع رگرسیون در یک نقطه مفروض، تعریف می‌کند.

۲.۳.۳ تعریف وزن‌های موثر

برای دستیابی به ارزیابی مربعی تابع رگرسیون در z ، می‌توانیم $f(x)$ را در z با یک چندجمله‌ای تیلور از درجه دو تقریب بزنیم

$$f(x) \approx f(z) + (x - z)^T \theta_z + \frac{1}{2} (x - z)^T H_z (x - z) \quad (13.3)$$

که در آن، $\theta_z := \nabla f(z)$ گرادیان $f(x)$ و $H_z := \nabla^2 f(z)$ هسی آنست که هر دو در z ارزیابی می‌شوند. اکنون، مساله یافتن $f(z), \theta_z, H_z$ بگونه‌ای است که، خطای مربعی وزنی محلی را در حوالی z برای اندازه‌گیری‌هایی در D کمینه‌سازی نماییم، یعنی

$$\min_{f(z), \theta_z, H_z} \sum_{i \in \Omega} \alpha_i(z) \left[y_i - f(z) - (x_i - z)^T \theta_z - \frac{1}{2} (x_i - z)^T H_z (x_i - z) \right]^2 \quad (14.3)$$

جایی که $\alpha_i(z)$ وزن مشاهده نام است که، می‌تواند به چند طریق مثلا، با یک تابع هسته [23]، [25] یا با تقریب تنک صورت گرفته با SPARROW تعریف شود.

اکنون، ابر بردار پارامتری زیر را تعریف می‌کنیم [47]

$$\Theta_z := [f(z), \theta_z, \text{vech}(H_z)]^T \quad (15.3)$$

که در آن $\text{vech}(H)$ نشان‌دهنده بردارسازی ماتریس $M \times M$ متقارن است، یعنی $M \times (M+1)/2$ بردار تشکیل شده با افزودن درآیه‌های قطری و مثلثی پایین H_z . ماتریس قطری A_z را تعریف کنید که، عنصر قطری i امش، وزن مشاهده $\alpha_i(z)$ است. با تعریف ماتریس

$$X_z := \begin{bmatrix} 1 & (x_1 - z)^T & \text{vech}^T[(x_1 - z)(x_1 - z)^T] \\ \vdots & \vdots & \vdots \\ 1 & (x_N - z)^T & \text{vech}^T[(x_N - z)(x_N - z)^T] \end{bmatrix} \quad (16.3)$$

می‌توانیم کمینه‌سازی در (14.3) را به صورت

$$\min_{\Theta_z} \|A_z^{1/2} [y - X_z \Theta_z]\|_2^2 \quad (17.3)$$

بیان کنیم، که بردار پاسخ $y = [y_1, \dots, y_N]^T$ است. پارامترهای تعریف شده با جواب کمترین مربعات عبارتند از [47]

$$\hat{\Theta}_z = (X_z^T A_z X_z)^{-1} X_z^T A_z y \quad (18.3)$$

که $X_z^T A_z X_z$ معکوس پذیر است. در نهایت، ارزیابی مربعی محلی تابع رگرسیون در z دقیقا عنصر نخست Θ_z ، یعنی

$$\hat{f}(z) = e_1^T (X_z^T A_z X_z)^{-1} X_z^T A_z y = \sum_{i \in \Omega} \beta_i y_i \quad (19.3)$$

است، که e_1 دارای یک در نخستین سطرش و صفر در جاهای دیگر است. پس، می‌بینیم که λ مین وزن موثر در (۱۲،۳) عبارت است از

$$l_i(z, D) = e_i^T A_z^T X_z (X_z^T A_z X_z)^{-1} e_1 \quad (20.3)$$

به‌صورت خلاصه، SPARROW تابع رگرسیون را در نقطه z با محاسبه (۱۲،۳) با وزن‌های موثر مفروض $(3, 20)$ ، ارزیابی می‌کند. اگر تنها از نخستین ستون X_z در (۲۰،۳) استفاده نماییم، یک ارزیابی ثابت محلی از $f(z)$ یعنی

$$\hat{f}(z) = (1^T A_z 1)^{-1} 1^T A_z y = \frac{\sum_{i \in \Omega} \alpha_i(z) y_i}{\sum_{k \in \Omega} \alpha_k(z)}. \quad (21.3)$$

تولید می‌شود، با استفاده از $M+1$ ستون از X_z یک ارزیابی خطی محلی را تولید می‌کنیم و با استفاده از همه X_z ها به یک ارزیابی مربعی محلی می‌رسیم. با استفاده از چندجمله‌ای‌های مرتبه بالاتر مانند مدل پارامتری محلی، انحراف ارزیابی را کاهش می‌دهیم [۲۵]، [۴۷]، اما این کار به بهای واریانس و افزایش زمان محاسباتی صورت می‌گیرد، زیرا تعداد پارامترهای محلی به‌صورت نمایی افزایش می‌یابد. علاوه بر این، چندجمله‌ای‌های مرتبه بالاتر، بهبود نه چندان چشمگیری را بر مدل مربعی ارائه می‌کنند، مگر زمانی که به دنبال ارزیابی گرادیان و هسی یعنی θ_z, H_z در (۱۴،۳) هستیم [۴۶].

۳.۳.۳ تعریف وزن‌های شهودی

چون وزن‌های موثر در (۲۰،۳)، تابعی از وزن‌های مشاهده در (۱۴،۳)، یعنی $\{\alpha_i(z) : i \in \Omega\}$ هستند، مساله باقی‌مانده، تعریف وزن‌های مشاهده است. اگر آنها را به‌صورت محلی در مدل ثابت (۲۱،۳) با یک تابع هسته، تعریف کنیم، ارزیابی رگرسیون نادارایی-واتسون^۱ (NWR) را تولید می‌کنیم. در این راستا، می‌توانیم وزن‌ها را با

$$\alpha_i(z) := K(S(z, x_i) / h) \quad (22.3)$$

تعریف کنیم، که در آن $K: R \rightarrow R_+$ یک تابع هسته، $h > 0$ پهنای باند و $S(z, x_i)$ فاصله

^۱ Nadaraya-Watson regression (NWR)

$$S(z, x_i) := (z - x_i)^T V^{-1} (z - x_i) \quad (23.3)$$

است، که V یا یک ماتریس قطری از ارزیابی های غیر انحرافی واریانس های مشاهده شده در ابعاد پیشگو در D (در این حالت $(23,3)$ فاصله اقلیدسی مدرج است)، یا ارزیابی غیرانحرافی کوواریانس پیشگو است (در این حالت $(23,3)$ فاصله ماهالانوبیس^۱ مد نظر است).

زمانی که وزن های مدل ثابت را به صورت محلی، تعریف می کنیم

$$\alpha_i(z) := \begin{cases} d(z, x_i), & i \in N_k(z) \subset \Omega \\ 0, & \text{else} \end{cases} \quad (24.3)$$

که در آن، $N_k(z)$ مجموعه اندیس k نزدیک ترین پیشگوی z در D است، سپس $(21,3)$ k -NNR را تولید می کند [23]. اگر $d(z, x_i) = 1$ باشد، سپس پهنای باند هسته ثابت از z حداقل به بزرگی بزرگ ترین فاصله بین زوج های مشاهدات یعنی، $h \geq \max_{i \in N_k(z)} S(z, x_i)$ است. در k -NNR (Wk-NNR)، این وزن را به عنوان معکوس فاصله $d(z, x_i) = 1/S(z, x_i)$ تعریف می کنیم.

در تقابل با NWR و k -NNR، در عوض، SPARROW وزن های مشاهده را از تقریب تنک z در D تعریف می کند. نخست، فرم ماتریسی از پیشگوی نرمال شده فرهنگ لغت را در نظر بگیرید

$$D := \left[\frac{x_1}{\|x_1\|_2}, \frac{x_2}{\|x_2\|_2}, \dots, \frac{x_N}{\|x_N\|_2} \right]. \quad (25.3)$$

برای یک ورودی z ، SPARROW جوابی را بر $z \approx Ds$ بگونه ای می یابد که $s = [s_1, s_2, \dots, s_N]^T$ دارای تعداد زیادی عنصر صفر باشد. راه های متنوعی برای تولید تقریب های تنک وجود دارد (برای بررسی بیشتر [6]، [16]، [51] را ببینید). در این کار، ما از اصل نویززدایی تعاقبی پایه ای^۲ (BPDN) [9] استفاده می کنیم که دارای مساله

$$\min_{s \in \mathbb{R}^N} \|s\|_1 \quad \text{subject to} \quad \frac{\|z - Ds\|_2^2}{\|z\|_2^2} \leq \varepsilon^2 \quad (26.3)$$

¹ Mahalanobis distance

² Basis Pursuit DeNoising

است که در آن، $\varepsilon^2 > 0$ سیگنال را به نسبت خطای تقریب محدود می‌سازد. در نهایت، SPARROW وزن مشاهده \mathbf{A} را با استفاده از وزن‌های تقریب تنک تعریف می‌کند

$$\alpha_i(\mathbf{z}) := \left[\frac{S(\mathbf{z}, \mathbf{x}_i)}{\min_{j \in \Omega} S(\mathbf{z}, \mathbf{x}_j)} \right]^{-1} \frac{s_i}{\|\mathbf{z}\|_2} \quad (27.3)$$

که در آن، s_i عنصر \mathbf{A} است. هدف ضریب نخست، وزن کردن با یک متغیر پاسخ نزدیک‌ترین پیشگو به \mathbf{z} است؛ و هدف تقسیم وزن تقریب تنک با $\|\mathbf{z}\|_2$ ، حذف تاثیر طول آنست. بنابراین، اگرچه مشابه با Wk-NNR، SPARROW، Wk-NNR بسختی یک مشاهده نزدیک تر به جستجو را وزن می‌کند، برخلاف Wk-NNR، این کار را زمانی انجام می‌دهد که، دارای ضریب ناصفری در تقریب تنکش با D باشد. زمانی که وزن‌های $\alpha_i(\mathbf{z})$ را از (27,3) در (21,3) جایگزین می‌کنیم، ارزیابی SPARROW (C-SPARROW) را بدست می‌آوریم، و زمانی که از این وزن‌ها را در (20,3) با تنها $M+1$ ستون نخست از \mathbf{X}_z بکار می‌بریم، (12,3) ارزیابی SPARROW خطی (L-SPARROW) را تولید می‌نماید. با استفاده از ستون‌های \mathbf{X}_z ارزیابی SPARROW مربعی (Q-SPARROW) را تولید می‌کنیم.

4.3 ارزیابی آزمایشی SPARROW

ما اکنون، کارایی SPARROW را در برابر چندین روش شناخته شده دیگر برای رگرسیون محلی، مقایسه می‌کنیم. در همه حالات، از فاصله اقلیدسی استاندارد شده در (23,3) استفاده می‌کنیم. NWR و همتای خطیش، رگرسیون هسته خطی محلی (LLKR) [23]؛ [25] را می‌آزماییم، که (18,3) را با استفاده از $M+1$ ستون نخست \mathbf{X}_z در (16,3) حل می‌کند. برای هر دوی NWR و LLKR، هسته گاوسی را در (22,3) به صورت زیر اتخاذ می‌کنیم

$$K(x) := \frac{1}{\sqrt{2\pi}} e^{-x^2/2}. \quad (28.3)$$

همچنین k -NNR، Wk -NNR [۲۳] را می‌آزماییم، که در آن، k را با اعتبارسازی گذری تودرتو، تنظیم می‌کنیم. برای یک خط مبنا، روش پارامتری کلی رگرسیون خطی چندگانه^۱ (MLR) [۲۴]، را آزمایش می‌کنیم، که در آن فرض بر آنست که، یک فرم خطی از تابع رگرسیون

$$f(x) = [1, x^T]b \quad (۲۹.۳)$$

و b برای کمینه‌سازی خطای مربعی میانگین

$$b = \frac{\arg \min_{b' \in \mathbb{R}^{M+1}} \|y - [1, X^T]b'\|_2^2}{\quad} \quad (۳۰.۳)$$

تعریف شده است، که در آن، ستون i ام X برابر با x_i است. برای تولید تقریب تنک برای یک نقطه آزمون در (۲۶،۳)، از روش گرادیان تصویری طیفی برای ℓ_1 -منظم‌سازی (SPGL1) [۵۲] با حداکثر ۲۰ تکرار و $\varepsilon = 10^{-6}$ استفاده می‌کنیم.

در این جا، ۴ مجموعه داده متمایز را به‌صورت مشترک در رگرسیون بکار می‌بریم (جدول ۳،۱ را ببینید).^۲

جدول ۱.۳ اطلاعات مجموعه داده‌ها. آخرین ستون پارامتر k تزار شده در آزمایش‌های مربوط به k -NNR و Wk -NNR را نشان می‌دهد.

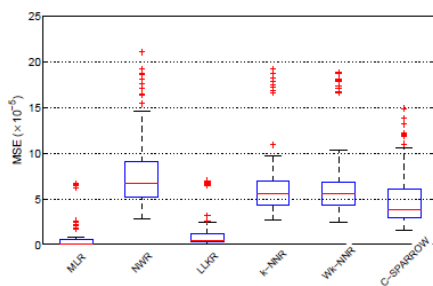
مجموعه داده‌ها	تعداد مشاهدات (N)	تعداد گرایش‌ها (M)	k
abalone	۴۱۷۷	۸	۹
bodyfat	۲۵۲	۱۴	۴
housing	۵۰۶	۱۳	۲
mpg	۳۹۲	۷	۴

^۱ Multiple Linear Regression

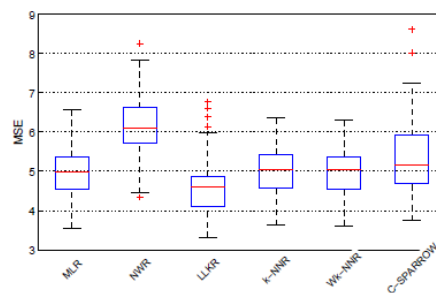
^۲ mpg, abalone and housing are from <http://archive.ics.uci.edu/ml/>; bodyfat is from <http://lib.stat.cmu.edu/dataset/>.

بجز برای bodyfat، هر مجموعه داده را بگونه‌ای استانداردسازی کردیم که، بعدشان صفر-میانگین و دارای واریانس یکسان باشد. شکل ۳،۳، ارزیابی های خطای مربع میانگین^۱ (MSE) از این الگوریتم‌ها را از ۱۰ دنباله مستقل از اعتبار گذری ۱۰ تایی، نشان می‌دهد. می‌بینیم، در حالی که MLR بخوبی برای bodyfat و abalone اجرا می‌شود، برای mpg و housing بسیار ضعیف عمل می‌کند. از سوی دیگر، LLKR بخوبی برای همه مجموعه‌های داده، اجرا می‌شود. این مطلب در کارایی حاصل از افزایش در محاسبات صورت گرفته توسط LLKR از (۱۸،۳) مشهود است. بجز برای abalone و housing، می‌بینیم که، C-SPARROW تقریباً مشابه با k-NNR و Wk-NNR اجرا می‌شود. این شگفت‌انگیز است، زیرا؛ (۱) C-SPARROW هیچ فرضی از تعداد همسایگی های بکار رفته برای هر نقطه آزمون را، مطرح نمی‌کند، و (۲) یک ارزیابی ثابت محلی را می‌سازد.

جدول ۲،۳، کارایی L-SPARROW را در مقایسه با C-SPARROW نشان می‌دهد. می‌توانیم انتظار داشته باشیم که L-SPARROW بهتر از C-SPARROW عمل کند، زیرا مدل مرتبه بالاتری است. با این حال، یک مساله با رگرسیون چندجمله‌ای محلی برای چندجمله‌ای های مرتبه بالاتر (یعنی مرتبه اول یا دوم) است، آن هم زمانی که ورودی به صورت محلی، رتبه ناقص است، و جواب بر (۱۸،۳) را ناپایدار می‌سازد.

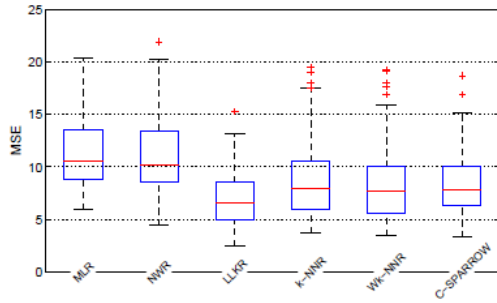


(ب) مجموعه داده bodyfat

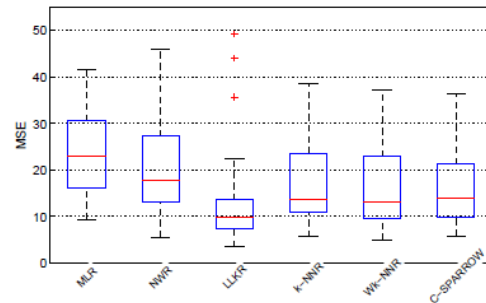


(الف) مجموعه داده abalone

¹ Mean Squared Error



mpg مجموعه داده (د)



housing مجموعه داده (ج)

شکل ۳.۳ ترسیم جعبه ای برای ارزیابی اعتبار سنجی متقابل ۱۰-تایی از خطای مربع میانگین (۱۰۰ بار اجرای مستقل) برای ۴ مجموعه داده گوناگون. هر جعبه ۲۵ تا ۷۵ درصد را معین می‌کند و خط قرمز، میانه را نشان می‌دهد. اکسترم با علامت + و برون هشته ها با ضرب‌آنها مشخص شده است.

مساله را با حل یک فرم منظم شده از بهینه‌سازی کمترین مربعات وزن شده در $(14,3)$ حل می‌کنیم. از ℓ_2 -نرم پارامترهای محلی به عنوان جمله منظم‌سازی استفاده می‌کنیم، و یک مساله رگرسیون مرزی را حل می‌کنیم [۲۸]، یعنی بجای حل $(17,3)$ ، مساله

$$\min_{\Theta_z, \lambda} \|A_z^{1/2}[y - X_z \Theta_z]\|_2^2 + \lambda \|\Theta_z\|_2^2 \quad (31.3)$$

را حل می‌کنیم، که در آن $\lambda \geq 0$ پارامتر مرزی^۱ است. برای یک λ داده شده، جواب بفرم [۲۴]

$$\hat{\Theta}(z) = (X_z^T A_z X_z + \lambda I)^{-1} X_z^T A_z y. \quad (32.3)$$

در می‌آید. ما λ را به همان طریق توصیف شده در بالا برای k تنظیم می‌کنیم. با این وجود، در حالی که می‌بینیم، کارایی L-SpARROW نسبت به استفاده از $(18,3)$ بهبود یافته است، همچنان تا سطح C-SPARROW پایین باقی می‌ماند.

¹ Ridge parameter

جدول ۲.۳ مقایسه‌ای از ارزیابی‌های MSE براساس ۴ مجموعه داده با ۱۰ توالی اعتبار گذری ۱۰ تایی از C-SPARROW و L-SPARROW با و بدون منظم‌سازی. آخرین ستون، نشان‌دهنده پارامتر مرزی بکار رفته برای دستیابی به ارزیابی L-SPARROW است.

λ	L-SPAR	L-SPAR. w/R	C-SPAR	مجموعه داده‌ها
10^{-3}	۹۸۸	۱۶	۵	abalone
10^{-6}	960×10^{-5}	5×10^{-5}	35×10^{-5}	bodyfat
10^{-4}	۴۳۰۵	۴۵	۱۰	housing
10^{-3}	۶۳۳۵	۸	۷	mpg

۱.۴.۳ نتیجه‌گیری

در این کار، انواعی سازگار از روش‌های رگرسیون چندجمله‌ای محلی را ارائه نمودیم: LLKR، NWR، k-NNR و Wk-NNR. LLKR و NWR از مجموعه داده، وزن، و متغیر پاسخ هر پیشگو با یک تابع هسته، استفاده می‌کنند. در عوض، k-NNR و Wk-NNR از متغیرهای پاسخ k نزدیک‌ترین پیشگو به نقطه بهره می‌گیرند، تا به صورت محلی تابع رگرسیون را ارزیابی نمایند. با SPARROW، ما استفاده از تقریب تنک برای انتخاب سازگاری را ارائه نمودیم که پیشگو را بکار می‌برد، و وزن‌های متغیرهای پاسخ برای ارزیابی تابع رگرسیون در نقطه داده شده، بکار می‌روند. آزمایش‌های ما، نشان می‌دهد که، SPARROW ثابت، می‌تواند یک الگوریتم رگرسیون رقابتی باشد. کار آینده ما وضعیت‌هایی را آنالیز می‌کند که، موجب توصیف داده‌ها به عنوان ترکیبی خطی (شامل وزن‌های منفی) از داده‌های برچسب گذاری شده می‌شود. علاوه بر این، می‌توانیم، دیگر الگوریتم‌های تقریب تنک از قبیل روش‌های حریصانه^۱ را بکار ببریم، که معمولاً، از نظر محاسباتی، کم هزینه‌تر از روش‌های بهینه‌سازی محدب از قبیل BPDN هستند.

^۱ Greedy approaches

۵.۳ مقایسه kNN با SRC

در جدول ۳،۳، مقایسه‌ای را از الگوریتم دسته‌بندی k نزدیک‌ترین همسایگی و دسته‌بندی بازنمایی تنک روی مجموعه داده‌های دسته‌بندی چند کلاسه، ارائه می‌کنیم. مجموعه داده‌ها از منزلگاه LIBSVM برای مجموعه داده‌های چندکلاسه انتخاب شده است.^۱ می‌بینیم که، دسته‌بندی بازنمایی تنک هیچ بهبود قابل توجهی را بر دسته‌بندی k نزدیک‌ترین همسایگی روی اکثریت مجموعه داده‌ها، ارائه نمی‌نماید.

جدول ۳.۳ مقایسه دقت بدست آمده توسط kNN و SRC روی ۵ مجموعه داده دسته‌بندی چندکلاسه.						
SRC	kNN	k	تعداد کلاس‌ها	p	n	مجموعه داده‌ها
۸۶	۸۶	۱۲۵	۳	۱۸۰	۲۰۰۰	دی ان آ
۶۵	۷۰	۲	۶	۹	۲۱۴	شیشه
۷۲	۹۵	۶	۳	۴	۱۵۰	عنابیه
۸۴	۹۴	۲	۱۱	۱۰	۵۲۸	مصوته
۹۹	۹۷	۷	۳	۱۳	۱۷۸	شراب

^۱ <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multiclass.html>

فصل چهارم

یک هم ارزی بین ϵ -SVR و BPDN

در مقاله ای با عنوان "یک هم ارزی بین تقریب تنک و ماشین های بردار پشتیبان"،¹ Girosi مساله برنامه ریزی مربعی رگرسیون بردار پشتیبان¹ (هسته) (SVR) را در چارچوب تئوری منظم سازی، نتیجه گرفت. این مطلب در تضاد با انحراف اصلی Vapnik از SVR با استفاده از اصل کمینه سازی ریسک ساختاری در تئوری یادگیری آماری است. Girosi نشان می دهد که، تحت شرط ها و فرضیات فرعی، روش نوپزدایی تعاقبی پایه ای (BPDN) منجر به مساله برنامه ریزی مربعی مشابهی می شود. Chen, Donoho و Saunders، نوپزدایی تعاقبی پایه ای را به عنوان یک روش مهارشدنی برای تقریب تنک ارائه نمودند.

۱.۴ فضای هیلبرت هسته تکثیری

فضای ضرب داخلی H را روی مجموعه X بر R در نظر بگیرید. اگر H نسبت به متریک القایی توسط ضرب داخلی، کامل باشد، سپس H یک فضای هیلبرت است. ما ضرب داخلی بردارهای $f, g \in H$ را با $\langle f, g \rangle$ نشان می دهیم.

تعریف ۱.۱،۴. می گوییم H یک فضای هیلبرت تکثیری است، اگر برای هر $x \in X$ ، تابع خطی $F_x: H \rightarrow R$ که در آن $F_x(f) = f(x)$ است، کران دار باشد.

چون F_x یک تابع خطی کران دار است، از قضیه نمایش ریس، در می یابیم که، برای هر $x \in X$ یک $K_x \in H$ یکتا با ویژگی بازیابی

$$F_x[f] = \langle K_x, f \rangle = f(x) \quad (۱.۴)$$

وجود دارد، که در آن، f تابعی در H است. برای $t \in X$ ، یک $K_t \in H$ وجود دارد و با ویژگی بازیابی در معادله (۱.۴)، داریم

$$K_t[x] = \langle K_x, K_t \rangle. \quad (۲.۴)$$

¹ Support Vector Regression

برای همه $y \in X$ تابع $K : X \times X \rightarrow R$ را به صورت

$$K(x, y) = K_y(x) \quad (۳.۴)$$

تعریف می کنیم و آن را هسته بازیابی^۱ (RK) می نامیم. از یکتایی K_y در می یابیم که K کاملاً توسط H تعیین می شود. قضیه زیر، ارتباطی را بین یک RK و هسته هیلبرت بازیابی متناظرش^۲ $(RKHS)$ برقرار می سازد:

قضیه ۱، ۲، ۴. برای هر $RKHS$ یک تابع معین مثبت متقارن یکتا (RK) وجود دارد و برعکس برای هر تابع معین مثبت متقارن K روی $X \times X$ ، یک $RKHS$ یکتا از توابع روی X با K به عنوان خود RK موجود است.

اثبات را می توانید در فصل نخست کتاب Wahba [۵۶] درباره اسپلاین ها بیابید. ما تنها بسراغ ساختارهای اثبات می رویم. در حالت خاص، اگر H یک $RKHS$ باشد، سپس RK عبارت است از

$$K(x, y) = \langle K_x, K_y \rangle \quad (۴.۴)$$

که در آن برای هر $x, y \in X$ ، توابع نماینده یکتا همان K_x, K_y هستند. برعکس، برای K مفروض، برای همه $x, y \in X$ ، ما تابع نماینده را به صورت

$$K_x(y) = K(x, y). \quad (۵.۴)$$

تعریف می کنیم. $RKHS$ را به عنوان مکمل فضای توابع H_0 تولید شده توسط $\{K_x, x \in X\}$ می سازیم، البته با ضرب داخلی تعریف شده به صورت

$$\langle f, g \rangle = \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j K(x_i, y_j) \quad (۶.۴)$$

که در آن، f, g توابعی در H_0 هستند،

^۱ Reproducing Kernel

^۲ Reproducing Kernel Hilbert Space

$$f = \sum_{i=1}^m \alpha_i K_{x_i} \quad (۷.۴)$$

$$g = \sum_{j=1}^n \beta_j K_{y_j}. \quad (۸.۴)$$

از این تعریف نتیجه می گیریم

$$K(x, x_i) = \langle K_x, K_{x_i} \rangle, \quad (۹.۴)$$

و نیز

$$\begin{aligned} \langle K_x, f \rangle &= \left\langle K_x, \sum_{i=1}^m \alpha_i K_{x_i} \right\rangle \\ &= \sum_{i=1}^m \alpha_i \langle K_x, K_{x_i} \rangle \\ &= \sum_{i=1}^m \alpha_i K(x, x_i) \\ &= \sum_{i=1}^m \alpha_i K_{x_i}(x) \\ &= f(x). \end{aligned} \quad (۱۰.۴)$$

آخرین معادله نشان می دهد که هسته، نماینده ارزیابی است.

تا کنون، می دانیم که، یک RKHS دارای هسته بازیابی یکتای K است. چون K یک تابع معین مثبت است، دارای یک آنالیز طیفی به شکل

$$K(x, y) = \sum_{n=1}^{\infty} \lambda_n \phi_n(x) \phi_n(y) \quad (۱۱.۴)$$

می باشد، که در آن، Φ_1, Φ_2, \dots دنباله متعامد یکه توابع ویژه در $L_2[X]$ و $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$ مقادیر ویژه متناظر هستند. توابع ویژه از RK، H را به عنوان RKHS تولید می کند. پس، هر تابع در H را می توان به صورت

$$\hat{f}(x) = \sum_{n=1}^{\infty} c_n \phi_n(x) \quad (12.4)$$

با نرم

$$\|\hat{f}\|_H^2 = \langle \hat{f}, \hat{f} \rangle_H = \sum_{n=1}^{\infty} \frac{c_n^2}{\lambda_n}. \quad (13.4)$$

نشان داد. تابع‌هایی بشکل (13,4)، به عنوان تابع‌های هموار متناظر با مقادیر کوچک‌تر نسبت به توابع هموارتر، شناخته می‌شوند.

۲.۴ ماشین‌های بردار پشتیبان برای رگرسیون

مجموعه داده تصادفی $D = \{(x_i, y_i) : x_i \in R^d, y_i \in R, i = 1, \dots, l\}$ با نمونه‌گیری تصادفی (در غیاب نویز) از یک تابع نامعلوم f بدست آمده است. هدف ما، پوشش f یا یک ارزیابی از D است. این مساله رگرسیون، جواب‌های بسیاری دارد، زیرا، توابع فراوانی از مجموعه نقاط داده شده، می‌گذرند. مساله را با این فرض که، درمیان تمامی توابع درون یابی، جواب مساله ما دارای بیشترین همواری است، محدود می‌سازیم (که در آن، همواری شامل نقاط نزدیکی است که مقادیر نزدیکی دارند). فرض کنید $\Phi[f]$ تابعی هموار باشد. در ϵ -SVR، هدف ما یافتن تابع \hat{f} با حداکثر انحراف ϵ از داده‌های آزمایشی است، یعنی برای $i = 1, \dots, l$ همواره $|\hat{f}(x_i) - y_i| \leq \epsilon$. مساله بهینه‌سازی عبارت است از

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \Phi[\hat{f}] \\ & \text{subject to} \quad \begin{cases} y_i - \hat{f}(x_i) \leq \epsilon \\ \hat{f}(x_i) - y_i \leq \epsilon \end{cases} \quad \text{for } i = 1, \dots, l. \end{aligned} \quad (14.4)$$

متغیرهای کمبود ζ_i, ζ_i^* را در ارتباط با امکان غیرعملی بودن مساله بهینه‌سازی محدود شده در (14,4) معرفی می‌کنیم. مساله بهینه‌سازی به‌صورت

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \Phi[\hat{f}] + C \sum_{i=1}^l (\varsigma_i + \varsigma_i^*) \\ & \text{subject to} \quad \begin{cases} y_i - \hat{f}(x_i) \leq \epsilon + \varsigma_i \\ \hat{f}(x_i) - y_i \leq \epsilon + \varsigma_i^* \\ \varsigma_i \geq 0 \\ \varsigma_i^* \geq 0 \end{cases} \quad \text{for } i = 1, \dots, l \end{aligned} \quad (15.4)$$

در می‌آید، که در آن، پارامتر آزاد $C > 0$ تعامل بین همواری \hat{f} و مقدار انحراف آن را -فراتر از ϵ - از داده‌های آزمایشی، کنترل می‌کند. مساله در (15,4) یک مساله برنامه‌ریزی محدب است و بنابراین دارای یک کمینه یکتاست. با معرفی ضرایب لاگرانژ (متغیرهای دوگان) برای افزودن محدودیت‌ها به تابع هدف، لاگرانژین های مساله را در (15,4) تشکیل می‌دهیم

$$\begin{aligned} L(f, \varsigma, \varsigma^*; \alpha, \alpha^*, r, r^*) = & \frac{1}{2} \Phi[\hat{f}] + C \sum_{i=1}^l (\varsigma_i + \varsigma_i^*) + \sum_{i=1}^l \alpha_i^* (y_i - f(x_i) - \epsilon - \varsigma_i^*) \\ & + \sum_{i=1}^l \alpha_i (\hat{f}(x_i) - y_i - \epsilon - \varsigma_i) - \sum_{i=1}^l (r_i \varsigma_i + r_i^* \varsigma_i^*) \end{aligned} \quad (16.4)$$

بگونه‌ای که α, α^*, r, r^* در شرایط نامنفی صدق کنند. لاگرانژین در (16,4) دارای یک نقطه زینی در جواب بهینه است. بنابراین، بهینه‌سازی با کمینه‌سازی (16,4) نسبت به متغیرهای اصلی $\hat{f}, \varsigma, \varsigma^*$ و بیشینه‌سازی نسبت به متغیرهای دوگان α, α^*, r, r^* سروکار دارد.

توجه داشته باشید که تا این جا، هیچ ساختاری را برای \hat{f} یا همواری تابع در نظر نگرفته ایم. اما همان‌طور که Smale و Cucker می‌گویند "فرآیند یادگیری در خلا بوقوع نمی‌پیوندد."

بمنظور یافتن تابع \hat{f} ، نیازمند تعیین فضای فرضیه هستیم که در جستجویمان آن را در نظر می‌گیریم. در این جا، فرض می‌کنیم که، تابع درون یابی \hat{f} متعلق به یک RKHS مانند H باشد. بنابراین، می‌توان آن را به صورت (12,4) با نرمی به شکل (13,4) نشان داد. با این حال، با دنبال کردن انحراف Vapnik از ϵ -SVR، انحراف صریحی را برای \hat{f} یعنی

$$\hat{f}(x) = \sum_{n=1}^{\infty} c_n \phi_n(x) + b. \quad (17.4)$$

در نظر می‌گیریم. با جایگذاری $(17,4)$ برای \hat{f} و $(13,4)$ برای $\Phi[\hat{f}]$ ، لاگرانژین به شکل

$$\begin{aligned} L(b, c, \varsigma, \varsigma^*; \alpha, \alpha^*, r, r^*) = & \frac{1}{2} \sum_{n=1}^{\infty} \frac{c_n^2}{\lambda_n} + C \sum_{i=1}^l (\varsigma_i + \varsigma_i^*) \\ & + \sum_{i=1}^l \alpha_i^* \left(y_i - \sum_{n=1}^{\infty} c_n \phi_n(x_i) - b - \epsilon - \varsigma_i^* \right) \\ & + \sum_{i=1}^l \alpha_i \left(\sum_{n=1}^{\infty} c_n \phi_n(x_i) + b - y_i - \epsilon - \varsigma_i \right) \\ & - \sum_{i=1}^l (r_i \varsigma_i + r_i^* \varsigma_i^*). \end{aligned} \quad (18.4)$$

در می‌آید. برای رسیدن به تابع هدف دوگان، نیازمند کمینه‌سازی لاگرانژین نسبت به متغیرهای اولیه و حذف آنها با جایگذاری هستیم. ما نیازمند این هستیم که، مشتقات جزئی L نسبت به متغیرهای اولیه $\hat{f}(b, c_n), \varsigma_i^*, \varsigma_i$ صفر باشد

$$\frac{\partial L}{\partial c_n} = 0 \Rightarrow c_n = \lambda_n \sum_{i=1}^l (\alpha_i^* - \alpha_i) \phi_n(x_i) \quad (19.4)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^l (\alpha_i^* - \alpha_i) = 0 \quad (20.4)$$

$$\frac{\partial L}{\partial \varsigma_i} = 0 \Rightarrow r_i = C - \alpha_i \quad (21.4)$$

$$\frac{\partial L}{\partial \varsigma_i^*} = 0 \Rightarrow r_i^* = C - \alpha_i^* \quad (22.4)$$

با جایگذاری $(19,4)$ در مدل ما برای درون یابی تابع داریم

$$\begin{aligned}
 \hat{f}(x) &= \sum_{n=1}^{\infty} c_n \phi_n(x) + b \\
 &= \sum_{n=1}^{\infty} \left(\lambda_n \sum_{i=1}^l (\alpha_i^* - \alpha_i) \phi_n(x_i) \right) \phi_n(x) + b \\
 &= \sum_{i=1}^l (\alpha_i^* - \alpha_i) \sum_{n=1}^{\infty} \lambda_n \phi_n(x_i) \phi_n(x) + b \\
 &= \sum_{i=1}^l (\alpha_i^* - \alpha_i) K(x, x_i) + b
 \end{aligned} \tag{۲۳.۴}$$

که در آن، آخرین تساوی از (۱۱،۴) می‌آید. به صورت مشابه، ما (۱۹،۴) را در معادله یمان برای نرم جایگذاری می‌کنیم

$$\begin{aligned}
 \|\hat{f}\|^2 &= \sum_{n=1}^{\infty} \frac{c_n^2}{\lambda_n} \\
 &= \sum_{n=1}^{\infty} \lambda_n \left(\sum_{i=1}^l (\alpha_i^* - \alpha_i) \phi_n(x_i) \right) \left(\sum_{j=1}^l (\alpha_j^* - \alpha_j) \phi_n(x_j) \right) \\
 &= \sum_{n=1}^{\infty} \lambda_n \sum_{i,j=1}^l (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) \phi_n(x_i) \phi_n(x_j) \\
 &= \sum_{i,j=1}^l (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) \sum_{n=1}^{\infty} \lambda_n \phi_n(x_i) \phi_n(x_j) \\
 &= \sum_{i,j=1}^l (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) K(x_i, x_j).
 \end{aligned} \tag{۲۴.۴}$$

با جایگذاری (۱۲،۴)، (۲۲،۴)، (۲۳،۴) و (۲۴،۴) در لاگرانژین، به

$$\begin{aligned}
 L(\alpha, \alpha^*) &= \frac{1}{2} \sum_{i,j=1}^l (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) K(x_i, x_j) + C \sum_{i=1}^l (\zeta_i + \zeta_i^*) \\
 &\quad + \sum_{i=1}^l \alpha_i^* \left(y_i - \sum_{j=1}^l (\alpha_j^* - \alpha_j) K(x_i, x_j) - b - \epsilon - \zeta_i^* \right) \\
 &\quad + \sum_{i=1}^l \alpha_i \left(\sum_{j=1}^l (\alpha_j^* - \alpha_j) K(x_i, x_j) + b - y_i - \epsilon - \zeta_i \right) \\
 &\quad - \sum_{i=1}^l (C - \alpha_i) \zeta_i - \sum_{i=1}^l (C - \alpha_i^*) \zeta_i^* \\
 &= -\frac{1}{2} \sum_{i,j=1}^l (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) K(x_i, x_j) - \epsilon \sum_{i=1}^l (\alpha_i^* - \alpha_i) + \sum_{i=1}^l y_i (\alpha_i^* - \alpha_i).
 \end{aligned} \tag{۲۵.۴}$$

می‌رسیم. با داشتن (۲۱،۴) و (۲۲،۴) و این واقعیت که، متغیرهای دوگان در محدودیت‌های نامنفی صدق می‌کنند، به محدودیت $0 \leq \alpha_i, \alpha_i^* \leq C, i=1, \dots, l$ می‌رسیم. بنابراین، مساله دوگان به شکل زیر است

$$\begin{aligned}
 &\text{minimize} \quad -L(\alpha, \alpha^*) \\
 &\text{subject to} \quad \begin{cases} 0 \leq \alpha_i, \alpha_i^* \leq C & \text{for } i=1, \dots, l \\ \sum_{i=1}^l (\alpha_i^* - \alpha_i) = 0 \\ \alpha_i \alpha_i^* = 0 & \text{for } i=1, \dots, l \end{cases}
 \end{aligned} \tag{۲۶.۴}$$

که در آن بجای کمینه‌سازی L ، ما $-L$ را بیشینه می‌کنیم. آخرین محدودیت به صورت خودکار برقرارست، در غیر این صورت، وجود کمبودهای ناصفر را در هر دو مسیر، نشان می‌دهد. ما آن را برای مقایسه‌های بعدی، در نظر می‌گیریم.

بطورکلی، در یک مساله کمینه‌سازی محدب، تابع هدف اولیه، همیشه، بزرگ‌تر یا مساوی با تابع هدف دوگان در مقادیر متغیر اولیه و دوگان است. تفاوت بین آنها، شکاف دوگانگی^۱ نامیده می‌شود. برای مساله برنامه‌ریزی مربعی در (۲۶،۴)، نشان داده شده است که این شکاف صفر است. بنابراین، جواب بهینه را با حل مساله دوگان می‌یابیم. این بدان معناست که، برای مساله کمینه‌سازی منظم شده در

^۱ Duality Gap

(۱۵،۴)، کمینه‌سازی بر فضای توابع هیلبرت به کمینه‌سازی بر R^l می‌رسد. این نتیجه با قضیه نمایش Kimeldorf و Wahba سازگار است [۲۹].

۳.۴ ϵ -SVR و تنک بودن

در یک کمینه محلی مساله بهینه‌سازی محدود شده، شرایط Karush-Kuhn-Tucker برقرار است. در میان این‌ها، شرط فقدان تمامیت، بیان می‌کند که در یک نقطه x_i ، حاصل ضرب بین محدودیت‌ها و متغیرهای دوگان، صفر است. برای این مساله در (۲۶،۴)، دو شرط از این دست عبارتند از

$$\begin{aligned}\alpha_i(f(x_i) - y_i - \epsilon - \varsigma_i) &= 0 \\ \alpha_i(y_i - f(x_i) - \epsilon - \varsigma_i) &= 0.\end{aligned}\quad (27.4)$$

براساس (۲۷،۴)، زمانی که $|f(x_i) - y_i| < \epsilon$ باشد، α_i, α_i^* باید صفر باشد. زمانی که $|f(x_i) - y_i| \geq \epsilon$ باشد، α_i یا α_i^* ممکن است ناصفر باشند. نقاطی با α_i یا α_i^* ناصفر، بردارهای پشتیبان نامیده می‌شوند. با افزایش پارامتر آزاد ϵ ، ما تعداد بردارهای پشتیبان را در حالی که تنک بودن جواب در (۲۳،۴) رخ می‌دهد، کاهش می‌دهیم. تاثیر پارامتر آزاد C روی تنک بودن جواب، ممکن است تنها به صورت تجربی نشان داده شود. زمانی که، مانند تنظیمات مساله فعلی ما، هیچ نویزی در داده‌ها موجود نباشد، مقدار بهینه C بی نهایت است. بنابراین، ϵ تنها پارامتر آزاد این فرمول‌بندی است (بدون احتساب پارامترهای هسته).

۴.۴ ارتباط با تقریب تنک

در تقریب تنک، هدف ما، تقریب یک تابع نامعلوم f با یک ترکیب خطی از مجموعه ثابت Φ از توابع

$$f(x; a) = \sum_{i=1}^n \alpha_i \phi_i(x) \quad (28.4)$$

است، که در آن، $\Phi = \{\Phi_i(x) : i = 1, \dots, n\}$ فرهنگ لغت نامیده می‌شود. فرهنگ لغت معمولاً بیش از اندازه کامل است، که بر این دلالت دارد که، a در $(28, 4)$ یکتا نیست، زیرا بعضی عناصر Φ معمولاً، ترکیباتی خطی از دیگر عناصر هستند. ما مساله را با نیاز به این که، a باید جوابی پراکنده باشد، محدود می‌سازیم؛ جوابی با حداقل تعداد عناصر ناصفر. تابع هزینه زیر، یک فرمول بندی را از این مساله نشان می‌دهد

$$E(a) = \frac{1}{2} \left\| f(x) - \sum_{i=1}^n \alpha_i \phi_i(x) \right\|_{L_2}^2 + \lambda \|a\|_{L_0}^p \quad (29.4)$$

که در آن، $\|\cdot\|_{L_0}$ تعداد عناصر ناصفر یک بردار را می‌شمارد، و $\|\cdot\|_{L_2}$ همان L_2 -نرم است. با این حال، بدلیل L_0 -نرم، کمینه‌سازی تابع هزینه در $(29, 4)$ NP-سخت است. برای رسیدگی کردن به مهارنشده بودن $(29, 4)$ ، [10] کمینه‌سازی تابع هزینه محدب

$$f(x; a) = \sum_{i=1}^n \alpha_i \phi_i(x) \quad (30.4)$$

را ارائه می‌نماید و آن را نویززدایی تعاقبی پایه ای می‌نامد. مولفان L_1 -نرم را به عنوان تقریبی بر L_0 -نرم بکار می‌برند. ما توابع فرهنگ لغت را به کمک هسته بازسازی K از یک RKHS، H تعیین می‌کنیم

$$\phi_i(x) = K(x, x_i), \quad i = 1, \dots, l \quad (31.4)$$

که در آن، $\{(x_i, y_i) : i = 1, \dots, l\}$ مجموعه داده حاصل از نمونه‌گیری f در غیاب نویز است. علاوه بر این، ما معیار L_2 را در $(30, 4)$ با H نرم جایگزین کرده و ϵ را به عنوان پارامتر منظم‌سازی تعریف می‌کنیم. به تابع هزینه زیر می‌رسیم

$$E(a) = \frac{1}{2} \|f(x) - \hat{f}(x, a)\|_H^2 + \epsilon \|a\|_{L_1} \quad (32.4)$$

که در آن، تابع تقریب عبارت است از

$$\hat{f}(x, a) = \sum_{i=1}^l \alpha_i K(x, x_i). \quad (33.4)$$

با فرض این که تابع هزینه f دارای میانگین صفر در H است، یعنی تصویرش روی تابع ثابت، صفر است (فرض بر این نیست که تابعی ثابت در H است)

$$\langle f, 1 \rangle_H = 0 \quad (34.4)$$

نیازمند این هستیم که تابع تقریب نیز دارای میانگین صفر در H باشد

$$\langle \hat{f}, 1 \rangle_H = 0. \quad (35.4)$$

به این منظور، ما K را بگونه‌ای نرمال می‌کنیم که

$$\langle 1, K(x, y) \rangle = 1. \quad (36.4)$$

با جایگذاری (33,4) در (35,4) و با استفاده از (36,4)، به محدودیت زیر می‌رسیم

$$\begin{aligned} \langle \hat{f}, 1 \rangle &= \left\langle \sum_{i=1}^l \alpha_i K(x, x_i), 1 \right\rangle \\ &= \sum_{i=1}^l \alpha_i \langle K(x, x_i), 1 \rangle \\ &= \sum_{i=1}^l \alpha_i = 0. \end{aligned} \quad (37.4)$$

نرم H را در معادله (32,4) از تابع هزینه گسترش می‌دهیم، یعنی

$$\begin{aligned} E(a) &= \frac{1}{2} \|f\|_H^2 - \sum_{i=1}^l a_i \langle f(x), K(x, x_i) \rangle_H \\ &\quad + \frac{1}{2} \sum_{i,j=1}^l a_i a_j \langle K(x, x_i), K(x, x_j) \rangle_H + \epsilon \|a\|_{L_1}. \end{aligned} \quad (38.4)$$

دو ویژگی زیر از یک هسته بازیابی را به یاد بیاورید

$$\langle K(x, x_i), K(x, x_j) \rangle_H = K(x_i, x_j) \quad (39.4)$$

$$\langle f(x), K(x, x_j) \rangle_H = f(x_j) \quad (40.4)$$

جایی که در آخرین معادله، $f(x_i) = y_i$ است، زیرا داده‌ها بدون نویز هستند. تابع هزینه به شکل

$$E(a) = \frac{1}{2} \|f\|_H^2 - \sum_{i=1}^l a_i y_i + \frac{1}{2} \sum_{i,j=1}^l a_i a_j K(x_i, x_j) + \epsilon \|a\|_{L_1}. \quad (41.4)$$

در می‌آید. بردار a می‌تواند به بخش‌های مثبت و منفی آنالیز شود. این منجر به این می‌شود که L_1 -نرم به صورت زیر نوشته شود

$$\|a\|_{L_1} = \sum_{i=1}^l |a_i| = \sum_{i=1}^l (a_i^+ + a_i^-) \quad (42.4)$$

به گونه‌ای که برای هر $i = 1, \dots, l$ داشته باشیم $a_i^+, a_i^- = 0$ و $a_i^+, a_i^- \geq 0$. با استفاده از معادله (42,4) و این واقعیت که $\|f\|_H^2$ نسبت به a_i^+, a_i^- ثابت است، به مساله برنامه‌ریزی مربعی

$$\begin{aligned} & \underset{a_i^+, a_i^-}{\text{minimize}} \quad - \sum_{i=1}^l (a_i^+ - a_i^-) y_i + \frac{1}{2} \sum_{i=1}^l (a_i^+ - a_i^-) (a_j^+ - a_j^-) K(x_i, x_j) + \epsilon \sum_{i=1}^l (a_i^+ + a_i^-) \\ & \text{subject to} \quad \begin{cases} a_i^+, a_i^- \geq 0 & \text{for } i = 1, \dots, l \\ \sum_{i=1}^l (a_i^+ - a_i^-) = 0 \\ a_i^+ a_i^- = 0 & \text{for } i = 1, \dots, l \end{cases} \end{aligned} \quad (43.4)$$

می‌رسیم؛ که در آن، محدودیت دوم با آنچه که در معادله (37,4) است، یکسان می‌باشد. با تغییر نام ضریب a_i^+ به α_i^+ و a_i^- به α_i^- ، واضح بنظر می‌رسد که، معادله (43,4) مساله برنامه‌ریزی مربعی مشابه با معادله (26,4) را مشخص کند. نتیجه می‌گیریم که اگر

- مجموعه داده‌ها بدون نویز باشد، یعنی $y_i = f(x_i)$ ،
- L_2 -نرم با H نرم در جمله برازش داده‌های BPDN جایگزین شود،
- تابع \hat{f} دارای میانگین صفر در RKHS، H باشد،
- اتم‌های فرهنگ لغت بکار رفته در BPDN به صورت (31,4) تعریف شوند،
- و پارامتر منظم‌سازی در SVR مایل به صفر باشد، زمانی که $C \rightarrow \infty$.

سپس، BPDN و ϵ -SVR هم‌ارز هستند، زیرا آنها به مساله برنامه‌ریزی مربعی یکسانی می‌رسند. آخرین شرط، به این واقعیت منتقل می‌شود که، ϵ -SVR منجر به تابع درون یابی خواهد شد که، محدود به بیش برآزش داده‌هاست (زیرا، تاثیر جمله منظم‌سازی تعدیل می‌شود). این مطلب برای آنالیز آنچه که درباره طرح تقریب تنک می‌گوییم، مهم است، این نسخه بروز شده BPDN است.

۵

فصل پنجم

نتیجه گیری و کارهای آینده

در این پایان نامه، نشان داده ایم که، کمینه سازی زیان مربعی ℓ_1 -منظم شده برای دسته بندی، یک موفقیت از هر دو نظر محاسباتی و آماری است. همچنین، نشان دادیم که کمینه سازی زیان مربعی دسته بندی ℓ_1 -منظم شده برای بازسازی، چندان ارزشی ندارد. روش های ساده تر، نظیر دسته بندی kNN و WkNNR حداقل خوب هستند.

چهار حوزه را ارائه نمودیم که، برای پژوهش های آینده نیز مناسب هستند. یکی، حالتی که طراحی با نگاشت های غیرخطی متغیرهای مشاهده شده، پر می شود. چگونه این بر مساله بهینه سازی تاثیر خواهد گذاشت و چگونه بر اجرای دسته بندی تاثیر گذار خواهد بود. [۴۴] این مساله را برای رگرسیون مرزی (یا دسته بندی کمترین مربعات منظم شده) پاسخ داده است. ما همچنین، به جستجوی نزدیک تر در مساله دوگان SVM علاقه مند هستیم. آیا این مساله می تواند بگونه ای موثر با افزودن شرط تنک بودن روی متغیر دوگان α حل شود؟ آیا این منجر به یک دسته بندی کننده سریعتر خواهد شد که، در کارایی با SVM قابل مقایسه باشد؟ یک نقص مهم اصل بهینه سازی لسو این است که، جواب های لسو، پایدار نیستند. در نهایت، ما علاقه مند به دانستن این هستیم که، آیا kNN یا WkNN برای یادگیری ویژگی در مقایسه با روش های تقریب تنک در یادگیری بازنمایی تنک به عنوان مثال در دسته بندی تصاویر، مناسب هستند. در ادامه به هر یک از این ایده ها، می پردازیم.

۱.۵ رگرسیون غیرخطی

مجموعه $\{\phi_i, i=1, \dots, p\}$ از p تابع از پیش تعیین شده را با حداقل یک عضو غیرخطی در نظر بگیرید. می توان این توابع را برای بیان متغیرهای اصلی $\{x_i, i=1, \dots, p\}$ بکار برد و به تابع رگرسیون غیرخطی حاصل برای ترکیب خطی این نگاشت ها، رسید

$$\sum_{i=1}^p a_i \phi_i(x_i) + a_0. \quad (1.5)$$

می توان این مدل را با استفاده از الگوریتم هایی یکسان برای برازش مدل های خطی، برازش نمود، البته با مزیت دستیابی به جوابی که یک ارتباط غیرخطی را بین متغیرهای پاسخ و تفسیر، مدل سازی می کند [۲۷]. الگوریتم هایی یکسان را می توان بکار برد، زیرا این مدل هنوز، در پارامترها (ضرایب رگرسیون) خطی است. این مدل تنها در متغیرهای تفسیری غیرخطی است.

مثالی از کاربرد این روش، تعاقب تطابق هسته است [۵۴]. در این چارچوب، $\phi_i(x) = k(x_i, x)$ است، که در آن، k تابعی دومتغیره است که، لزوماً شرایط مرسر^۱ را ندارد. مثالی دیگر، فرمول بندی لسوی هسته ای شده است [۵۷].

۲.۵ منظم سازی α : متغیر دوگان SVM

بیاد بیاورید که، انگیزه ما برای ℓ_1 -منظم سازی، رسیدن به یک دسته بندی کننده است که، سریع تر از SVM باشد؛ البته در این معنا که، ارزیابی های هسته کمتری در زمان آزمون وجود دارد

$$f(x) = \sum a_i k(x_i, x). \quad (۲.۵)$$

همچنین، بیاد بیاورید که، می گوییم، کنترل مستقیمی بر تنک بودن α در مساله بهینه سازی SVM نداریم. هدف ما از مستقیماً پراکنده ساختن چیست؟ بیایید نگاهی دیگر به مساله بهینه سازی دوگان SVM بیندازیم.

۳.۵ ناپایداری و غیریکتا بودن جواب های لسو

جواب های لسو زمانی که $rank(X) \neq p$ باشد، یکتا نیستند [۵۰]. این حالت بویژه زمانی اتفاق می افتد که $p > n$ باشد. در این حالت، جواب های چندگانه بر مساله بهینه سازی لسو وجود دارد. جواب مساله بهینه سازی لسو باید به صورت

$$a \in \arg \min_a \frac{1}{2} \|y - Xa\|^2 + \lambda \|a\|_1. \quad (۳.۵)$$

باشد، این امر منجر به دو مساله بزرگ می شود؛ نخست این که، یک جواب می تواند یک ضریب α مثبت باشد، در حالی که، دیگری دارای یک ضریب α منفی است. علاوه بر این، دو جواب متمایز می تواند دارای پشتیبانهای متمایزی باشند.

^۱ Mercer 's conditions

این مطلب، زمانی که کار انتخاب متغیر پیش‌بینی نیست، مساله ساز است. با این حال، ما علاقه‌مند به مطالعه تاثیرات ناسازگاری جواب‌های لسو روی کارایی پیش‌بینی و پایداریش هستیم.

۴.۵ kNN یا تعاقب تطابقی برای یادگیری ویژگی

اگر تقریب تنک برای کدگذاری تنک خوب باشد، ترجیح می‌دهیم که بدانیم، آیا تعاقب تطابقی یا kNN می‌تواند موثر باشد. مزیت تعاقب تطابقی و kNN این است که، آنها از نظر محاسباتی کمتر از تقریب تنک با روش‌های آزادسازی محدب نظیر SPGL1، هزینه‌بر هستند.

واژه‌نامه

الف

soft thresholding	آستانه‌سازی نرم
posterior probability	احتمال پسین
cross-validation	اعتبارسنجی متقابل
shrinkage	انقباض

ب

sparse representation	بازنمایی تنک
coefficient vector	بردار ضریب
weight vector	بردار وزن
feature vector	بردار ویژگی
quadratic programming	برنامه‌ریزی مربعی

پ

regressor	پیشگو
-----------	-------

ت

basis pursuit	تعاقب پایه
---------------	------------

matching pursuit	تعاقب تطابقی
sparse approximation	تقریب تنک
sparse	تنک

د

sparse representation classification	دسته‌بندی بازنمایی تنک
support vector classification	دسته‌بندی بردار پشتیبان
plug-in classifier	دسته‌بندی‌کننده‌ی جانشین
accuracy	دقت (دسته‌بندی‌کننده)

ر

ridge regression	منظم l_2 رگرسیون
sparse regression	رگرسیون تنک
sparse approximation weighted regression (SPARROW)	رگرسیون وزن‌دار مبتنی بر تقریب تنک

ز

hinge loss	زیان محوری
square loss	زیان مربعی
Down sampling	زیرنمونه‌گیری

ش

sparsity constraint

شرط تنک بودن

ف

reproducing kernel Hilbert space

فضای هیلبرت هسته تکثیری

Fisher consistent

فیشر سازگار

ق

generalization performance

قابلیت تعمیم (دسته‌بندی کننده)

ک

efficient

کارآمد

coordinate descent

کاهش مختصاتی

auto-encoder

کدگذار خودکار

sparse coding

کدگذاری تنک

kernel

کرنل

م

restricted Boltzmann machine (RBM)	ماشین بولتزمن تحدیدی
additive models	مدل‌های جمعی
linear inverse problem	مسئله‌ی معکوس خطی
training samples	مشاهدات آموزشی
regularization	منظم‌سازی

ن

Riesz representation theorem	نظریه‌ی بازنمایی ریس
------------------------------	----------------------

و

observational weights	وزن‌های شهودی
feature	ویژگی

ی

feature learning	یادگیری ویژگی
------------------	---------------

منابع و مراجع

1. Robert Andersen. Robust regression for the linear model. In Modern Methods for Robust Regression, pages 47-70. Sage Publications, 2008.
2. S. Arora, L. Babai, J. Stern, and Z. Sweedyk. The hardness of approximate optima in lattices, codes, and systems of linear equations. In Proceedings of the IEEE 34th Annual Foundations of Computer Science, pages 724-733, Washington, DC, USA, 1993. IEEE Computer Society.
3. Peter L. Bartlett, Michael I. Jordan, and Jon D. McAuliffe. Convexity, classification, and risk bounds. Technical Report 638, Department of Statistics, U.C. Berkeley, 2003.
4. Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for optimal margin classifiers. In Proceedings of the Fifth Annual Workshop on Computational Learning Theory (COLT), 1992.
5. Stephane Boucheron, Olivier Bousquet, and Gabor Lugosi. Theory of classification : A survey of some recent advances. ESAIM: Probability and Statistics, 9:323-375, 2005.
6. Alfred M. Bruckstein, David L. Donoho, and Michael Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. SIAM Review, 51(1): 34-81, 2009.
7. Peter Buhlmann and Bin Yu. Boosting with the L2 loss: regression and classification. Journal of the American Statistical Association, 98(462):324-339, 2003.
8. Andreas Buja, Trevor Hastie, and Robert Tibshirani. Linear smoothers and additive models. The Annals of Statistics, 17(2):435-555, 1989.
9. S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. SIAM J. Sci. Comput., 20(1):33-61, Aug. 1998.
10. Scott S. Chen, David L. Donoho, and Michael A. Saunders. Atomic decomposition by basis pursuit. Technical Report 479, Department of Statistics, Stanford University, May 1995.

11. W. S. Cleveland and S. J. Devlin. Locally weighted regression: an approach to regression analysis by local fitting. *Journal of the American Statistical Association*, 83(403):596- 610, 1988.
12. Adam Coates and Andrew Ng. The importance of encoding versus training with sparse coding and vector quantization. In *International Conference on Machine Learning (ICML)*, pages 921-928, 2011.
13. Adam Coates, Andrew Y. Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. *Journal of Machine Learning Research – Proceedings Track*, 15:215-223, 2011.
14. Felipe Cucker and Steve Smale. On the mathematical foundations of learning. *Bulletin of the American Mathematical Society*, 39:1-49, 2002.
15. L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, 1996.
16. M. Elad. *Sparse and redundant representations: From theory to applications in signal and image processing*. Springer, 2010.
17. John Fox. *Robust regression: Appendix to an R and S-PLUS companion to applied regression*, 2002.
18. Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1, part 2):119-139, 1997.
19. Jerome H. Friedman. Fast sparse regression and classification. In Paul Eilers, editor, *Proceedings of the 23rd International Workshop on Statistical Modelling*, pages 27-57. Statistical Modelling Society, 2008.
20. Jerome H. Friedman and Werner Stuetzle. Projection pursuit regression. *Journal of the American Statistical Association*, 76(376):817-823, 1981.
21. Jort F. Gemmeke, Tuomas Virtanen, and Antti Hurmalainen. Exemplar-based sparse representations for noise robust automatic speech recognition. *IEEE Transactions on Audio, Speech and Language Processing*, 19(7):2067-2080, 2011.
22. Frederico Girosi. An equivalence between sparse approximation and support vector machines. *Neural Computation*, 10:1455-1480, August 1998.

23. W. Hardle and O. Linton. Applied nonparametric methods. Technical Report 1069, Yale University, 1994.
24. T. Hastie, R. Tibshirani, and J. Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer-Verlag, 2 edition, 2009.
25. T. J. Hastie and C. Loader. Local regression: Automatic kernel carpentry. *Statistical Science*, 8(2):120-129, 1993.
26. Trevor Hastie and Ji Zhu. Comment. *Statistical Science*, 21:352-357, 2006.
27. Tim C. Hesterberg, Nam H. Choi, Lukas Meier, and Chris Fraley. Least angle and ℓ_1 - penalized regression: A review. *Statistics Surveys*, 2:61-93, 2008.
28. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12:55-67, 1970.
29. George S. Kimeldorf and Grace Wahba. Some results on Tchebycheffian spline functions. *Journal of Mathematical Analysis and Applications*, 33:82-95, 1971.
30. Su-In Lee, Honglak Lee, Pieter Abbeel, and Andrew Y. Ng. Efficient L_1 regularized logistic regression. In *Proceedings of the 21st National Conference on Artificial Intelligence (AAAI)*, 2006.
31. M. S. Lewicki. Efficient coding of natural sounds. *Nature Neuroscience*, 5(4):356-363, Mar. 2002.
32. Yi Lin. A note on margin-based loss functions in classification. Technical report, Department of Statistics, University of Wisconsin, Madison, 2002.
33. Julien Mairal. Sparse Coding for Machine Learning, Image Processing and Computer Vision. PhD thesis, Ecole Normale Supérieure de Cachan, 2010.
34. S. Mallat. *A Wavelet Tour of Signal Processing: The Sparse Way*. Academic Press, Elsevier, Amsterdam, 3rd edition, 2009.
35. J. Marron and M. Todd. Distance weighted discrimination. Technical report, School of Operations Research and Industrial Engineering, Cornell University, 2002.

36. Hosein Mohimani, Massoud Babaie-Zadeh, and Christian Jutten. A fast approach for overcomplete sparse decomposition based on smoothed ℓ_0 norm. *Transactions on Signal Processing*, 57:289-301, January 2009.
37. E. Nadaraya. On estimating regression. *Theory of Probability and its Applications*, 9 (1):141-142, 1964.
38. B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607-609, 1996.
39. Mark D. Plumbley, Thomas Blumensath, Laurent Daudet, Remi Gribonval, and Mike E. Davies. Sparse Representations in Audio and Music: from Coding to Source Separation. *Proceedings of the IEEE*, 98(6):995-1005, 2010.
40. Tomaso Poggio, Lorenzo Rosasco, and Andre Wibisono. Sufficient conditions for uniform stability of regularization algorithms. Technical Report CBCL-284, Center for Biological and Computational Learning, MIT, December 2009.
41. Pradeep Ravikumar, John Lafferty, Han Liu, and Larry Wasserman. Sparse additive models. *Journal of The Royal Statistical Society (Series B)*, 71(5):1009-1030, 2009.
42. Ryan Rifkin. Everything old is new again: a fresh look at historical approaches in machine learning. PhD thesis, Sloan School of Management, Massachusetts Institute of Technology, 2002.
43. Ryan Rifkin and Aldebaro Klautau. In defense of one-vs-all classification. *Journal of Machine Learning Research*, 5:101-141, December 2004.
44. Ryan Rifkin, Gene Yeo, and Tomaso Poggio. Regularized least-squares classification. In *Advances in Learning Theory: Methods, Model and Applications*, volume 190, pages 131-153. IOS Press, 2003.
45. Lorenzo Rosasco, Ernesto De Vito, Andrea Caponnetto, Michele Piana, and Alessandro Verri. Are loss functions all the same? *Neural Computation*, 16(5):1063-107, 2004.

46. D. Ruppert. Local polynomial regression and its applications in environmental statistics. Technical report, Cornell University, 1996.
47. D. Ruppert and M. P. Wand. Multivariate locally weighted least squares regression. *The Annals of Statistics*, 22:1346-1370, 1994.
48. Alex J. Smola and Bernhard Scholkopf. A tutorial on support vector regression. *Statistics and Computing*, 14:199-222, August 2004.
49. Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society (Series B)*, 58:267-288, 1996.
50. Robert Tibshirani. Regression shrinkage and selection via the lasso: a retrospective. *Journal of the Royal Statistical Society (Series B)*, 73(3):273-282, 2011.
51. Ryan J. Tibshirani. The lasso problem and uniqueness. arXiv:1206.0313v1, 2012. J. A. Tropp and S. J. Wright. Computational methods for sparse solution of linear inverse problems. *Proceedings of the IEEE*, 98(6):948-958, June 2010.
52. E. van den Berg and M. P. Friedlander. Probing the Pareto frontier for basis pursuit solutions. *SIAM Journal on Scientific Computing*, 31(2):890-912, 2008.
53. Vladimir N. Vapnik. The nature of statistical learning theory. *Statistics for Engineering and Information Science*. Springer-Verlag New York, Inc., New York, NY, USA, 1995.
54. Pascal Vincent and Yoshua Bengio. Kernel matching pursuit. *Machine Learning*, 48: 165-187, September 2002.
55. Wagner, J. Wright, A. Ganesh, Z. Zhou, H. Mobahi, and Y. Ma. Towards a practical face recognition system: Robust alignment and illumination via sparse representation. To appear in *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2011.
56. Grace Wahba. Spline models for observational data, volume 59 of CBMS-NSF Regional Conference Series in Applied Mathematics. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1990.

57. Gang Wang, Dit-Yan Yeung, and Frederick H. Lochovsky. The kernel path in kernelized lasso. In International Conference on Artificial Intelligence and Statistics (AISTATS), 2007.
58. Sijian Wang, Bin Nan, Saharon Rosset, and Ji Zhu. Random Lasso. *Annals of Applied Statistics*, 5(1):468-485, 2011.
59. John Wright, Allen Y. Yang, Arvind Ganesh, S. Shankar Sastry, and Yi Ma. Robust face recognition via sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31:210-227, 2009.
60. Huan Xu, Constantine Caramanis, and Shie Mannor. Sparse algorithms are not stable: A no-free-lunch theorem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(1):187-193, 2012.
61. J. Yang, K. Yu, Y. Gong, and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
62. Guo-Xun Yuan, Kai-Wei Chang, Cho-Jui Hsieh, and Chih-Jen Lin. A comparison of optimization methods and software for large-scale ℓ_1 -regularized linear classification. *Journal of Machine Learning Research*, 11:3183-3234, 2010.
63. Guo-Xun Yuan, Chia-Hua Ho, and Chih-Jen Lin. Recent advances of large-scale linear classification. Submitted, 2011.
64. Tong Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32:56-134, March 2004.
65. Tong Zhang and Frank J. Oles. Text categorization based on regularized linear classification methods. *Information Retrieval*, 4:5-31, 2001.

Abstract*

The ℓ_1 -regularized square loss minimization problem has recently gained much attention. This optimization principle has two main applications in the machine learning literature. Specically, the lasso or basis pursuit de-noising (although the two are not entirely equivalent in practice) is the optimization principle used for solving the linear inverse problem: $y = Xa$, under convex sparsity constraints. When the lasso is used for regression and classification, y is a vector of outputs. When it is used for sparse coding and feature learning, or in the context of sparse representation classification, y is the feature vector or signal itself.

The use of lasso for regression is already well-established. In this thesis, we argue that the use of lasso for classification also has its advantages. One might think that the square loss is not appropriate for the classification task, however, theoretical results show that all convex loss functions are Fisher consistent. Additionally, square loss minimization, like logistic loss minimization, and unlike hinge loss minimization, gives estimates of the posterior probability. The value of the posterior probability at a point tells us about the confidence of the classifier in its prediction. Another benefit of the lasso for classification is that ℓ_1 -regularization leads to a sparse classifier, that once trained, can be evaluated quickly. Additionally, one has direct control over the sparsity of the solution through the regularization parameter. The only problem with the lasso is the stability of its solutions (Wang et al., 2011).

The second part of the thesis, is on the use of the lasso for signal and feature representation. The lasso or basis pursuit de-noising is also an integral part of the sparse representation classification method. We extend this method to the regression setting. Through experimental results we argue that one can easily achieve the same or even better results using simpler methods like k-nearest neighbor classification which is also better motivated theoretically. We conclude that ℓ_1 -regularized square loss minimization is not worth it.

Key Words: square loss minimization, ℓ_1 -regularization, binary classification, sparse representation.



**Amirkabir University of Technology
(Tehran Polytechnic)**

Department of Computer Engineering and Information Technology

MSc Thesis

**Title
Efficient Classification Based on
Sparse Regression**

**By
Pardis Noorzad**

**Supervisor
Dr. Mohammad Rahmati**

**Advisor
Dr. Nasrollah Moghaddam Charkari
Dr. Mohammad Mehdi Ebadzadeh**

July 2012