



نکات قابل توجه

- پاسخ‌های خود را به آدرس ml941.sharif@gmail.com ارسال کنید.
- نام فایل ارسالی تان به صورت STDNUM-P2.zip باشد.
- در انتخاب محیط و کتابخانه‌ها برای پیاده‌سازی محدودیتی وجود ندارد.
- مشاهده تقلب در پروژه‌ی ارسالی پیامدهای جالبی نخواهد داشت.

۱ مقدمه

در فاز اول پروژه تعدادی از روش‌های یادگیری نظارتی را دیدیم. در این فاز قرار است تعدادی از روش‌های غیرنظارتی را ببینید و مقایسه‌ای روی آن‌ها داشته باشید. بدین منظور دو مساله‌ی کاهش بعد و خوشه‌بندی در نظر گرفته شده است. همان‌طور که می‌دانید در مساله‌ی خوشه‌بندی داده‌ها برچسب نخورده‌اند و هدف ما این است که این برچسب‌ها را یاد بگیریم. تاکنون با تعدادی روش احتمالی و غیر احتمالی آشنا شده‌اید که باید از آن‌ها در این بخش استفاده کنید. همچنین در بسیاری موارد خوب است که با کاهش بعد سعی کنیم نویز و همبستگی بین ابعاد را پایین بیاوریم تا به نتیجه‌ی بهتری برسیم. در قسمت اول مساله‌ی خوشه‌بندی مطرح می‌شود و روش‌های kmeans و مدل مخلوط با هم مقایسه می‌شوند. در قسمت دوم مساله‌ی کاهش بعد در دسته‌بندی چند کلاس مطرح است. چالشی که در این جا وجود دارد این است که آیا روی داده‌های هر کلاس به طور مستقل کاهش بعد را داشته باشیم یا اینکه از راه دیگری استفاده کنیم. بدین منظور سه روش کاهش بعد روی داده‌ها مطرح شده است و می‌خواهیم مقایسه‌ای روی آن‌ها داشته باشیم.

۲ خوشه‌بندی

در این بخش زیرمجموعه‌ای از مجموعه داده‌گان MNIST بدون برچسب را در اختیار دارید و قرار است روی تصاویر آن خوشه‌بندی انجام دهید. این زیرمجموعه شامل ارقام صفر تا پنج است یعنی از ارقامی که شباهت زیادی با هم دارند تنها یکی‌شان در این مجموعه داده وجود دارد. در قسمت اول ابتدا با گذاشتن آستانه روی شدت روشنایی هر نقطه، تصاویر را دودویی کنید. ابتدا از روش kmeans برای خوشه‌بندی استفاده کنید. سپس فرض کنید نقاط تصاویر از یک مدل مخلوط

برنولی^۱ تولید شده اند.^۲ روش EM را خودتان پیاده سازی کنید و خوشه بندی را انجام دهید. حال آستانه را بردارید و تصاویر اولیه را در نظر بگیرید. این بار فرض کنید نقاط تصاویر از یک مدل مخلوط چندجمله^۳ تولید شده اند.^۴ روابط EM را بنویسید و خوشه بندی را انجام دهید. پس از اتمام کار، بعد داده ها را با استفاده از PCA کاهش دهید و دوباره سه روش یاد شده را تست کنید و اثر کاهش بعد را گزارش کنید. کلیه روش های پیاده شده را با هم مقایسه کنید.

روش ارزیابی

برای مقایسه روش های که پیاده کرده اید برچسب داده های آموزش به شما داده می شود. دقت کنید که شما این برچسب ها را نمی توانید در هنگام خوشه بندی استفاده کنید و باید فرض کنید آن ها را در اختیار ندارید. از معیارهای RandIndex، Micro-F₁ و Macro-F₁ برای ارزیابی استفاده کنید.

۳ کاهش ابعاد

شرح مسئله

یک مسئله دسته بندی با چند کلاس را در نظر بگیرید که بعد داده های آن بسیار زیاد است و قبل از اینکه یک الگوریتم دسته بندی را بر آن ها اعمال کنیم، باید به طریقی بعد آن ها را کم کنیم. برای مثال فرض کنید ۱۰ کلاس وجود دارد و می خواهیم داده ها را به ۱۰۰ بعد کاهش دهیم. شاید اولین گزینه برای این کار PCA باشد. پس راه اول این است که روی کل داده ها PCA انجام دهیم. ولی در این حالت از برچسب داده ها استفاده ای نشده است. یک پیشنهاد دیگر این است که این ۱۰۰ بعد این گونه انتخاب شوند که PCA با $k = 10$ روی داده ها هر کدام از کلاس ها اعمال شوند تا ویژگی هایی که واریانس داده های هر کلاس در آن راستا بیشینه است، استخراج شوند. در نهایت این پایه های ۱۰ تایی برای هر کلاس را جمع می کنیم و یک داده ی تست جدید را روی فضایی ۱۰۰ بعدی حاصل تصویر می کنیم. ایرادی که ممکن است به راه حل بالا گرفته شود این است مسئله برای هر کلاس بدون توجه به پایه هایی که برای کلاس های دیگر پیدا شده حل می شود. چه بسا بخاطر شباهت ذاتی داده ها بعضی از پایه هایی که واریانس را بیشینه می کند بین کلاس ها مشابه یا حتی مشترک باشد^۵، در این حالت تعدادی از بعد هایی که استفاده می کنیم در حقیقت «به هدر می روند» و مهمتر از آن، ویژگی های پیدا شده اصلاً جدا کننده^۶ نیستند.

تحت تاثیر این انتقاد روش سومی پیشنهاد می شود. برای شرح این روش فرض کنید که داده ها با x و برچسب ها با y نشان داده می شوند. هم چنین $C_i = \mathbb{E}[xx^T | y = i]$ ماتریس گشتاور دوم داده های کلاس i م را نشان می دهد. در این

^۱ Bernouli Mixture Model

^۲ 9.3.3 from Bishop

^۳ multinomial Mixture Model

^۴ Example 9.19 from Bishop

^۵ این مسئله می تواند مانع استقلال خطی ۱۰۰ بردار پیدا شده بعد از جمع مولفه های اصلی هر کلاس شود. شما هنگام پیاده سازی باید تدبیری برای این حالت اندیشیده باشید.

^۶ discriminative

روش بجای روش دوم که به دنبال یافتن پایه‌هایی مثل v بودیم که $v^T C_i v$ را بیشینه کند، دنبال بیشینه کردن نسبت این کمیت برای دو کلاس هستیم یعنی

$$R_{ij}(v) = \frac{\mathbb{E}[(v^T x)^2 | y = i]}{\mathbb{E}[(v^T x)^2 | y = j]} = \frac{v^T C_i v}{v^T C_j v}$$

نشان دهید v ‌هایی که بیشینه محلی این تابع هستند، بردار ویژه‌های تعمیم یافته خواهند بود. یعنی v ‌هایی که $C_i v = \lambda C_j v$. دقت کنید که در روش دوم اگر ۱۰ کلاس داشته باشیم، ۱۰ سری مولفه اصلی^۷ پیدا می‌شود که از هر سری همه یا تعدادی که اصلی‌تر هستند انتخاب می‌شوند. اما در روش سوم به ازای هر i و j یک سری جواب خواهیم داشت یعنی $(\cdot)^j$ مجموعه از v که از هر کدام تعدادی را انتخاب می‌کنیم تا فضای کاهش بعد را تشکیل دهند.

در ادامه فرض کنید بردار ویژه‌ها طوری نرمال شده اند که $v^T C_j v = 1$ و $v^T C_i v = \lambda$. همچنین یک جمله منظم‌سازی به C_j اضافه می‌کنیم تا ماتریس در مخرج همیشه از مرتبه کامل باشد؛ به این ترتیب که

$$R_{ij}^\beta = \frac{v^T C_i v}{v^T (C_j + \beta I) v}$$

پیاده‌سازی

شما باید سه روش بالا را پیاده‌سازی کنید تا روشی که نتیجه‌ی بهتری می‌دهد پیدا شود. نتیجه بهتر یعنی خطای کمتر در دسته‌بندی چند کلاسه. برای مقایسه، روی خروجی هر روش یک دسته‌بند svm با هسته گاوسی واریانس ۱۰ و هزینه ۱ یاد بگیرید. (اگر از libsvm استفاده می‌کنید پارامترهای بیان شده معادل -g 0.1 -t 2 -c 1 svm-train خواهد بود)

تعیین تعداد ابعادی که در کاهش بعد استفاده می‌کنید به عهده‌ی شماست. همچنین در پیاده‌سازی روش سوم یک آستانه α در نظر بگیرید و بردار ویژه‌هایی که مقدار ویژه متناظر آنها کوچکتر از آستانه است یعنی $\alpha < \lambda$ را دور بریزید. مقدار این آستانه و مقدار ضریب منظم‌سازی یعنی β باید با cross validation تعیین شوند. فرض کنید β ضریبی از میانگین مقدار ویژه‌های C_j است و در cross validation مقدار آن ضریب را تعیین کنید.

علاوه بر این، روش سوم را روی داده‌های MNIST اجرا کرده و چهار بردار اول برای بهینه‌سازی‌های $R_{۱۷}$ ، $R_{۶۸}$ ، $R_{۰.۸}$ و $R_{۴۹}$ را به صورت تصویر نشان دهید. در مورد کیفیت و توانایی جداسازی این بردارها برای کلاس‌های مربوطه‌شان بحث کنید.

مجموعه داده

داده‌ای که در این قسمت استفاده می‌کنید، زیر مجموعه‌ای از 20 news groups است. این مجموعه داده متشکل از متون ۲۰۰۰۰ خبر است که تقریباً به طور مساوی زیر عنوان یکی از ۲۰ موضوع خبری است. در این فاز شما با یک زیرمجموعه از این داده‌ها کار خواهید که داده‌های مربوط به ۵ موضوع از ۲۰ موضوع کل داده‌هاست. همچنین متن‌ها پیش‌پردازش شده‌اند و ویژگی tf-idf از آن‌ها استخراج شده اند. در نهایت ۲۹۱۹ بردار ویژگی ۴۶۳۹۷ بعدی به همراه برچسب هر کدام تحویل شما داده شده است.

^۷Principal Component

۴ مستندات

مستند ارسالی حداقل باید شامل موارد زیر باشد:

- توضیحاتی در مورد نحوه‌ی پیاده‌سازی هر قسمت و ابزارهای مورد استفاده.
- جداول نتایج به دست آمده در هر قسمت. به همراه مقادیر استفاده شده در cross validation و جواب گرفته شده برای هر کدام.
- محاسبات یا اثبات‌های لازم برای بدست آوردن روابط مورد استفاده.

موفق باشید.