

به نام خدا

نام درس: هوش مصنوعی و سیستم های خبره

پروژه SVM (ماشین های بردار پشتیبان)

کد دانشجویی : ۴۰۲۳۰۰۰۳۲

نام و نام خانوادگی: سروش میرشکاری

پاسخهای پارت اول: آشنایی با SVM خطی

سوال ۱: چرا الان جداساز خطی جواب نمیدهد؟

جواب ساده‌اش اینه که جداساز خطی، همونطور که از اسمش پیداست، فقط بلده یک خط راست برای جدا کردن کلاس‌ها بکشه. وقتی ما از make_blobs استفاده می‌کردیم، داده‌ها مثل دو تا خوشی یا تووه بودن که به راحتی می‌شد یک خط صاف بینشون کشید و از هم جداشون کرد. برای همین دقت مدل نزدیک به ۱۰۰٪ بود.

اما وقتی دیتاست رو به make_moons تغییر دادیم، شکل داده‌ها کاملاً عوض شد. الان داده‌ها شبیه دو تا هلال ماه در هم تنیده هستن. شما هر چقدر هم تلاش کنی، نمی‌توانی یک خط راست پیدا کنی که این دو تا هلال رو به درستی از هم جدا کنه. خط راست نهایتاً می‌تونه یک بخش کوچکی از هر هلال رو درست طبقه‌بندی کنه و بخش بزرگی رو اشتباه می‌گیره. به همین دلیله که دقت مدل به شدت افت می‌کنه (مثلاً نزدیک ۵۰٪ یا حتی کمتر)، چون مدل خطی اساساً برای حل چنین مسئله غیرخطی طراحی نشده و ابزار لازم برای کشیدن یک مرز منحنی رو نداره.

سوال ۲: به نظرت چه جداسازی برای این دیتا مناسب‌تره؟ چرا؟

برای داده‌های غیرخطی مثل make_moons، ما به یک جداساز غیرخطی احتیاج داریم؛ یعنی به جای یک خط راست، به یک مرز تصمیم منحنی نیاز داریم که بتوونه خودش رو با شکل پیچیده‌ی داده‌ها وفق بده.

در SVM، این کار رو با استفاده از کرنل‌ها (Kernels) انجام می‌دیم. به جای کرنل linear، باید از کرنل‌های قوی‌تری مثل:

RBF (Radial Basis Function): این معروف‌ترین و پرکاربردترین کرنل برای مسائل غیرخطی هست. کارش اینه که انگار داده‌ها رو به یک بعد بالاتر می‌بره؛ جایی که جداسازی اونها با یک صفحه (که معادل خط در ابعاد پایین‌تره) ممکن می‌شه. نتیجه‌اش در فضای اصلی ما، یک مرز تصمیم نرم و کاملاً انعطاف‌پذیره که می‌تونه هر شکلی، از دایره گرفته تا شکل‌های خیلی پیچیده‌تر، به خودش بگیره. برای داده‌های هلالی ما، این کرنل عالی عمل می‌کنه.

چندجمله‌ای (Polynomial): این کرنل هم می‌تونه مرزهای منحنی ایجاد کنه. انگار که به مدل اجازه می‌دیم از معادلات درجه ۲، ۳ یا بالاتر برای کشیدن مرز استفاده کنه.

پس دلیل اینکه این جداسازها مناسب‌ترن اینه که انعطاف‌پذیری لازم برای یادگیری الگوهای پیچیده و غیرخطی رو دارن، کاری که از دست یک مدل خطی ساده برنمی‌یاد.

پاسخهای پارت دوم: بازی با کرنل‌ها

سوال ۳: هر کدام از این فیچرها چه چیزی رو نشون میدن؟

دو ویژگی‌ای که از دیتاست سرطان سینه انتخاب کردیم این‌ها هستن:

(mean smoothness) میانگین نرمی: این ویژگی میگه که کانتور یا حاشیه سلول های تومور چقدر صاف و یکنواخته. اگه تغییرات شعاع سلول در نقاط مختلف کم باشه، یعنی سلول صافتر و گردد. معمولاً سلول های سرطانی حاشیه های نامنظمتر و ناهموارتری دارن.

(mean compactness) میانگین فشردگی: این یک معیار برای سنجش "گرد بودن" شکل سلوله. هر چی سلول به دایره نزدیکتر باشه، فشردگی اش بیشتره. این ویژگی از طریق فرمولی شامل محیط و مساحت سلول محاسبه میشه و نشون میده که شکل سلول چقدر پیچیده و نامنظمه.

سوال ۴: کدام کرنل بهترین عملکرد را داشت؟ چرا؟

روی دیتاست سرطان سینه (با این دو ویژگی)، کرنل RBF به طور واضح بهترین عملکرد رو داشت. دلیلش اینه که وقتی ما این دو ویژگی رو روی نمودار رسم میکنیم، میبینیم که مرز بین تومورهای خوش خیم و بد خیم یک خط راست و ساده نیست. بلکه یک رابطه پیچیده تر و غیرخطی بین نرمی و فشردگی سلول ها برای تشخیص بیماری وجود داره.

کرنل خطی چون فقط میتونه یک خط راست بکشه، نمیتونه این دو گروه رو که کمی در هم تنیده شدن، به خوبی جدا کنه و دقتش پایین تره.

کرنل RBF با انعطاف پذیری بالایی که داره، میتونه یک مرز تصمیم منحنی و نرم ایجاد کنه که به بهترین شکل ممکن دور خوش های هر کلاس میپیچه و اونها رو از هم جدا میکنه. برای همین دقت بالاتری به دست میاره.

سوال ۵: آیا مدل خطی کافی بود؛ اگر نه شکل داده چه چیزی را نشان می دهد؟

قطعاً نه، مدل خطی کافی نبود. دقتش پایین تر شنیده این رو ثابت میکنه. اما مهمتر از اون، شکل پراکندگی داده ها روی نمودار به ما نشون میده که چرا مدل خطی کافی نیست. وقتی به نمودار نگاه میکنیم، میبینیم که نقاط آبی (خوش خیم) و قرمز (بد خیم) طوری پخش شدن که نمیشه با یک خطکش یک خط صاف بینشون کشید که کمترین خط را رو داشته باشه. داده ها یک الگوی غیرخطی دارن. این یعنی برای جداسازی موفق، ما به یک مرز تصمیم نیاز داریم که بتوونه بپیچه و خم بشه، کاری که فقط از عهده کرنل های غیرخطی مثل RBF برミارد.

سوال ۶: اگر تعداد ویژگی ها بیشتر بود چه تغییری در انتخاب کرنل رخ می داد؟

این سوال خیلی مهمیه. اگه به جای ۲ ویژگی، از تمام ۳۰ ویژگی دیتاست سرطان سینه استفاده میکردیم، انتخاب ما ممکن بود تغییر کنه. با افزایش تعداد ابعاد (ویژگی ها)، یک پدیده جالب به نام "Curse of Dimensionality" (Curse of Dimensionality) رخ میده. در فضاهای با ابعاد خیلی بالا، داده ها خیلی پراکنده تر میشن و به طرز عجیبی، احتمال اینکه بشه اونها رو با یک صفحه (یعنی یک مرز خطی) از هم جدا کرد، بیشتر میشه.

بنابراین:

برای تعداد ویژگی های بسیار بالا (مثلآ هزاران ویژگی مثل مسائل تحلیل متن): اغلب کرنل خطی سریع تر و حتی گاهی دقیق تر عمل میکنه. چون هم محاسباتش خیلی سُبک تر و هم در اون ابعاد بالا، معمولاً یک جداساز خطی پیدا میشه.

برای تعداد ویژگی های متوسط (مثل ۳۰ ویژگی دیتاست ما): هنوز هم کرنل RBF یک انتخاب خیلی قوی و مطمئنه، چون روابط بین این ۳۰ ویژگی ممکنه پیچیده و غیرخطی باشه. اما همیشه ارزش داره که کرنل خطی رو هم به عنوان یک گزینه سریع و ساده امتحان کنیم.

پس به طور خلاصه، هر چی تعداد ویژگی‌ها بیشتر بشه، تمایل ما به استفاده از کرنل خطی (به دلیل سرعت و کارایی) بیشتر می‌شه.

پاسخهای پارت سوم: بازی با نویز

سوال ۶ و ۷: دقت مدل را با حالت بدون نویز مقایسه کن و تاثیر افزایش نویز را بررسی کن.

همونطور که انتظار می‌رفت، اضافه کردن نویز به داده‌ها مثل اینه که عینک‌مون رو کثیف کنیم؛ دیدن الگوها سخت‌تر می‌شه.

مقایسه با حالت بدون نویز: وقتی نویز (حتی به مقدار کم مثل $\frac{1}{3}$) اضافه کردیم، دقت تمام مدل‌ها کم شد. دلیلش اینه که نویز بعضی از نقاط رو از جای اصلی‌شون هل می‌ده و ممکنه اون‌ها رو به سمت مرز یا حتی اون طرف مرز کلاس مقابله ببره. این کار مرز بین کلاس‌ها رو مبهم می‌کنه و کار مدل رو برای پیدا کردن یک جداساز تمیز سخت می‌کنه.

تاثیر افزایش نویز: وقتی سطح نویز رو از $\frac{1}{5}$ به $\frac{3}{5}$ افزایش دادیم، دقت مدل‌ها باز هم کمتر شد. نویز بیشتر یعنی هرج‌ومرج و همپوشانی بیشتر بین کلاس‌ها. این کار باعث می‌شه مدل در یادگیری الگوی اصلی ضعیفتر عمل کنه و خطای بیشتری داشته باشه.

سوال ۸: کدام کرنل بهتر با نویز کنار آمد؟ چرا؟

در این آزمایش، معمولاً کرنل RBF بهترین مقاومت رو در برابر نویز نشون می‌ده. دلیلش به ماهیت انعطاف‌پذیر اما هوشمندش برمی‌گرده:

کرنل خطی چون انعطاف نداره، با اضافه شدن نویز و غیرخطی‌تر شدن داده‌ها، به شدت عملکردش افت می‌کنه.

کرنل چندجمله‌ای (Poly) می‌تونه خیلی حساس باشه. اگه درجه‌اش بالا باشه، ممکنه سعی کنه خودش رو به تک‌تک نقاط نویزی هم Fit کنه و یک مرز خیلی عجیب و غریب و پر پیچ و خم بکشه که به این حالت بیش‌برازش (Overfitting) روی نویز می‌گیره و باعث می‌شه روی داده‌های جدید خیلی بد عمل کنه.

کرنل RBF یک تعادل خوب داره. به اندازه کافی انعطاف‌پذیره که بتوونه یک مرز غیرخطی مناسب پیدا کنه، اما اونقدر هم حساس نیست که بخواهد خودش رو به هر نقطه نویزی بچسبونه. انگار که سعی می‌کنه یک نمای کلی و نرم از مرز رو یاد بگیره، که همین باعث می‌شه در حضور نویز، استوارتر و قابل اعتمادتر از بقیه باشه.

پاسخهای پارت چهارم: داده‌های درهم‌تنیده

سوال ۹: کدام کرنل بهتر با داده‌ای در هم کنار آمد؟

وقتی داده‌ها به شدت در هم تنیده و دارای همپوشانی هستن (مثل وقتی که $class_sep=0.8$ بود)، باز هم کرنل RBF بهترین عملکرد رو داره. در چنین شرایطی، پیدا کردن یک مرز عالی که هیچ خطایی نداشته باشه غیرممکنه. هدف اینه که مرزی پیدا کنیم که بهترین توازن رو بین طبقه‌بندی درست داده‌ها برقرار کنه.

کرنل خطی که کاملاً شکست می‌خوره. اما کرنل RBF می‌تونه یک مرز نرم و غیرخطی ایجاد کنه که سعی می‌کنه تا جای ممکن خودش رو از وسط ناحیه همپوشانی عبور بده و بهترین جداسازی ممکن رو انجام بده.

سوال ۱۰: آیا با تغییر $n_informative$ یا $class_sep$ عملکرد مدل تغییر کرد؟

بله، این پارامترها به شدت روی عملکرد مدل تأثیر دارند:

(class_sep) میزان جدایی کلاس‌ها: این پارامتر مستقیماً سختی مسئله را تعیین می‌کنه. هر چی این عدد بزرگ‌تر باشد، کلاس‌ها از هم دورتر و جداسازی‌شون راحت‌تره و در نتیجه دقت مدل بالاتر میره. هر چی کوچک‌تر باشد، همپوشانی بیشتر و کار مدل سخت‌تر می‌شود و دقت پایین می‌دارد.

(n_informative) تعداد ویژگی‌های اطلاعاتی: این پارامتر می‌گه چند تا از ویژگی‌ها واقعاً برای تشخیص کلاس مفید هستند. اگه این عدد رو کم کنیم، یعنی اطلاعات مفید کمتری به مدل می‌دمیم و کارش برای پیدا کردن الگو سخت می‌شود و دقتش کم می‌شود. اگه زیادش کنیم، اطلاعات بیشتری در اختیار مدل قرار می‌گیره و می‌تونه بهتر یاد بگیره (البته تا جایی که تعداد ویژگی‌ها از حد معقولی بیشتر نشوند).

پاسخ‌های پارت پنجم: کنترل پیچیدگی با C

سوال ۱۱: بهترین مقدار C کدام بود؟ چرا؟

پارامتر C در SVM مثل یک اهرم کنترل بین دو هدف متضاد عمل می‌کنه:

داشتمن یک مرز ساده و با حاشیه امنیت بالا: این باعث می‌شود مدل روی داده‌های جدید خوب عمل کند (یعنی تعمیم‌پذیری بالایی داشته باشد).

طبقه‌بندی درست تمام نقاط آموزشی: این باعث می‌شود خطای مدل روی داده‌هایی که دیده، کم باشد. C خیلی کوچک (مثلًا ۰/۰۵): مدل به هدف اول اهمیت بیشتری می‌دهد. یعنی یک مرز خیلی ساده و صاف می‌کشد و حاضره چند تا از نقاط آموزشی را اشتباه طبقه‌بندی کند. به این حالت کم‌برازش (Underfitting) می‌گذرد.

C خیلی بزرگ (مثلًا ۱۰۰۰): مدل به هدف دوم وسوس نشون می‌دهد. سعی می‌کند به هر قیمتی شده تمام نقاط آموزشی رو درست جا بده، حتی اگه مجبور بشو یک مرز خیلی پیچیده و پر پیچ و خم بکشد. این مرز روی داده‌های آموزشی عالی عمل می‌کند، اما روی داده‌های جدید افتضاح خواهد بود. به این حالت بیش‌برازش (Overfitting) می‌گذرد.

بهترین مقدار C، مقداری که بهترین توازن رو بین این دو هدف برقرار کند و در نهایت بالاترین دقت روی داده‌های تست (Test) به ما بده. در این آزمایش، معمولاً یک مقدار میانی مثل $C=10$ یا $C=1$ بهترین عملکرد رو داشت چون نه اونقدر ساده بود که دچار کم‌برازش بشود و نه اونقدر پیچیده بود که دچار بیش‌برازش بشود.

سوال ۱۲: نوع غالب خطای مدل در Confusion Matrix چه بود؟ آیا قابل بهبود است؟

برای جواب دادن به این سوال، باید به ماتریس درهم‌ریختگی (Confusion Matrix) نگاه کنیم. این ماتریس به ما می‌گه مدل چه نوع اشتباهاتی کرده. دو نوع خطای اصلی داریم:

(Mثبت کاذب): نمونه در واقعیت منفی (مثلًا خوش‌خیم) بوده، اما مدل به اشتباه گفته مثبته (بدخیم).

(منفی کاذب): نمونه در واقعیت مثبت (بدخیم) بوده، اما مدل به اشتباه گفته منفیه (خوش‌خیم). این خطأ در مسائل پژوهشی معمولاً خطرناک‌تره.

اینکه کدام خطأ غالب بود، به مقدار C و ماهیت داده‌ها بستگی دارد. اما آیا قابل بهبود است؟ بله، قطعاً راه‌های بهبودش این‌هاست:

تنظیم دقیق‌تر هایپرپارامترها: می‌توانیم با GridSearchCV بهترین مقادیر C و gamma را پیدا کنیم. مهندسی ویژگی: شاید با اضافه کردن ویژگی‌های بهتر یا تغییر ویژگی‌های موجود، بشه به مدل کمک کرد. استفاده از class_weight='balanced': اگه یک نوع خطا (مثلاً منفی کاذب) برآمون مهم‌تره، می‌توانیم با این پارامتر به مدل بگیم که برای اون کلاس اهمیت بیشتری قائل بشه و جریمه اشتباه کردن روی اون کلاس رو سنگین‌تر کنه. این کار باعث می‌شه مدل در تشخیص اون کلاس خاص، دقت بیشتری به خرج بده.

پاسخ‌های پارت ششم: SVM چندکلاسه

سوال ۱۳: مرز بین کدام دو کلاس سخت‌تر بود؟ چرا؟

برای فهمیدن این موضوع، باید به نمودار مرز تصمیم نگاه کنیم. معمولاً مرز بین دو کلاسی که از نظر ویژگی‌ها به هم نزدیک‌تر هستن و نقاطشون در نمودار بیشترین همپوشانی رو با هم دارن، سخت‌ترین مرز برای یادگیریه. در دیتاست Wine، معمولاً کلاس ۱ و ۲ ویژگی‌های شبیه‌تری نسبت به کلاس ۰ دارن. در نتیجه، مدل برای کشیدن یک خط مرزی تمیز و قاطع بین این دو کلاس نزدیک به هم، با چالش بیشتری رویرو می‌شه و احتمالاً ناحیه مرزی بینشون کمی مبهم‌تر و باریک‌تر از مرزهای دیگه است.

سوال ۱۴: آیا مدل توانست سه کلاس را به درستی از هم جدا کند؟ چه کرنلی بهتر جواب داد؟

بله، SVM با استفاده از استراتژی‌های داخلی خودش مثل (OvO) یا One-vs-Rest به خوبی می‌توانه مسائل چندکلاسه رو هم مدیریت کنه. در استراتژی OvR، مدل سه تا طبقه‌بند باینری یاد می‌گیره: یکی برای "کلاس ۰ در مقابل بقیه"، یکی برای "کلاس ۱ در مقابل بقیه" و یکی برای "کلاس ۲ در مقابل بقیه". در نهایت برای یک نقطه جدید، هر سه طبقه‌بند نظر می‌دان و هر کدام که با اطمینان بیشتری گفت این نقطه مال منه، برنده می‌شه.

در این مسئله هم، مثل اکثر مسائل پیچیده، کرنل RBF به دلیل انعطاف‌پذیری بالایی که در ایجاد مرزهای غیرخطی داره، بهترین گزینه بود و توانست با دقت خوبی سه کلاس رو از هم جدا کنه.