

گزارش فاز دوم پروژه ماشین لرنینگ پیدا کردن ستارگان متغیر کهکشان M33

امیررضا بهرامی مقدس

یوسف جمشیدی

سوگل سنجری پور

در این فاز قصد داریم الگوریتم های متفاوت ماشین لرنینگ را روی داده هایمان که پیشتر تهیه کرده ایم به دست آوریم. ابتدا داده هایمان را که در فاز قبلی کاهش داده و داده های پرت را شناسایی و حذف کرده بودیم، مجدد به ترتیب کاهشی ایندکس L (معیاری از متغیر بودن یا نبودن ستاره) مرتب کردیم. در این مرحله دیتا را محدودتر میکنیم. به این ترتیب که هر 812 ستاره متغیر را انتخاب میکنیم. از بین بقیه ستارگان نامتغیر 4000 ستاره انتخاب میکنیم. حالا کل مراحل تست و آموزش را با این دیتا ست تمیزتر شده انجام میدهیم. از کل داده ها 20% برای تست، 20% از باقی مانده برای اعتبارسنجی (validation) و بقیه برای آموزش ماشین استفاده شد. در این روند، از stratify استفاده کردیم تا در تقسیم بندی نسبت بین تعداد ستارگان متغیر و نامتغیر حفظ شود.

الگوریتم هایی که برای این فاز استفاده کردیم به ترتیب الگوریتم KNN, Random Forest و SVM و Naive Bayes.

الگوریتم KNN:

در بازشناخت الگو کی-نزدیکترین همسایه (k-nearest neighbors algorithm): یک متد آمار پارمتری است که برای طبقه بندی آماری و رگرسیون استفاده می شود. در هر دو حالت k شامل نزدیک ترین مثال آموزشی در فضای داده ای می باشد و خروجی آن بسته به نوع مورد استفاده در طبقه بندی و رگرسیون متغیر است. در حالت طبقه بندی با توجه به مقدار مشخص شده برای کی، به محاسبه فاصله نقطه ای که میخواهیم برچسب آن را مشخص کنیم با نزدیک ترین نقاط میپردازد و با توجه به تعداد رای حداکثری این نقاط همسایه، در رابطه با برچسب نقطه مورد نظر تصمیم گیری می کنیم. برای محاسبه این فاصله میتوان از روش های مختلفی استفاده کرد که یکی از مطرح ترین این روش ها، فاصله اقلیدسی است. در حالت رگرسیون نیز میانگین مقادیر بدست آمده از کی خروجی آن می باشد. از آنجا که محاسبات این الگوریتم بر اساس فاصله است نرمال سازی داده ها می تواند به بهبود عملکرد آن کمک کند.

فاز یادگیری (training phase) الگوریتم، شامل ذخیره سازی بردارهای ویژگی و برچسب دسته نمونه های اولیه است. در فاز طبقه بندی، k یک ثابت توسط کاربر تعریف می شود و بردار بدون برچسب (نقطه تست) از دسته ای است که بیشترین تعداد را در k نزدیکترین همسایه آن نقطه داشته باشد. به این ترتیب برچسب نقطه تست نیز مشخص می شود.

داده هایمان را به دو کلاس متغیر و نامتغیر تقسیم میکنیم و آنها را به صورت بردارهای 12 بعدی میبینیم. مسئله عملاً 12 درجه آزادی دارد که تمام فاکتورها مستقل از هم هستند. میانگین بردارهای هر کدام از دو گروه را به دست می آوریم. در نهایت دو بردار 12 بعدی داریم به عنوان مرکز کلاس ها. سپس فاصله داده های تستمان تا این دو را میسنجیم و بقیه نمونه ها را برچسب میزنیم.

الگوریتم Random Forest:

در کل، معمولاً درخت تصمیمی که بیش از حد عمیق باشد الگوی دقیق نخواهد داشت: دچار بیش برارزش شده، و دارای سوگیری پایین و واریانس بالا میباشد. جنگل تصادفی روشی است برای میانگین گیری با هدف کاهش واریانس با استفاده از درخت های تصمیم عمیقی که از قسمت های مختلف داده آموزشی ایجاد شده باشند. در این روش معمولاً افزایش جزئی سوگیری و از دست رفتن کمی از قابلیت تفسیر اتفاق افتاده اما در کل عملکرد مدل را بسیار افزایش خواهد داد. در این روش با بهره گیری از معیارهایی مثل gini و entropy و به کمک مینیم کردن آنها، پارامترهای ایده آل برای تصمیمات را انتخاب کردیم. عملاً سوالهایی که کمترین احتمال ایجاد خطا در فرآیند را دارند در نظر میگیریم.

الگوریتم Naive bayes:

در یادگیری ماشین به گروهی از دسته بندی کننده های ساده بر پایه احتمالات گفته می شود که با فرض استقلال متغیرهای تصادفی و براساس قضیه بیز ساخته می شوند. به طور ساده روش بیز روشی برای دسته بندی پدیده ها، بر پایه احتمال وقوع یا عدم وقوع یک پدیده است. این روش از ساده ترین الگوریتم های پیش بینی است که دقت قابل قبولی هم دارد.

الگوریتم SVM:

این روش از جمله روش های نسبتاً جدیدی است که در سال های اخیر کارایی خوبی نسبت به روش های قدیمی تر برای طبقه بندی نشان داده است. مبنای کاری دسته بندی کننده SVM دسته بندی خطی داده ها است و در تقسیم خطی داده ها سعی می کنیم خطی را انتخاب کنیم که حاشیه اطمینان بیشتری داشته باشد. حل معادله پیدا کردن خط بهینه برای داده ها به وسیله روش های QP که روش های شناخته شده ای در حل مسائل محدودیت دار هستند صورت می گیرد. قبل از تقسیم خطی برای اینکه ماشین بتواند داده های با پیچیدگی بالا را دسته بندی کند داده ها را به وسیله تابع phi به فضای با ابعاد خیلی بالاتر می بریم. برای اینکه بتوانیم مسئله ابعاد خیلی بالا را با استفاده از این روش ها حل کنیم از قضیه دوگانگی لاگرانژ برای تبدیل مسئله مینیم سازی مورد نظر به فرم دوگانگی آن که در آن به جای تابع پیچیده phi که ما را به فضایی با ابعاد بالا می برد، تابع ساده تری به نام تابع هسته که ضرب برداری تابع phi است ظاهر می شود استفاده می کنیم. از توابع هسته مختلفی از جمله هسته های نمایی، چند جمله ای و سیگموید می توان استفاده نمود. ماتریس الگو را آماده می کنیم تابع کرنلی را برای استفاده انتخاب می کنیم. پارامتر تابع کرنل و مقدار C را انتخاب می کنیم. برای محاسبه مقادیر α_i الگوریتم آموزشی را با استفاده از حل کننده های QP

اجرا می‌کنیم. داده‌های جدید با استفاده از مقادیر α_i و بردارهای پشتیبان می‌توانند دسته بندی شوند. آموزش نسبتاً ساده است. برخلاف شبکه های عصبی در ماکزیمم محلی گیر نمی‌افتد. برای داده‌های با ابعاد بالا تقریباً خوب جواب می‌دهد. مصالحه بین پیچیدگی دسته‌بندی‌کننده و میزان خطا به‌طور واضح کنترل می‌شود. به یک تابع کرنل خوب و انتخاب پارامتر C نیاز دارد.

همچنین از این الگوریتم برای انتخاب اینکه کدام یک از الگوریتم های ما دقت بهتری داشتند استفاده کردیم. نتیجه این بود که تمامی الگوریتم ها دقت 100 داشتند به جز KNN که دقتش 99.7% بود.