

سوال اول

۱: پیمانه جمع آوری داده ای که صدق داشته باشد، می توان پیمانه داده گشته را با ویژگی های ساختاری تحت عنوان Feature engineering و ویژگی های متریک را استخراج کرد و همچنین بعنوان راه حل اصلی و کلی این مشکل Imputation می توان گفت که در بخش عملی نیز انجام دادیم، مقادیر گشته را با راهکارهای مانند میانگین، میان یا KNN پیم می کنیم.

۲: می توان با oversampling نقاط داده از کلاس اقلیت را تکرار کرد تا با مجموعه داده، تعادل بیشتر برقرار بشود، همچنین می توان نقاط داده را از کلاس اکثریت کم کرد تا با سایر کلاس اقلیت سازگاری داشته باشد؛ همچنین می توان با دستاوت دادن وزن ها، وزن بالاتری به کلاس اقلیت اختصاص داد تا نمونه های کم تکرار تر هم افزایش پیدا کند.

۳: با روش های آماری یا دانش دامنه کار می توان موارد پرت را حذف کرد، همچنین استفاده از اوتلین های (smoothing) نظیر میانگین می توان نویز را هموار سازی کرد، راه دیگر هم استفاده از روش های است که نویز در آن ها تاثیر کمتری دارد مثل درخت های تصمیم یا Random forest.

۴: می توان از تکنیک های کاهش بعد نظیر PCA استفاده کرد همچنین با Feature selection می توان کلاس ساز باشد که زیر مجموعه ای از مفیدترین ها که ارتباط چندانی با یکدیگر ندارند را انتخاب کرد و بعد به کاهش ابعاد می شود. همچنین روش های منظم سازی که به مضاطر داشتن روابط پیچیده و متکی به داده ها این وابستگی را کاهش می دهد.

سوال ۲:

$$\text{آزمون دلتا} = a_1 x + y_0 + \text{میانگین معادلات} = \text{نمره آزمون}$$

می توان به صورت فوق آزمون دادن را با ضریب a_1 به آن افزود و برای یافتن a_1 از روش Least square method به این صورت عمل می کنیم که a_1 را یا ضریب که مجموع مجذور خطای میان نمرات امتحانی پیش بین شدن و نمرات واقعی به حداقل برسد و نحوه انجام آن، نخست محاسبه مجذور خطا برای هر نقطه سپس مجذور خطا برای همه نقاط را محاسبه کرد و ضرایب را تنظیم کرد و این بار مجموع مربعات خطا ها را محاسبه می کنیم و این فرآیند را تکرار می شود تا ترکیبی از ضرایب یا ضریب تا مجذور خطاها کم شود. در روش گزاردان نزولی ضریب از راه های بهینه سازی بوده. بهت دیگرادان را حساب کرده به گونه ای که مجموع مربعات خطاها کمترین افزایش را دارد و ضرایب را در بهت محاسبه می کنند این روش تا رسیدن به کمینه ادامه میابد از روش های دیگر منظم سازی را میتوان بیان کرد مانند بهینه سازی با بارش چنالی به تابع هدف اضافه کرده که با پتانسیل دادن به مقادیر ضرایب بزرگ از بهت از بهت بزرگ آزمون گیری می کنند.

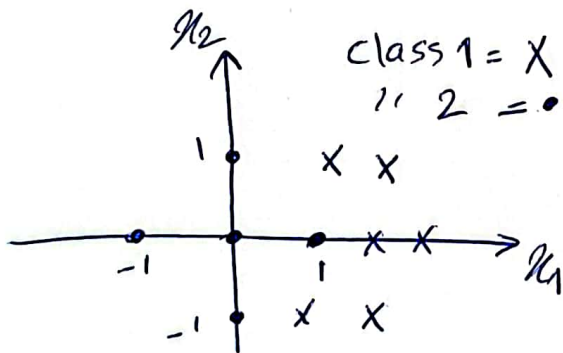
		Spam	Not Spam
Predicted {	Spam	۳۰۰	۳۰
	Not Spam	۲۰	۲۰۰
		Actual	

$$\text{Recall} = \frac{300}{(300+20)} = 93.7\%$$

$$\text{Precision} = \frac{300}{(300+30)} = 90.9\%$$

$$\text{Accuracy} = \frac{300+200}{(300+20+30+200)} = 91\%$$

$$\text{F1-score} = \frac{2 \times (0.909 \times 0.937)}{(0.909 + 0.937)} = 92\%$$



KNN : سوال اول:
داده تست (0.15, 0)
فاصله اقلیدسی:

کلاس 1

$$\begin{cases} (1, 1) = 1/5 \\ (1, 2) = 1/5 \\ (2, 1) = 1/5 \\ (2, 2) = 1/1 \checkmark \\ (3, 1) = 1/5 \\ (3, 2) = 2/5 \end{cases}$$

خواص اقلیدسی نقطه تست را از هر کدام از نقاط کلاس های تابع
با توجه به $k=3$
نزدیک ترین ها با \checkmark مشخص شدند
در مجموع دو تای آنها متعلق به
کلاس 2 بود پس کلاس 2 طبقه بندی می شود.

کلاس 1

$$\begin{cases} (1, 1) = 1/5 \\ (1, 2) = 2 \\ (2, 1) = 2/5 \\ (2, 2) = 2 \checkmark \\ (3, 1) = 2/5 \\ (3, 2) = 2 \end{cases}$$

کلاس 2

$$\begin{cases} (0, 2) = 0 \checkmark \\ (-1, 2) = 1/5 \\ (0, 1) = 1 \\ (1, 0) = 0/5 \checkmark \\ (0, -1) = 1 \end{cases}$$

نقاط \checkmark موارد 3 نقطه با نزدیک ترین فاصله هستند که معیار این نیز به کلاس 2 طبقه بندی می گردد.

سوال اول:

— Support Vector ها نزدیک ترین نقطه به مرز تقسیم (hyperplane) هستند که بین کلاس‌های مختلف

قرار دارند این نقاط برای تعیین مرز تقسیم گیری و تعیین حاشیه‌های طبقه‌بندی ضروری هستند.

مثلاً در همان مثال اسیل Ham و spam اسیل‌های هزینه‌ها + و - و غیره آن - در SVM یک ابرصفحه می‌باشد که دو کلاس را با جدا کردن حاشیه جدا می‌کند، این حاشیه فاصله بین ابرصفحه نزدیک ترین نقاط داده از هر کلاس است این نزدیک ترین نقاط که حاشیه را تعریف می‌کنند Support Vector هستند.

— برای داده‌های بسیار بزرگ - داده‌های با ابعاد زیاد - داده‌های غیر خطی با مرزهای تقسیم‌گیری پیچیده - داده‌های نامتعادل که یک کلاس از کلاس دیگر بسیار بیشتر است - داده‌های نویز

— آنها داده‌های غیر خطی قابل جداسازی را طبقه‌بندی می‌کنند. یک تابع کرنل برای نگاشت داده‌های ورودی از فضای ویژگی اصلی به فضایی با ابعاد بالاتر استفاده می‌شود و کلاس‌ها به صورت خطی قابل تفکیک می‌باشند و این تبدیل به آن اجازه می‌دهد که یک ابرصفحه با جداسازی آنها باید اکثر حاشیه‌ها را

— hard SVM ها هدف کمینه کردن نرم weight vector است و با فرض ابرصفحه با کمترین حاشیه کلاس‌ها را در داده‌های آموزشی جدا می‌کند بدون در نظر گرفتن Margin Violation

Soft SVM : (الگوریتم بخش قبلی را با اجازه نقض برخی حاشیه و طبقه‌بندی

تفاوت اصلی آنها در برخورد با حاشیه‌ها و طبقه‌بندی‌ها است در SVM ها حاشیه‌ها را نمی‌پذیرد و سعی می‌کند ابرصفحه‌ای که کلاس‌ها را بدون هیچگونه طبقه‌بندی نادرست از هم جدا می‌کند ولی Soft با صرفی یک پا ماست برخی موارد نقض و اشتباه را می‌پذیرد در hard به نویزها و موارد پیرامون تر است و Soft انعطاف پذیر است.

