

سوال اول

1) اگر تعداد این داده‌های ناموجود کم باشد می‌توان آن‌ها را با چیزهایی مانند میانگین و میان برکرد یا اگر در سطرهای که تعداد این ویژگی‌ها ناموجود زیاد باشد، آن سطر را حذف کرد.

ولی اگر تعداد داده‌های حذف زیاد باشد، این ویژگی را حذف می‌کنیم و یا آن را جایگزین نمی‌کنیم.

2) این مشکل زمانی پیش می‌آید که تعداد نمونه‌ها در هر کلاس نامتوازن باشد. برای حل این مشکل می‌توان از دو راه زیر استفاده کرد:

① Oversampling : تعداد نمونه‌های مربوط به کلاس کم نمونه را افزایش دهیم (با روشی مانند تولید رندوم داده)

② Undersampling : نمونه‌های مربوط به کلاس با نمونه زیاد را حذف کنیم (برای این روش هم می‌توان از حذف رندوم استفاده کرد)

3) می‌توان نمونه‌های نویزی را حذف کرد و یا داده‌ها با استفاده از فیلترهای میانگین‌گیری تقطیع شوند.

همچنین با استاندارد کردن و یا نرمالایز کردن می‌توان تأثیر نویز را کاهش داد.

همچنین تقطیع داده‌ها را با بردن به فضای فوری (تبدیل فوری) می‌توان انجام داد.

4) با توجه به کورولیشن می‌توان ~~از بین~~ ویژگی‌ها می‌توان ~~از بین~~ ویژگی‌ها را نگه داشت که تأثیر بیشتری در مدل ما دارند و بقیه را حذف کنیم.

سوال دوم :

برای درک ارتباط دو متغیر و یا تأثیر هر متغیر بر روی target ، کورولیشن گیری همیشه کمک می کند.

برای (least square method) بدین صورت ضرایب بدست می آید که باید طوری ضرایب

انتخاب شوند که ~~اختلاف~~ مربع های اختلاف مقادیر اصلی و مقادیر پیش بینی شده کمینه شود.

برای Gradient descent ، اینکار را آنقدر تکراری کنیم که به کمترین گرادیان ممکن برسیم :

- سب را (مقدار گرادیان) را ~~را~~ برای هر پارامتر حساب می کنیم.

- مقادیر ~~پارامترها~~ پارامترها را در جهت کمترین مقدار گرادیان تغییر می دهیم.

از روش های دیگری توان Least Absolute Deviation و روش نیوکمن را نام برد که

اولی مجموع اختلاف مقادیر اصلی و پیش بینی شده را کم می کند و با این روش ضرایب

را بدست می آورد و روش دوم نیز گفته می شود که بهترین ضرایب را پیدا خواهد کرد.

سوال سوم :

$$\text{Accuracy} = \frac{300 + 200}{300 + 200 + 20 + 30} = \frac{500}{550} = \frac{10}{11} \approx 0,91$$

$$\text{Precision} = \frac{300}{300 + 30} = \frac{300}{330} = \frac{10}{11} \approx 0,91$$

$$\text{Recall} = \frac{300}{300 + 20} = \frac{300}{320} = \frac{15}{16} \approx 0,94$$

$$F_1\text{-score} = 2 \times \frac{\frac{10}{11} \times \frac{15}{16}}{\frac{10}{11} + \frac{15}{16}} = \frac{2 \times 176 \times 150}{88 \times 325} = \frac{12}{13} \approx 0,92$$

مقادیر پیش بینی شده	$Tp=300$	$Fp=30$
	$FN=20$	$TN=200$
	مقادیر اصلی	

KNN : فاصله‌های اقلیدسی :

$$d_1(x, y) \Rightarrow \text{فاصله اقلیدسی } (x, y) \\ (1/2, 0)$$

$$d_1(0, 0) = 1/2 \text{ و } d_1(0, 1) = \frac{\sqrt{5}}{2} \text{ و } d_1(1, 0) = 1/2 \text{ و } d_1(1, 1) = \frac{\sqrt{5}}{2} \text{ و} \\ d_1(0, -1) = \frac{\sqrt{5}}{2} \text{ و } d_1(1, -1) = \frac{\sqrt{5}}{2} \text{ و } d_1(2, 0) = 3/2 \text{ و } d_1(3, 0) = 5/2$$

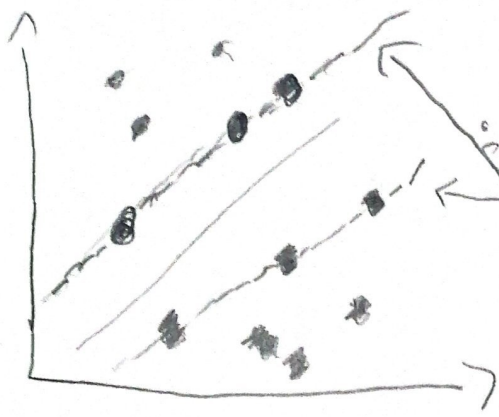
بنابراین به علت وجود $(0, 1)$ و $(1, 0)$ و $(0, -1)$ که این نمونه به کلاس دوم تعلق دارد.

فاصله‌های منجهن :

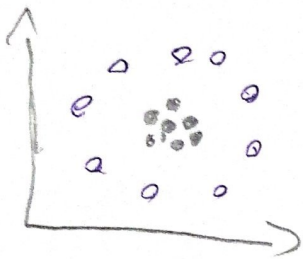
$$d_2(x, y) \Rightarrow \text{فاصله منجهن } (x, y) \\ (1/2, 0)$$

$$d_2(0, 0) = 1/2 \text{ و } d_2(0, 1) = 3/2 \text{ و } d_2(1, 0) = 1/2 \text{ و } d_2(1, 1) = 3/2 \text{ و} \\ d_2(0, -1) = 3/2 \text{ و } d_2(1, -1) = 3/2 \text{ و } d_2(2, 0) = 3/2 \text{ و } d_2(3, 0) = 5/2$$

بنابراین به علت وجود نقاط $(0, 1)$ و $(1, 0)$ و یکی از نقاطی که فاصله‌ی $3/2$ دارند که با زهم در کلاس دوم قرار می‌گیرد.



• به فضای گفته می شود که به حداقل کننده نزدیکترین مثال \leftarrow support vector



• برای داده های که صورت خطی تشکیل پذیر نیستند مانند

و همچنین برای داده های خیلی بزرگ یا نویزی یک مناسب نیست

• یکی یک تابع ریاضی است که داده ها را به بعد بالاتر تبدیل می کند.

این تبدیل این امکان را می دهد که در فضای جدید، الگوریتم خطی را روی داده ها اجرا کنیم. به عبارت دیگر کنترل ها با اجازه می دهند که با استفاده از الگوریتم های خطی، مسائل غیر خطی را حل کنیم.

• هدف اصلی Hard SVM، حاشیه ی بین دو کلاس را حداکثر کردن است و

حیثی اعتبار طاق بندی در آن معیار نیست و تمام داده ها باید بدرستی دسته بندی

شوند اما در Soft SVM، به داده ها اجازه می دهیم در حاشیه ی

~~Soft SVM~~ قرار بگیرند یا حتی در داخل حاشیه باشند.

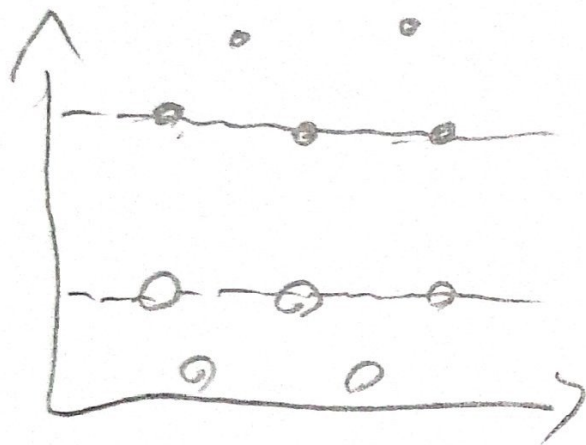
Soft SVM به ما امکان می دهد که با اجازه ی اعتبار طاق بندی، مدلی با خصوصیت بیشتر داشته باشیم و برای داده های noisy یا دارای نویز مناسب تر است.

• خطی را پیدا می‌کند که دو کلاس را از هم جدایی کند و از Support Vector های دو طرف

بیشترین فاصله را داشته باشد. اگر Support Vector ها و خط پیدا شده هر سه مواردی

باشند، خط پیدا شده در وسط Support Vector ها قرار خواهد گرفت تا بیشترین

فاصله را داشته باشد.



خط پیدا شده