

Q1)

وجود نداشتن ویژگی‌ها: روش‌های پر کردن داده‌های ناقص مانند میانگین‌گیری،

جایگزینی با صفر، الگوریتم‌های پیش‌بینی مانند شبکه‌های عصبی و رگرسیون و یا حذف

داده‌های ناقص از دیتاست

ناستونز بودن توزیع داده‌گان: استفاده از روش‌هایی که مبتنی بر هزینه‌اند و به نمونه‌

های مختلف وزن‌های متفاوتی می‌دهند؛ یا روش‌های ستادسازی مانند oversampling

و undersampling می‌توانند به نرسال کردن توزیع داده‌گان در کلاک‌ها کمک کند.

دبود نویز: استفاده از روش‌های پیش‌پردازش مانند فیلترهای نویزی، فیلتر میانگین،

نیلتر گوس و یا بهره بردن از الگوریتم‌های ستاد به نویز مانند رگرسیون و یا NN

و ویژگی‌های همبسته: استفاده از روش‌های کاهش بُعد شد PCA و یا حذف یکی از

ویژگی‌های همبسته می‌تواند تا حد این Correlate را کم کند.

Q21

با انفاذ شدن یک ویژگی جدید می‌توان از رگرسیون چند جمله‌ای استفاده کرد:

$$\text{نمره آزمون} = B_0 + (B_1 \times \text{آزمون داده}) + (B_2 \times \text{مطالعه})$$

برای ضرایب مناسب دو روش وجود دارد:

1) Least Square method $\rightarrow B = X^T Y^{-1} (X^T X)$

X = ماتریس ویژگی‌ها ، Y = بردار مقادیر واقعی ، B = بردار ضرایب

2) decent Gradient \Rightarrow

درین روش از یک نقطه تعادلی شروع به update کردن ضرایب می‌کنیم تا مجموع

مربعات اختلاف را کمینه کنیم! در روش اول با مشتق‌گیری نسبت به ویژگی‌ها

به دنبال کمینه شدن مجموع هستیم ولی در روش دوم با آپدیت کردن مقادیر ویژگی

ها به این minimum می‌رسیم.

روش‌های دیگر:

روش‌های رگرسیونی وجود دارد که علاوه بر کمینه کردن مجموع مربعات اختلاف،

به منظور پیشگیری از overfitting برای ضرایب بزرگ جریمه مناسب می‌شود.

Q3) درستی درک عدد

صفتی عدد	200	30
صفتی درست	20	300

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{300}{300 + 20} = 0.94$$

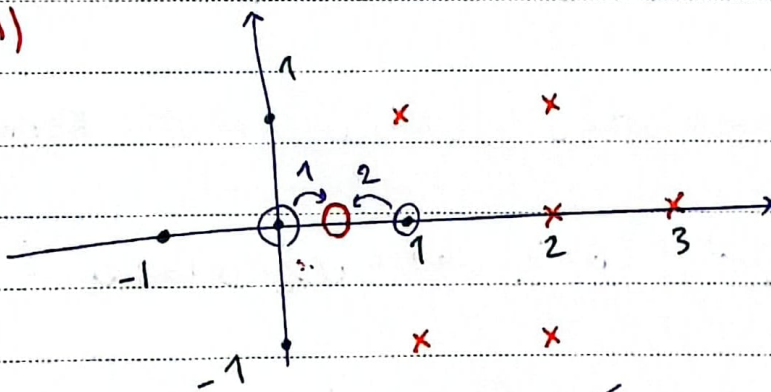
$$\text{Precision} = \frac{TP}{TP + FP} = \frac{300}{300 + 30} = 0.91$$

$$\text{accuracy} = \frac{TP + TN}{\text{all}} = \frac{300 + 200}{550} = 0.91$$

$$F_1\text{-score} = 2 \frac{Pr \times Re}{Pr + Re} = 2 \frac{0.91 \times 0.94}{0.91 + 0.94} = 2 \frac{0.18206}{1.85} = 0.920$$

⊕ KNN

Q1)



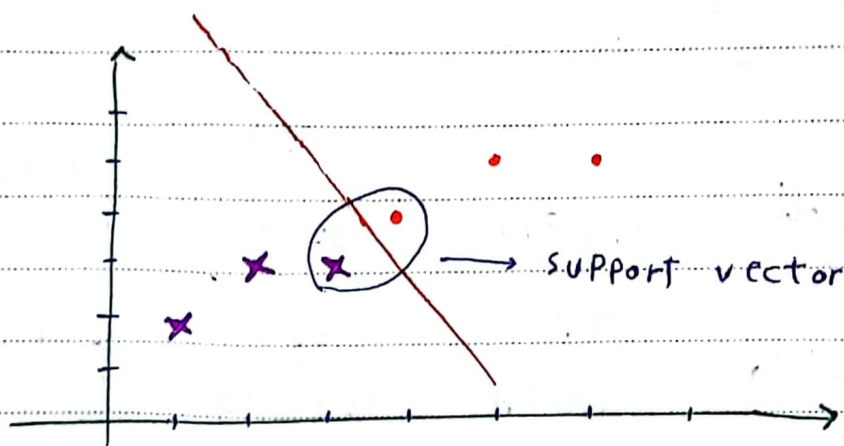
⊕ سوین همسایه هر چه باشد، چون دوستی اول در کلاسی دایره

هستند، داده ما نیز در دسته 2 یا همان • قرار می گیرد!

⊕ SVM

به نقاشی که برای تعیین margin بین دسته های مختلف مورد استفاده

قرار می گیرند Support vector می گوئیم:



در داده های پرتراکم با جمع زیاد، داده های نامتوازن در توزیع، دادگان دایرهای

و اثرات همبسته و پیچیده و در جاهایی که کاهش بعد منجر به از دست رفتن

دیتا شود، این روش مناسب نمی باشد!

Kernel (=) به ما این امکان را می دهد که با استفاده از توابع مشخص، داده ها

را به یک فضای ویژگی با ابعاد بالاتر نگاشت کنیم به طوری که دادگان در این فضا

بصورت خطی جدا پذیر شوند! سپس SVM را روی داده های نگاشت شده

اعمال می کنیم. در واقع یک خط در ابعاد بالا، برابر یک نفر در فضای 3 بعدی است!

① Hard-SVM \Rightarrow هیچ داده‌ای در margin قرار نمی‌گیرد و زنی می‌شود

دادگان بصورت خطی جدا پذیرند؛ اگر واقعاً خطی جدا نشوند

این نوع SVM دچار overfitting خواهد شد.

② Soft-SVM \Rightarrow به خلاف قبل اجازه می‌دهد در margin داده قرار

بگیرد و مصدق آنرا طوری تغییر می‌دهد که پس از

classification مقدار خطا و overfitting به حداقل برسد و طبقاً یک سری داده را

نویز و داند margin تشخیص می‌دهد.

