

(6) مزایای weighted :

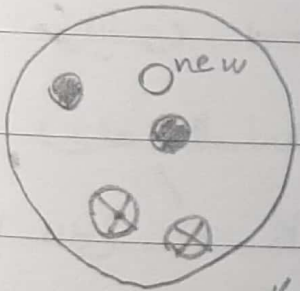
در روش weighted، record هایی که به new record نزدیکترند وزن بیشتری می‌گیرند، در حالی که در روش unweighted تفاوتی میان رکوردهای نزدیک

و دور وجود ندارد (همانند شکل زیر)، رکورد new فارغ از فاصله اش با رکورد A به B predict می‌شود.



A ●
B ⊗

در روش weighted احتمال بروز tie ضعیف‌تر از unweighted است.



در روش weighted، رکورد new به روشی که فاصله اش از A و B یکسان است، به B predict می‌شود.

اما در روش unweighted، چنین حالتی را tie می‌نامیم و new به هیچ کدام از آنها A و B predict نمی‌شود.

که منطقی هم هست

انبار سوال کنجایب weighted

البته در روش weighted هم به دست حالتی پیش بیاید که فاصله یک رکورد

train با new record صفر باشد، در نتیجه وزن بی نهایت برای آن یک



رکورد بدست می آید (چون سوال بپرسد، لباس A برای new).

در این حالت بقیه رکوردها در نظر گرفته نمی شوند، هر چند که

فردا حتی شان زیر باشد (عیب روش weighted)

(7) الوریتم DM باید بتواند ^(ی) rare (نایاب) را هم ^(ی) تشخیص دهد

با تعداد رکورد کم

اند، برای همین dataset باید balance شود. حقیقت database

s.a.

هم محدود است و برخی موارد دسترسی به برخی رکوردها نیاز است.

(9) معایب استفاده از K کوچک: نتیجه prediction نتایج ضعیف

noise و outlier هستند مراد بگیرد. برای مثال اگر K را مساوی 1

انتخاب کنیم و الگوریتم فقط یک رکورد از train را بعنوان نتیجه برگرداند

(over fitting) یا حفظ کردن (memorizing) حافظه می دهد.

مزایای استفاده از K کوچک: هزینه محاسباتی الگوریتم در مقایسه

با حالتی که K مقادیر بزرگی می گیرد کاهش می یابد (دسترسد KNN الگوریتم

instance-based است، یعنی به ازای هر رکورد مجموعه test و الگوریتم درباره

شان می شود.)

(10) برای اولویت دادن به ویژگی (attribute) هایی که ارتباط بیشتری

با target variable دارند. زیرا KNN به صورت default همه ویژگی ها را با
اعمال می کند
s.a.m

انجام ۱۵) بنابراین امکانش هست که چندین رکورد در attribute های

مهم به new record بسیار نزدیک باشند، اما در سایر ویژگی‌ها غیر ضروری

فاصله داشته باشند. در این صورت اگر stretch axes برای متغیرهای مهم رخ ندهد.

نتیجه prediction اشتباه می‌شود.

810198358

امیرمهدی انصاری پور

12) با درست کردن dataset از رکورد های R2:R10 و ران کردن الگوریتم KNN به ازای $k = 2$:

```
6 R3 = c(0.13, 0.19, 1)
7 R4 = c(0.65, 0, 1)
8 R5 = c(0.06, 0.92, 0)
9 R6 = c(0.38, 0.39, 0)
10 R7 = c(0.72, 0.19, 0)
11 R8 = c(0.75, 1, 1)
12 R9 = c(0.63, 0.9, 1)
13 R10 = c(1, 0.48, 1)
14
15 data = rbind(R2, R3, R4, R5, R6, R7, R8, R9, R10)
16 dimnames(data) = list(c("R2 (Bad loss)", "R3 (Bad loss)", "R4 (Bad loss)", "R5 (Bad loss)", "R6 (Good risk)", "R7 (Good risk)", "R8 (Good risk)", "R9 (Good risk)", "R10 (Good risk)"),
17                        c("Age(MMN)", "Income(MMN)", "Marital(MMN)"))
18
19 trueClass = c("Bad loss", "Bad loss", "Bad loss", "Bad loss", "Good risk", "Good risk", "Good risk", "Good risk", "Good risk")
20 knn(data, new, cl= trueClass, k = 2, prob = TRUE)
21
22
23
```

22:50 (Top Level) R Script

Console Terminal Background Jobs

R 4.2.2 ~ /

```
> data
  Age(MMN) Income(MMN) Marital(MMN)
R2 (Bad loss)      0.25      0.01      1
R3 (Bad loss)      0.13      0.19      1
R4 (Bad loss)      0.65      0.00      1
R5 (Bad loss)      0.06      0.92      0
R6 (Good risk)     0.38      0.39      0
R7 (Good risk)     0.72      0.19      0
R8 (Good risk)     0.75      1.00      1
R9 (Good risk)     0.63      0.90      1
R10 (Good risk)    1.00      0.48      1
> knn(data, new, cl= trueClass, k = 2, prob = TRUE)
[1] Good risk
attr(,"prob")
[1] 0.5
Levels: Bad loss Good risk
> |
```

همانطور که دیده میشود، کلاس رکورد 1 Good risk پیش بینی شده است.