

تمرین سری دوم - یادگیری ماشین

(1)

(الف)

در این سوال برای استفاده از LOOCV باید هر داده را یک بار به عنوان تست قرار داده و باقی را برای آموزش در نظر بگیریم. اگر مقدار k برابر تمام داده ها باشد تمام مقادیر تست منفی در نظر گرفته میشود که مقدار خطای 0.4 دارد.

اگر مقدار $k = 1$ در نظر بگیریم داده های سمت راست که همه منفی هستند و نزدیک ترین داده به آنها نیز منفی است همگی در حالت تست مقدار منفی میگیرند اما در مورد داده های سمت چپ ، اگر داده منفی تست باشد طبق آموزش مقدار آن مثبت در نظر گرفته شده و اگر داده های مثبت تست باشند ، نزدیک ترین مقدار به آنها منفی است پس مقدار آن ها نیز منفی میشود پس نیمی از حالات خطا دارند. مقدار خطای 0.5 را دارد.

اگر $k = 3$ باشد در این صورت داده های سمت راست درست تخمین زده میشوند اما در مورد داده های سمت چپ، مقادیر مثبت به علت اینکه در تمام حالات نزدیک 2 مثبت و 1 منفی هستند مقدار مثبت میگیرند و درست هستند اما در مورد داده منفی ، مقدار آن مثبت در نظر گرفته میشود. پس مقدار خطا برای این حالت 0.1 است.

از حالت $k = 3$ تا $k = 6$ مقدار برای مقادیر مثبت درست محاسبه میشود و برای مقدار منفی در سمت چپ تا $k = 9$ مقدار خطا در نظر گرفته میشود. پس از $k = 3$ تا $k = 6$ مقدار خطا ثابت است و برابر 0.1 است و بعد از آن به سمت 0.4 زیاد میشود.

اما از میان این k ها ، مقدار $k = 3$ بهترین k است. زیرا با افزایش k به سمت بایاس بالاتر میرویم و پیچیدگی نیست بیشتر میشود زیرا به داده ها با فواصل بیشتر میرویم. پس بهترین k برای این سوال مقدار $k = 3$ است.

(ب)

برای یافتن بهترین مقدار k میتوانیم از الگوریتم Grid search استفاده کرد. در روش جستجوی شبکه، شبکه ای از مقادیر ممکن برای هایپرپارامترها ایجاد می کنیم. هر تکرار ترکیبی از فرایپارامترها را به ترتیب خاصی امتحان می کند. این مدل را بر روی هر ترکیبی از ابرپارامترهای ممکن برازش می دهد و عملکرد مدل را ثبت می کند. در نهایت، بهترین مدل را با بهترین هایپرپارامترها برمی گرداند.

(2)

KNN یک الگوریتم تمایز است زیرا احتمال مشروط یک نمونه متعلق به یک کلاس معین را مدل می کند.

درخت های تصمیم متمایز هستند زیرا مرزهای صریح بین کلاس ها را یاد می گیرند.

(3)

(الف)

$$\begin{aligned} \text{3 الف)} \quad \frac{d\sigma(a)}{da} &= \sigma(a)(1-\sigma(a)) \Rightarrow \frac{1}{1+e^{-a}} \left(\frac{1+e^{-a}-1}{1+e^{-a}} \right) \quad (1) \\ \sigma(a) &= \frac{1}{1+e^{-a}} \Rightarrow \frac{d\sigma(a)}{da} \Rightarrow - \frac{1}{(1+e^{-a})^2} (e^{-a}) \\ &\Rightarrow \frac{e^{-a}}{(1+e^{-a})^2} \xrightarrow[\text{! اینجا فرمول داریم}]{\text{مربع فرض ابر به صورت}} \frac{1+e^{-a}-1}{(1+e^{-a})(1+e^{-a})} \quad (2) \\ \text{1} &= \text{2} \Rightarrow \boxed{\sigma(a)(1-\sigma(a))} \end{aligned}$$

(ب)

$$\Rightarrow \prod P(y_i=1|x)^{y_i} P(y_i=0|x)^{1-y_i} \quad (1) \quad (3)$$

$$P(y_i=1|x) = \sigma(\omega^T x) \quad (2) \quad \text{likelihood} = \prod \sigma(\omega^T x)^{y_i} (1 - \sigma(\omega^T x))^{1-y_i}$$

$$-\log \rightarrow -\log(\text{likelihood}) = -\sum y_i \log(\sigma(\omega^T x)) - \sum \log(1 - \sigma(\omega^T x))^{1-y_i}$$

$$\Rightarrow -\sum y_i \log(\sigma(\omega^T x)) - \sum (1-y_i) \log(1 - \sigma(\omega^T x)) \quad (3)$$

$$\sigma(\omega^T x) = \frac{1}{1+e^{-\omega^T x}}, \quad 1 - \sigma(\omega^T x) = \frac{e^{-\omega^T x}}{1+e^{-\omega^T x}} \quad (4)$$

$$(3,4) \Rightarrow \left[\sum y_i \log(1+e^{-\omega^T x}) - \sum (1-y_i) \log\left(\frac{e^{-\omega^T x}}{1+e^{-\omega^T x}}\right) \right]$$

$$\Rightarrow \sum \omega^T x_i + y_i \omega^T x_i + \log(1+e^{-\omega^T x})$$

$$\text{likelihood} = \sum y_i \log(1 + e^{-w^T x_i}) - \sum (1 - y_i) \log\left(\frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}}\right) \quad (1)$$

$$\Rightarrow \sum y_i \log(1 + e^{-w^T x_i}) - \sum (1 - y_i) (\log e^{-w^T x_i} - \log(1 + e^{-w^T x_i}))$$

$$(2) \Rightarrow \sum (1 + y_i) (-w^T x_i - \log(1 + e^{-w^T x_i})) \Rightarrow \sum (w^T x_i + y_i (-w^T x_i) + \log(1 + e^{-w^T x_i}) + y_i \log(1 + e^{-w^T x_i}))$$

$$(1), (2) \Rightarrow \sum w^T x_i - y_i w^T x_i + \log(1 + e^{-w^T x_i}) = -\log(\text{likelihood}) = L_{H_{\text{new}}}$$

$$\frac{\partial L_{H_{\text{new}}}}{\partial w} \Rightarrow \frac{\partial [\sum w^T x_i - y_i w^T x_i + \log(1 + e^{-w^T x_i})]}{\partial w} \Rightarrow \sum x_i - y_i x_i + \frac{-x_i e^{-w^T x_i}}{1 + e^{-w^T x_i}}$$

$$\Rightarrow -\sum \left(y_i + 1 + \frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}} \right) x_i = -\sum \left(y_i + \frac{1 + e^{-w^T x_i} - e^{-w^T x_i}}{1 + e^{-w^T x_i}} \right) x_i$$

$$\Rightarrow -\sum \left(y_i - \frac{1}{1 + e^{-w^T x_i}} \right) x_i \quad (3) \Rightarrow -\sum (y_i - \hat{y}_i) x_i$$

$$\sigma(w^T x) = \frac{1}{1 + e^{-w^T x}} \quad (3) = \hat{y}_i$$

$$(3) \text{ تابع سلیموره: } 1 + e^{-\omega x}$$

چراغ تابع پراستری به برعکس تابعی تأثیر می‌گذارد پراستر ω است. هرچه مقدار ω بیشتر شود

تابع سلیموره بیشتر به سمت بیش پرازش می‌رود. هرچه مقدار ω کمتر باشد حالت تابع به صورت (ک)

خواهد بود. اجازه آخرش می‌دهد اما اگر مقدار ω بیشتر شود این معنی به سمت حالت (ل)

می‌رود به بیش پرازش زیادی دارد.

$$P(\text{buy} = \text{yes}) = \frac{9}{14}$$

$$P(\text{buy} = \text{No}) = \frac{5}{14}$$

(الف.4)

$$X_1 \Rightarrow \begin{cases} P(\text{youth} | \text{yes}) = \frac{2}{9} \\ P(\text{high} | \text{yes}) = \frac{2}{9} \\ P(\text{yes} | \text{yes}) = \frac{6}{9} \\ P(\text{fair} | \text{yes}) = \frac{6}{9} \end{cases}$$

$$\begin{cases} P(\text{youth} | \text{No}) = \frac{3}{5} \\ P(\text{high} | \text{No}) = \frac{2}{5} \\ P(\text{yes} | \text{No}) = \frac{1}{5} \\ P(\text{fair} | \text{No}) = \frac{2}{5} \end{cases}$$

$$\left(\frac{2}{9}\right)^2 \times \left(\frac{6}{9}\right)^2 \times \frac{9}{14} = 0.0141 \checkmark$$

$$\frac{3}{5} \times \frac{4}{25} \times \frac{1}{5} \times \frac{5}{14} = 0.0115$$

$$X_2 \Rightarrow \begin{cases} P(\text{senior} | \text{yes}) = \frac{3}{9} \\ P(\text{low} | \text{yes}) = \frac{3}{9} \\ P(\text{No} | \text{yes}) = \frac{3}{9} \\ P(\text{ex} | \text{yes}) = \frac{3}{9} \end{cases}$$

$$\begin{cases} P(\text{senior} | \text{No}) = \frac{2}{5} \\ P(\text{low} | \text{No}) = \frac{1}{5} \\ P(\text{No} | \text{No}) = \frac{4}{5} \\ P(\text{ex} | \text{No}) = \frac{3}{5} \end{cases}$$

$$\left(\frac{3}{9}\right)^4 \times \frac{9}{14} = 0.079$$

$$\frac{5}{14} \times \frac{2}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} = 0.0137 \checkmark$$

$$X_3 \Rightarrow \begin{cases} P(\text{middle} | \text{yes}) = \frac{4}{9} \\ P(\text{medium} | \text{yes}) = \frac{4}{9} \\ P(\text{No} | \text{yes}) = \frac{3}{9} \\ P(\text{fair} | \text{yes}) = \frac{6}{9} \end{cases}$$

$$\begin{cases} P(\text{middle} | \text{No}) = \frac{0}{5} \\ P(\text{medium} | \text{No}) = \frac{2}{5} \\ P(\text{No} | \text{No}) = \frac{4}{5} \\ P(\text{fair} | \text{No}) = \frac{2}{5} \end{cases}$$

$P(\text{middle} | \text{No})$ is 0, but smoothed likelihood is not zero Prob. Error

$$\frac{0+L}{5+L \times K} \Rightarrow L=1, K=4 \Rightarrow \frac{1}{9}$$

if Prob. of 0 is Error

$$\left(\frac{4}{9}\right)^2 \times \frac{3}{9} \times \frac{6}{9} \times \frac{9}{14} = 0.281 \checkmark$$

$$\frac{1}{5} \times \left(\frac{2}{5}\right)^2 \times \frac{4}{5} \times \frac{5}{14} = 0.005$$

$$S(\text{age}) \Rightarrow S[\text{youth}] = [2+, 3-], S[\text{middle}] = [4+, 5-], S[\text{senior}] = [3+, 2-] \quad (\text{ب.4})$$

$$\Rightarrow E[\text{youth}] = -\left(\frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5}\right) = 0.97, \quad E[\text{middle}] = 0, \quad E[\text{senior}] = 0.97$$

$$E(\text{All}) = 0.94, \quad \text{gain}(\text{age}) = 0.94 - \frac{5}{14}(0.97) - \frac{5}{14}(0.97) = 0.226$$

$$S[\text{income}] \Rightarrow S[\text{high}] = [2+, 2-], S[\text{medium}] = [4+, 2-], S[\text{low}] = [3+, 1-]$$

$$E[\text{high}] = 1, \quad E[\text{medium}] = 0.917, \quad E[\text{low}] = 0.811$$

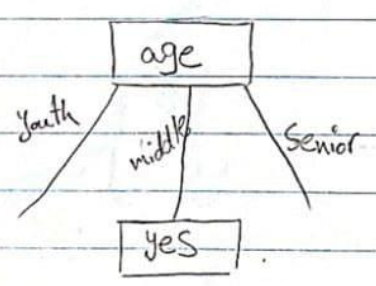
$$\text{gain}[\text{income}] = 0.94 - (0.917) \times \frac{6}{14} - (0.811) \times \frac{4}{14} = 0.31$$

$$S[\text{stu}] \Rightarrow S[\text{yes}] = [6+, 1-], S[\text{no}] = [3+, 4-] \Rightarrow E[\text{yes}] = 0.591, \quad E[\text{no}] = 0.984$$

$$\text{gain}[\text{stu}] = 0.94 - \frac{7}{14}(0.591) - \frac{7}{14}(0.984) = 0.153$$

$$S[\text{credit}] \Rightarrow S[\text{fair}] = [6+, 2-], S[\text{ex}] = [3+, 3-] \Rightarrow E[\text{fair}] = 0.811, \quad E[\text{ex}] = 1$$

$$\text{gain}[\text{credit}] = 0.94 - \frac{8}{14}(0.811) - \frac{6}{14} = 0.48$$



النتيجة النهائية هي age

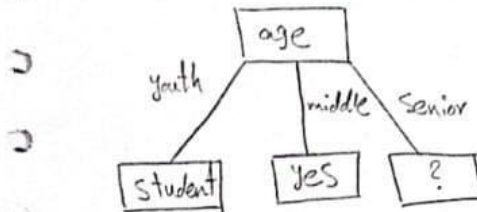
$$E(All) = .1970$$

با در نظر گرفتن $age = youth$

$$S[income] \Rightarrow E[high] = 0, E[medium] = 1, E[low] = 0 \Rightarrow \boxed{gain[income] = .57}$$

$$S[stu] \Rightarrow E[yes] = 0, E[no] = 0 \Rightarrow \boxed{gain[stu] = .1970} \checkmark$$

$$S[credit] \Rightarrow E[fair] = .197, E[ex] = 1 \Rightarrow \boxed{gain[credit] = .1972}$$



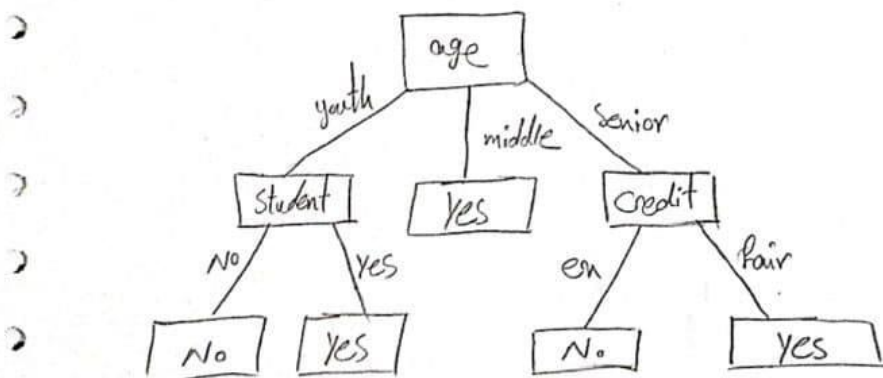
$$E(All) = .1970$$

با در نظر گرفتن $age = senior$

$$S[income] \Rightarrow E[medium] = .197, E[low] = 1 \Rightarrow \boxed{gain[income] = .1020}$$

$$S[stu] \Rightarrow E[yes] = .197, E[no] = 1 \Rightarrow \boxed{gain[stu] = .1020}$$

$$S[credit] \Rightarrow E[fair] = 0, E[ex] = 0 \Rightarrow \boxed{gain[credit] = .1970}$$



چون $age = senior$ است، $Student$ در سطح برابر به اشتباه طبقه‌بندی می‌شود، پس $income$ و stu را بررسی می‌کنیم.

پس $income$ و stu را بررسی می‌کنیم و $credit$ را بررسی می‌کنیم.

(5

الف)

```
=== Classifier model (full training set) ===

J48 pruned tree
-----

wage-increase-first-year <= 2.5: bad (15.27/2.27)
wage-increase-first-year > 2.5
| statutory-holidays <= 10: bad (10.77/4.77)
| statutory-holidays > 10: good (30.96/1.0)

Number of Leaves :    3
Size of the tree :    5

Time taken to build model: 0.02 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      42           73.6842 %
Incorrectly Classified Instances    15           26.3158 %
Kappa statistic                     0.4415
Mean absolute error                 0.3192
Root mean squared error             0.4669
Relative absolute error             69.7715 %
Root relative squared error        97.7888 %
Total Number of Instances          57

=== Detailed Accuracy By Class ===

      TP Rate  FP Rate  Precision  Recall   F-Measure  MCC      ROC Area  PRC Area  Class
      -----  -
      0.700    0.243    0.609    0.700    0.651     0.444    0.695    0.559    bad
      0.757    0.300    0.824    0.757    0.789     0.444    0.695    0.738    good
Weighted Avg.   0.737    0.280    0.748    0.737    0.740     0.444    0.695    0.675
```

ب)

=== Confusion Matrix ===

```
  a  b  <-- classified as
14  6 |  a = bad
 9 28 |  b = good
```

(ج

=== Classifier model (full training set) ===

J48 unpruned tree

```
wage-increase-first-year <= 2.5
|  education-allowance = yes
|  |  wage-increase-first-year <= 2.1
|  |  |  pension = none: bad (2.43/0.43)
|  |  |  pension = ret_allw: bad (0.0)
|  |  |  pension = empl_contr: good (3.16/1.5)
|  |  wage-increase-first-year > 2.1: bad (2.04/0.04)
|  education-allowance = no
|  |  contribution-to-health-plan = none: bad (3.39)
|  |  contribution-to-health-plan = half: good (0.18/0.05)
|  |  contribution-to-health-plan = full: bad (4.06)
wage-increase-first-year > 2.5
|  longterm-disability-assistance = yes
|  |  statutory-holidays <= 10
|  |  |  wage-increase-first-year <= 3: bad (2.0)
|  |  |  wage-increase-first-year > 3: good (3.99)
|  |  statutory-holidays > 10: good (25.67)
|  longterm-disability-assistance = no
|  |  contribution-to-health-plan = none: bad (4.07/1.07)
|  |  contribution-to-health-plan = half: bad (3.37/1.37)
|  |  contribution-to-health-plan = full: good (2.62)
```

Number of Leaves : 13

Size of the tree : 22

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

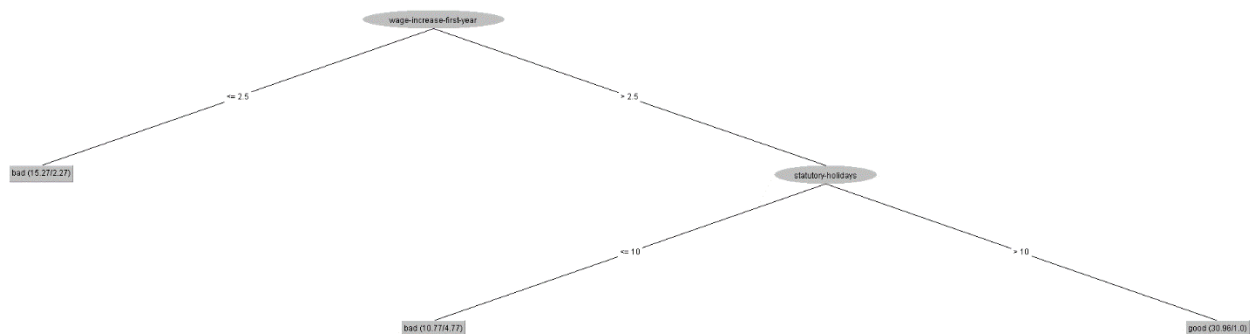
Correctly Classified Instances	45	78.9474 %
Incorrectly Classified Instances	12	21.0526 %
Kappa statistic	0.5378	
Mean absolute error	0.2677	
Root mean squared error	0.432	
Relative absolute error	58.5226 %	
Root relative squared error	90.4708 %	
Total Number of Instances	57	

(د)

بدون هرس :



با هرس :



در این دیتاست ، در حالت با هرس عمق درخت 2 است در صورتی که بدون هرس عمق 4 است ، همچنین تعداد ویژگی ها در با هرس برابر 2 است در صورتی که برای بدون هرس مقدار آن 6 است که نشان دهنده پیچیدگی است. برای داده های کمتر از 2.5 برای حالت بدون هرس ممکن است حالت خوبی هم باشد در صورتی که در حالت با هرس برای داده ها کمتر از 2.5 حالت بد است.

پیاده سازی:

(1)

(2)

(الف)

برای حل این مشکل ما میتوانیم به جای الگوریتم ID3 از الگوریتم C4.5 استفاده کنیم که در این الگوریتم مقادیر نامعلوم قابل حل هستند.

الگوریتم C4.5 با برگرداندن توزیع احتمال برچسب ها در زیر شاخه مشخصه ای که مقدار آن وجود ندارد، با مقادیر گمشده سروکار دارد. فرض کنید نمونه ای در داده های آزمایشی خود داشتیم که آب و هوا را آفتابی نشان می داد اما مقداری برای ویژگی رطوبت نداشت. همچنین، فرض کنید که داده های آموزشی ما دارای 2 نمونه بود که آب و هوا آن آفتابی، رطوبت زیر 75 و برچسب Play بود. علاوه بر این، فرض کنید داده های آموزشی 3 مورد داشتند که در آن آب و هوا آفتابی، رطوبت بالای 75 بود، و دارای برچسب Don't Play بود. بنابراین برای نمونه آزمایشی با ویژگی رطوبت از دست رفته، الگوریتم C4.5 توزیع احتمال $[0.4, 0.6]$ مربوط به $[Play, Don't Play]$ را برمی گرداند.

(ج)

داده های تست شده و مقادیر محاسبه شده برای در فایل Final.csv در پوشه قرار داده شده است.

```

For K = 1 Average of Validations 96.00000000000001
For K = 3 Average of Validations 94.66666666666666
For K = 5 Average of Validations 93.33333333333334
For K = 7 Average of Validations 93.33333333333334
For K = 9 Average of Validations 93.33333333333334
For K = 11 Average of Validations 93.33333333333334
For K = 13 Average of Validations 94.0
For K = 17 Average of Validations 92.66666666666667
For K = 21 Average of Validations 92.0
For K = 23 Average of Validations 92.66666666666667
For K = 25 Average of Validations 91.33333333333333
For K = 27 Average of Validations 90.66666666666667
For K = 29 Average of Validations 90.66666666666667
For K = 31 Average of Validations 90.0
For K = 33 Average of Validations 89.33333333333333
For K = 35 Average of Validations 87.99999999999999
For K = 37 Average of Validations 88.0
For K = 39 Average of Validations 87.33333333333334
For K = 41 Average of Validations 86.66666666666667
For K = 43 Average of Validations 87.33333333333333
For K = 45 Average of Validations 86.66666666666667
For K = 47 Average of Validations 86.66666666666667
For K = 49 Average of Validations 86.0
For K = 51 Average of Validations 86.0

```

مقدار هر k و همچنین میانگین accuracy برای cross validation هر k محاسبه شده است و طبق این مقادیر بدست آمده مقدار $k = 1$ بهترین حالت برای این دیتاست است.

مقادیر ماتریس درهم ریختگی برای هر داده در فایل پیوست قرار داده شده است.