



دانشکده مهندسی کامپیوتر

به نام خدا



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

پاسخنامه تمرین اول

استاد درس:

دکتر مریم امیر مزلقانی

نیم سال اول ۱۴۰۲-۱۴۰۳

راه ارتباطی:

Aut.DataMining.Fall@gmail.com



سوال ۱.

الف) $H(\text{Target class}) = H\left(\frac{4}{9}, \frac{5}{9}\right) = -\frac{4}{9} \log_2\left(\frac{4}{9}\right) - \frac{5}{9} \log_2\left(\frac{5}{9}\right) = 0.9911$

Target class

ب) $IG(TC, a_1) = H(TC) - H(TC|a_1) = 0.9911 -$

$\frac{4}{9} H\left(\frac{3}{8}, \frac{1}{8}\right) - \frac{5}{9} H\left(\frac{1}{8}, \frac{7}{8}\right) = 0.2294$

$IG(TC, a_2) = H(TC) - H(TC|a_2) = 0.9911 - \frac{5}{9} H\left(\frac{4}{8}, \frac{1}{8}\right) - \frac{4}{9} H\left(\frac{2}{8}, \frac{6}{8}\right) = 0.0072$

ج) از آنجایی که مقدار آستانه در اختیار ما است باید آن را به گونه‌ای انتخاب کنیم که برای ما بهتر باشد. IG را برای آستانه‌های مختلف حساب می‌کنیم.

* همچنین آستانه باید جای باشد.

a_r	۱	۳	۴	۵	۶	۷	۸
TC	+	-	+	-	+	+	-
	۲	۵	۵	۵	۵	۷/۱۰	۷/۱۰

class تغییر کند.

$۲ \rightarrow IG(TC, a_2) = H(TC) - \frac{1}{9} H(1, 0) - \frac{8}{9} H\left(\frac{3}{8}, \frac{5}{8}\right) = 0.1427$

$۳ \rightarrow IG(TC, a_2) = H(TC) - \frac{4}{9} H\left(\frac{1}{8}, \frac{1}{8}\right) - \frac{5}{9} H\left(\frac{5}{8}, \frac{3}{8}\right) = 0.0029$

$۴ \rightarrow IG(TC, a_2) = H(TC) - \frac{3}{9} H\left(\frac{2}{8}, \frac{2}{8}\right) - \frac{6}{9} H\left(\frac{2}{8}, \frac{6}{8}\right) = 0.0072$

$۵ \rightarrow IG(TC, a_2) = H(TC) - \frac{5}{9} H\left(\frac{4}{8}, \frac{1}{8}\right) - \frac{4}{9} H\left(\frac{1}{8}, \frac{7}{8}\right) = 0.0072$

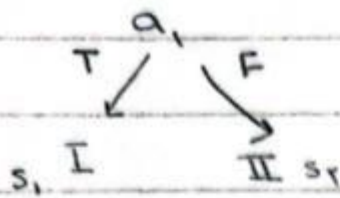
$۶ \rightarrow IG(TC, a_2) = H(TC) - \frac{6}{9} H\left(\frac{3}{8}, \frac{1}{8}\right) - \frac{3}{9} H\left(\frac{1}{8}, \frac{7}{8}\right) = 0.0183$

$۷ \rightarrow IG(TC, a_2) = H(TC) - \frac{1}{9} H\left(\frac{4}{8}, \frac{4}{8}\right) - \frac{1}{9} H(1, 0) = 0.1022$

→ می‌بینیم که مقدار بهتر از a_1 نداریم. آستانه ۲ را انتخاب می‌کنیم.

تعداد آستانه: ۱، ۲، ۳، ۴، ۵، ۶، ۷، ۸

a_r	۲	۳	۴	۵	۶	۷	۸
TC	+	-	+	-	+	+	-



بنا ابتدا a_1 را در بیشترین مقدار می دهیم.

$$I : H(s_1) = -\frac{1}{4} \log_2 \left(\frac{1}{4}\right) - \frac{3}{4} \log_2 \left(\frac{3}{4}\right) = 0.8113$$

$$IG(s_1 | a_2) = H(s_1) - \frac{1}{4} H(1,0) - \frac{3}{4} H\left(\frac{1}{3}, \frac{2}{3}\right) = 0.3113$$

$$IG(s_1 | a_3) = H(s_1) - \frac{1}{4} H(1,0) - \frac{3}{4} H\left(\frac{2}{3}, \frac{1}{3}\right) = 0.1229$$

← a_2 را انتخاب می کنیم

$$II : H(s_2) = -\frac{1}{8} \log_2 \frac{1}{8} - \frac{4}{8} \log_2 \frac{4}{8} = 0.7219$$

$$IG(s_2 | a_2) = H(s_2) - \frac{1}{8} H(1,0) - \frac{4}{8} H\left(\frac{1}{4}, \frac{3}{4}\right) = 0.3219$$

$$IG(s_2 | a_3) = H(s_2) - 1 \times H\left(\frac{1}{8}, \frac{7}{8}\right) = 0$$

← a_3 را انتخاب می کنیم

در صورت I بیان اینکه $IG(s_1 | a_3)$ بیشترین مقدار را بگیرد باید استاندارد ۵،۵ می گرفتیم

$$IG(s_1 | a_3) = H(s_1) - \frac{1}{4} H\left(\frac{1}{4}, \frac{3}{4}\right) - \frac{3}{4} H(1,0) = 0.3113$$

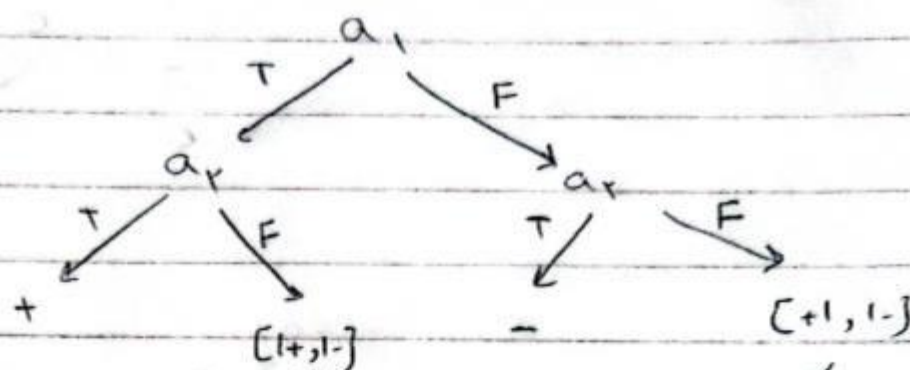
← برابر با مقدار $IG(s_1 | a_2)$ می شود و قدرتی ندارد کدام را انتخاب می کردیم

بیان صفت II نیز باید استاندارد ۴،۵ قرار می دادیم

$$IG(s_2 | a_3) = H(s_2) - \frac{1}{8} H\left(\frac{1}{8}, \frac{7}{8}\right) - \frac{4}{8} H(1,0) = 0.3219$$

← باز هم قدرتی نداشت

← در نهایت استاندارد ۲ در نظر می گیریم و برای s_2 و s_1 را در نظر می گیریم

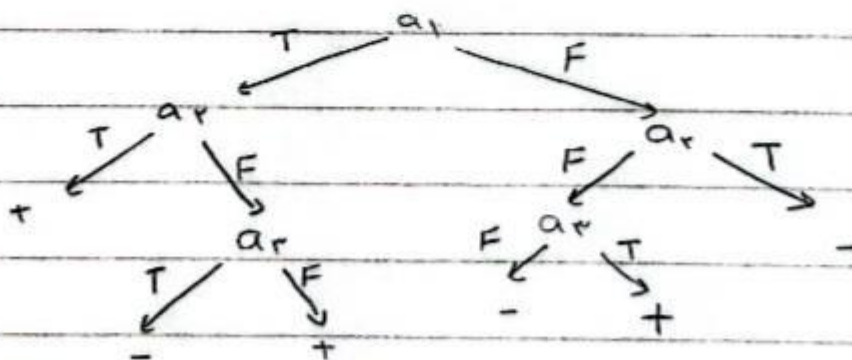


با بررسی این درخت متوجه می شویم اگر آستانه ۰.۵ قرار دهیم بطور کامل دسته بندی خواهر شد:

$$I: IG(S, |a_1) = 0.3113$$

$$II: IG(S, |a_2) = H(S) - \frac{3}{8} H(\frac{1}{4}, \frac{3}{4}) - \frac{2}{8} H(1, 0) = 0.1710$$

a_2	1	2	3	4	5	6	7	8
	+	-	+	-	+	-	+	-



* البته اینجا ممکن است overfitting روی داده آموزش شود



سوال ۲.

(الف)

اگر به جای مربع خطا از قدرمطلق استفاده گردد مسئله بهینه‌سازی به $\min ||X\beta - y||$ تغییر میکند که باید آن را بدست بیاوریم و قدرمطلق در مقایسه با توان ۲ به مقادیر کوچک تفاضل $X\beta - y$ اهمیت بیشتری داده و به مقادیر بزرگ $X\beta - y$ اهمیت کمتری میدهد. لذا در مقابل داده‌های پرت مقاوم‌تر است و همچنین به سمت صفر و تنک کردن مقادیر $X\beta - y$ می‌رویم.

(ب)

در روش gradient descent از گرادیان و مشتق‌گیری استفاده میشود، اما با در نظر گذاشتن قدرمطلق دیگر تابع مشتق‌پذیر نخواهد بود لذا نمیتوان از روش gradient descent استفاده کرد. اما با توجه به محدب بودن تابع میتوان از subgradient استفاده کرد و عیناً روشهای حل قبلی را با subgradient جلو برد.

(ج)

به شکل زیر تشکیل میدهیم X برای هر کلاس داده‌ها یک ماتریس

$$X = \begin{bmatrix} x_{A11} & x_{A21} & x_{A31} & x_{A41} & x_{A51} \\ x_{A12} & x_{A22} & x_{A32} & x_{A42} & x_{A52} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{A1k} & x_{A2k} & x_{A3k} & x_{A4k} & x_{A5k} \end{bmatrix}$$

ماتریس داده‌های کلاس A

ماتریس y را نیز همان داده جدید y در نظر میگیریم. و با حل مسئله بهینه‌سازی \min خطا را بدست میآوریم. این کار را برای هر کلاس تکرار میکنیم. (با X متناسب با همان کلاس) سپس هرکدام که داری خطای مینیمم کمتری باشد کلاس داده جدید خواهد بود.

(البته میتوان از رگرسیون شماره کلاس بر حسب داده‌ها هم استفاده کرد که کمتر کلی است)



سوال ۳.

$$x = \begin{bmatrix} 1 & 42 \\ 1 & 74 \\ 1 & 48 \\ 1 & 35 \\ 1 & 56 \\ 1 & 26 \\ 1 & 60 \end{bmatrix}$$

$$y = \begin{bmatrix} 98 \\ 130 \\ 120 \\ 88 \\ 182 \\ 80 \\ 135 \end{bmatrix}$$

$$x^T x \beta = x^T y$$

$$x^T x = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 42 & 74 & 48 & 35 & 56 & 26 & 60 \end{bmatrix} \times \begin{bmatrix} 1 & 42 \\ 1 & 74 \\ 1 & 48 \\ 1 & 35 \\ 1 & 56 \\ 1 & 26 \\ 1 & 60 \end{bmatrix} = \begin{bmatrix} 7 & 341 \\ 341 & 18181 \end{bmatrix}$$

$$x^T y = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 42 & 74 & 48 & 35 & 56 & 26 & 60 \end{bmatrix} \times \begin{bmatrix} 98 \\ 130 \\ 120 \\ 88 \\ 182 \\ 80 \\ 135 \end{bmatrix} = \begin{bmatrix} 833 \\ 42948 \end{bmatrix}$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} 7 & 341 \\ 341 & 18181 \end{bmatrix}^{-1} \times \begin{bmatrix} 833 \\ 42948 \end{bmatrix} = \begin{bmatrix} 44.34 \\ 1.532 \end{bmatrix}$$

$$y = 1.532x + 44.34$$

تخمین فشار خون برای سن ۴۰:

$$y = 1.532 \times 40 + 44.34 = 105.62$$



سوال ۴.

۱.

دقت آموزش بسیار بالا (۱۰۰٪) است در حالی که دقت آزمون پایینتر است (۵۰٪). این نشان دهنده احتمال زیاد **overfitting** است. در واقع، مدل به خوبی به داده‌های آموزشی هماهنگ شده است (دقت بالا)، اما نمی‌تواند با داده‌های جدید (آزمون) به خوبی عمل کند، که نشان دهنده **overfitting** است. بنابراین، در احتمال **overfitting** بیشتر است، زیرا دقت آموزش بسیار بالا و دقت آزمون بسیار پایین است، که نشان دهنده عدم تعمیم‌پذیری خوب مدل به داده‌های جدید است.

۲.

دقت آموزش و دقت آزمون به نسبت هم نزدیک‌تر هستند (۸۰٪ در مقابل ۷۰٪)، که نشان دهنده یک تطابق بهتر بین عملکرد مدل در داده‌های آموزش و آزمون است. این حالت ممکن است نشانگر یک مدل بهتر و کمتر **overfit** شده باشد.

سوال ۵.

منظم سازی برای جلوگیری از برازش بیش از حد داده‌ها، به ویژه زمانی که اختلاف زیادی بین عملکرد مجموعه آموزش و مجموعه تست وجود دارد. با منظم سازی، تعداد ویژگی‌های مورد استفاده در تمرین ثابت نگه داشته می‌شود، اما مقدار ضرایب (w) کاهش می‌یابد.

Lasso Regression

این یک تکنیک منظم‌سازی است که در انتخاب ویژگی با استفاده از روش انقباض استفاده می‌شود که به آن روش رگرسیون جریمه شده نیز گفته می‌شود. Lasso. مخفف Least Absolute Shrinkage and Selection Operator است که هم برای منظم سازی و هم برای انتخاب مدل استفاده می‌شود. این روش عبارت جریمه را به تابع هزینه اضافه می‌کند. که این عبارت جریمه مجموع مطلق ضرایب است که باعث کاهش مقدار ضرایب به منظور کاهش ضرر می‌شود. این روش تمایل دارد که ضرایب را به صفر مطلق میل دهد.

$$L_{lasso} = \operatorname{argmin}_{\hat{\beta}} \left(\|Y - \beta * X\|^2 + \lambda * \|\beta\|_1 \right)$$



Ridge Regression

تفاوت این روش با روش قبلی در این است که مقدار جریمه‌های که به تابع هزینه اضافه می‌شود برابر با مجذور ضرایب است. بر خلاف Lasso این روش هیچگاه ضرایب را به سمت صفر مطلق سوق نمی‌دهد.

$$L_{ridge} = \operatorname{argmin}_{\hat{\beta}} \left(\|Y - \beta * X\|^2 + \lambda * \|\beta\|_2^2 \right)$$

سوال ۶.

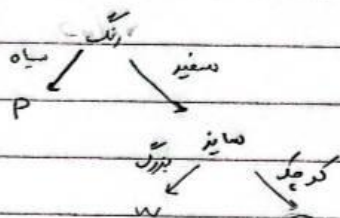
Instance	رنگ	سایز	class
۱	سفید	کوچک	w
۲	سفید	کوچک	w
۳	سفید	کوچک	w
۴	سفید	بزرگ	w
۵	سفید	کوچک	p
۶	سیاه	کوچک	p
۷	سیاه	بزرگ	p
۸	سیاه	بزرگ	p

$$H(class) = 1$$

$$IG(class | رنگ) = H(class) - \frac{3}{8} H\left(\frac{4}{8}, \frac{4}{8}\right) - \frac{5}{8} H\left(\frac{1}{8}, \frac{7}{8}\right) = 0.5211$$

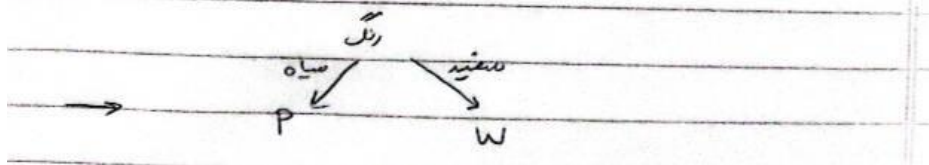
$$IG(class | سایز) = H(class) - \frac{1}{8} H\left(\frac{1}{8}, \frac{7}{8}\right) - \frac{3}{8} H\left(\frac{1}{8}, \frac{7}{8}\right) = 0.1887$$

← ابتدا رنگ را انتخاب می‌کنیم



در اینجا ۳ داده در کلاس w قرار می‌گیرند و ۱ داده در کلاس p قرار می‌گیرد

← plastic در نهایت این قسمت را w و p در نظر می‌گیریم





سوال ۷.

(الف)

برای هر کدام از ویژگی‌های درون جدول مقدار Gini index را محاسبه میکنیم.

Gini Index for Heavy :

$$Gini(Heavy=NO) = 1 - \left(\left(\frac{12}{18} \right)^2 + \left(\frac{6}{18} \right)^2 \right) = \frac{12}{18}$$

$$Gini(Heavy=Yes) = 1 - \left(\left(\frac{1}{9} \right)^2 + \left(\frac{2}{9} \right)^2 \right) = \frac{8}{9}$$

$$Gini(Heavy) = \frac{6}{18} \times \frac{12}{18} + \frac{12}{18} \times \frac{8}{9} = \frac{3}{10} + \frac{1}{4} = 0.475$$

Gini Index for Spotted :

$$Gini(Spotted=NO) = 1 - \left(\left(\frac{6}{10} \right)^2 + \left(\frac{4}{10} \right)^2 \right) = \frac{12}{25}$$

$$Gini(Spotted=Yes) = 1 - \left(\left(\frac{1}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right) = \frac{8}{25}$$

$$Gini(Spotted) = \frac{6}{10} \times \frac{12}{25} + \frac{4}{10} \times \frac{8}{25} = 0.28$$

Gini Index for Smooth :

$$Gini(Smooth=NO) = 1 - \left(\left(\frac{4}{12} \right)^2 + \left(\frac{2}{12} \right)^2 \right) = \frac{1}{3}$$

$$Gini(Smooth=Yes) = 1 - \left(\left(\frac{1}{8} \right)^2 + \left(\frac{3}{8} \right)^2 \right) = \frac{3}{8}$$

$$Gini(Smooth) = \frac{1}{3} \times \frac{1}{3} + \frac{1}{3} \times \frac{3}{8} = \frac{1}{2} + \frac{1}{12} = \frac{7}{12} = 0.583$$

لذا $\boxed{\text{Smooth}}$ انتخاب گردد چرا که Gini آن کمترین است.



(ب)

$$I(\text{parent}) = -\frac{3}{8} \log \frac{3}{8} - \frac{5}{8} \log \frac{5}{8} = .5625 + .4375 = .9999$$

$$I(\text{smooth} = \text{NO}) = -\frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} = 1$$

$$I(\text{smooth} = \text{Yes}) = -\frac{3}{8} \log \frac{3}{8} - \frac{1}{8} \log \frac{1}{8} = .5625 + .125 = .6875$$

~~Gain~~

$$\text{Gain}(\text{smooth}) = .9999 - \left(\frac{1}{4} \times 1 + \frac{1}{4} \times .6875 \right) = .9999 - .7656 = .2343$$

data smooth = NO

Heavy	spotted	Poisson
NO	NO	NO
NO	Yes	NO
NO	Yes	Yes
Yes	NO	Yes

Gain for Heavy:

$$\text{Gain}(\text{Heavy} = \text{NO}) = 1 - \left(\left(\frac{1}{4} \right)^2 + \left(\frac{3}{4} \right)^2 \right) = \frac{3}{4}$$

$$\text{Gain}(\text{Heavy} = \text{Yes}) = 0$$

$$\text{Gain}(\text{Heavy}) = \frac{3}{4} \times \frac{3}{4} + \frac{1}{4} \times 0 = \frac{9}{16} = .5625$$

Gain for spotted:

$$\text{Gain}(\text{spotted} = \text{NO}) = 1 - \left(\left(\frac{1}{4} \right)^2 + \left(\frac{1}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right) = \frac{5}{8}$$

$$\text{Gain}(\text{spotted} = \text{Yes}) = 1 - \left(\left(\frac{1}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right) = \frac{3}{4}$$

$$\text{Gain}(\text{spotted}) = \frac{1}{4} \times \frac{5}{8} + \frac{1}{4} \times \frac{3}{4} = .75$$

Heavy انتخاب

data smooth = Yes

Heavy	spotted	Poisson
Yes	NO	NO
Yes	NO	Yes
NO	Yes	Yes
NO	NO	Yes

$$\text{Gain}(\text{Heavy} = \text{NO}) = 1 - 1 = 0$$

$$\text{Gain}(\text{Heavy} = \text{Yes}) = 1 - \left(\left(\frac{1}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right) = \frac{1}{2}$$

$$\text{Gain}(\text{Heavy}) = \frac{1}{4} \times \frac{1}{2} + \frac{1}{4} \times 0 = .125$$

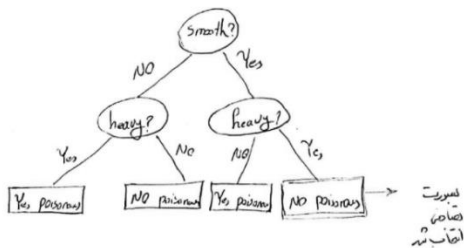
$$\text{Gain}(\text{spotted} = \text{NO}) = 1 - \left(\left(\frac{1}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right) = \frac{3}{4}$$

$$\text{Gain}(\text{spotted} = \text{Yes}) = 0$$

$$\text{Gain}(\text{spotted}) = \frac{3}{4} \times \frac{3}{4} + \frac{1}{4} \times 0 = .75$$

Heavy انتخاب

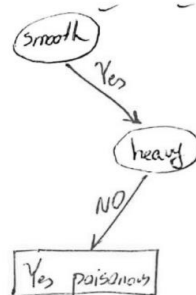
لذا داده برای تصمیم‌گیری:





(د)

داده‌های نمونه‌ای (۱) عبارت است از $smooth = Yes, spotted = No, Heavy = No$
 مثال در درخت تصمیم زیر می‌بینیم:



با توجه به درخت تصمیم این نوع ماکارون سمی خواهد بود.