

به نام خدا
تمرین سری دوم

مریم رضوانی
شماره دانشجویی: ۹۹۲۱۱۶۰۰۱۹

(1) به سؤالات زیر پاسخ دهید.

الف) $f1\text{-score}$ و Accuracy یک معیار ارزیابی برای مدل‌های categorical است. برای هر دو معیار، هر چه مقدار بالاتر باشد، یک مدل بهتر می‌تواند مشاهدات را به کلاس‌ها طبقه‌بندی کند. $f1\text{-score}$ به این صورت عمل می‌کند که میانگین هارمونیک دقت ($precision$) و ($recall$) است. دقت نسبت تعداد پیش‌بینی‌های مثبت صحیح به تعداد کل پیش‌بینی‌های مثبت است.

$$f1\text{-score} = 2 * recall * precision / recall + precision$$

$recall$ نسبت تعداد پیش‌بینی‌های مثبت صحیح به تعداد کل مثبت‌های واقعی است. این معیار نشان می‌دهد که چه درصد از داده‌های مثبت را مدل به درستی شناسایی کرده است. این معیار زمانی مناسب است که داده‌ها توزیع نامتوازنی از کلاس‌ها داشته باشند و هزینه‌ی خطای منفی و مثبت متفاوت باشد.

اما $accuracy$ نسبت تعداد پیش‌بینی‌های درست به تعداد کل پیش‌بینی‌ها است.

$$Accuracy = \frac{TN+TP}{TN+TP+FN+FP}$$

این معیار نشان می‌دهد که چه درصد از داده‌ها را مدل به درستی دسته‌بندی کرده است. این معیار زمانی مناسب است که داده‌ها توزیع متوازنی از کلاس‌ها داشته باشند و هزینه‌ی خطای منفی و مثبت یکسان باشد.

به عنوان یک قاعده کلی: ما اغلب زمانی از Accuracy استفاده می‌کنیم که کلاس‌ها متعادل هستند و هیچ نقطه ضعف عمده‌ای برای پیش‌بینی منفی‌های کاذب وجود ندارد. ما اغلب از زمانی $f1\text{-score}$ استفاده می‌کنیم که کلاس‌ها نامتعادل هستند و پیش‌بینی منفی‌های کاذب جنبه منفی جدی دارد.

به عنوان مثال، اگر از یک مدل رگرسیون لجستیک برای پیش‌بینی اینکه آیا فردی سرطان دارد یا نه استفاده کنیم، منفی‌های کاذب واقعاً بد هستند (مثلاً پیش‌بینی اینکه کسی سرطان ندارد در حالی که واقعاً سرطان دارد) بنابراین $f1\text{-score}$ مدل‌هایی را که دارای منفی کاذب بیش از حد هستند جریمه می‌کند. بیش از Accuracy.

ب) بین sensitivity و specificity در یک confusion matrix این است که:

- sensitivity: تقسیم پیش‌بینی‌های درست رده مثبت (TP) بر تعداد کل نمونه‌های مثبت موجود در مجموعه داده‌ها (TP+FN) به دست می‌آید. برای مثال برای تشخیص بیماران قلبی در یک مدل؛ نسبت بیماران قلبی که درست تشخیص داده شده‌اند (TP) به مجموع تعداد بیماران که درست تشخیص داده شده‌اند (TP) و بیماران که درست تشخیص داده نشده‌اند (FN).

همچنین به عنوان نرخ مثبت واقعی (TPR) شناخته می‌شوند.

- Specificity: پیش‌بینی درست تعداد منفی‌های رده منفی (TN) بر تعداد کل منفی‌های واقعی (TN+FP) است. برای نمونه در آزمایش افراد بیمارانی قلبی؛ تعداد افرادی که مدل توانسته سالم بودنشان (بیماری ندارند) را تشخیص دهد (TN). به نسبت مجموع افرادی که درست تشخیص دادیم (بیمار نیستند) سالم اند (TN) و افرادی که مدل نتوانسته سالم بودنشان را تشخیص دهد (به اشتباه بیمار تشخیص داده) است (FP).

همچنین به عنوان نرخ منفی واقعی (TNR) شناخته می‌شوند.

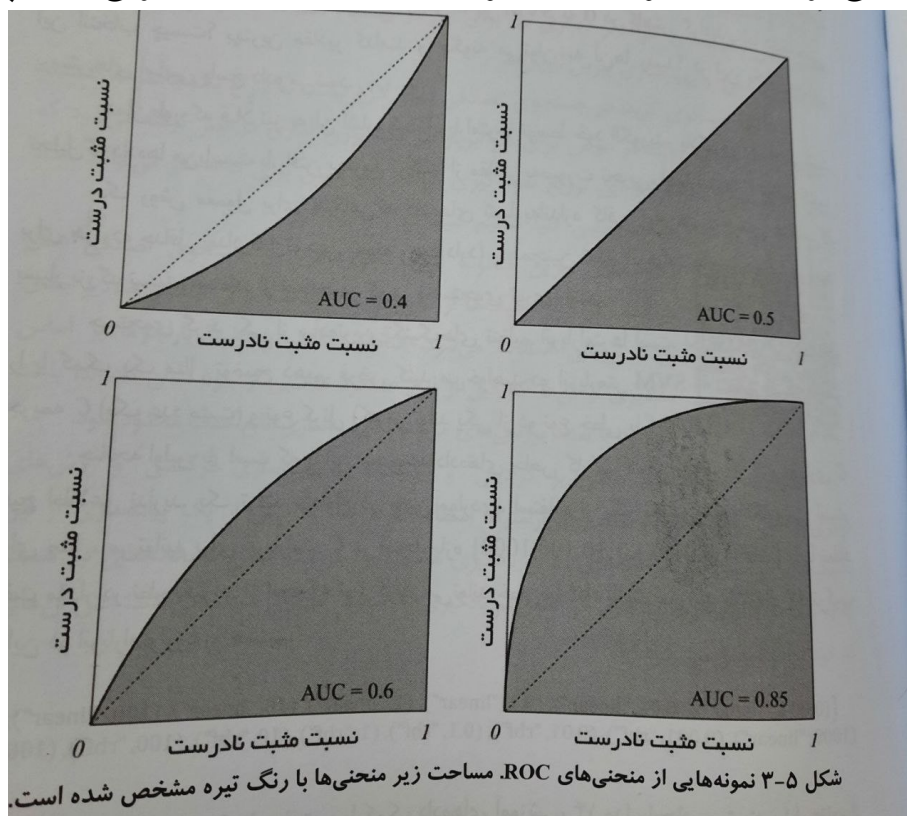
- نرخ مثبت کاذب (FPR): یعنی چه تعداد از افراد سالم توسط آزمایش به اشتباه بیمار شناسایی شده‌اند. $(1 - \text{Specificity})$
- نرخ منفی کاذب (FNR): یعنی چه تعداد از افراد بیمار توسط آزمایش به اشتباه سالم شناسایی شده‌اند. $(1 - \text{Sensitivity})$

ج) برای ارزیابی یک مدل طبقه‌بندی، می‌توانیم از منحنی مشخصه عملکرد گیرنده (ROC) و سطح زیر منحنی (AUC) استفاده کنیم. این دو معیار به ما نشان می‌دهند که چگونه مدل ما در تشخیص کلاس‌های مثبت و منفی عمل می‌کند. منحنی‌های ROC ترکیبی از نسبت مثبت درست (TPR) که برابر با Recall است، و نسبت مثبت نادرست (نسبتی از نمونه‌های منفی که به درستی پیش‌بینی نشده‌اند) را برای خلق تصویری خلاصه از عملکرد رده‌بندی استفاده می‌کنند.

نسبت مثبت درست (TPR) برابر است با: $TPR = TP / (TP + FN)$ و نسبت مثبت نادرست (FPR) برابر است با: $FPR = FP / (FP + TN)$ منحنی‌های ROC تنها برای طبقه‌بندی‌هایی استفاده می‌شوند که همراه با پیش‌بینی خود امتیازاتی (با یک احتمال) را برمی‌گردانند. برای مثال طبقه‌بندی‌هایی نظیر درخت تصمیم و شبکه‌های عصبی (و مدل‌های تلفیقی از درختان تصمیم) را می‌توان با استفاده از منحنی‌های ROC ارزیابی کرد.

برای رسم یک منحنی ROC ابتدا باید بازه خروجی مدل گسسته سازی شود. چنانچه بازه مزبور برای یک مدل برابر با $[0,1]$ باشد آن‌گاه می‌توان این بازه را بصورت $[0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9]$ گسسته سازی کرد. سپس از هر مقدار گسسته به عنوان آستانه پیش‌بینی استفاده می‌شود و برچسب‌های نمونه‌های موجود در داده‌ها با استفاده از مدل و این آستانه پیش‌بینی می‌شوند. برای مثال اگر می‌خواهید TPR و FPR را برای حد آستانه ای برابر با 0.7 محاسبه کنید؛ مدل را روی هر نمونه اعمال کرده و خروجی مدل را دریافت کنید. چنانچه مقدار خروجی مدل عددی بزرگ‌تر یا مساوی 0.7 بود شما رده مثبت را پیش‌بینی می‌کنید و در غیر اینصورت پیش‌بینی شما رده منفی است. در شکل زیر به سادگی می‌توان دید که با تنظیم آستانه با مقدار صفر؛ تمام پیش‌بینی‌ها مثبت خواهند بود. و بنابراین مقادیر TPR و FPR نیز برابر با یک خواهد بود. (شکل بالا سمت راست). از طرف دیگر چنانچه مقدار آستانه برابر با یک باشد هیچ یک از پیش‌بینی‌ها مثبت نخواهد بود و مقادیر TPR و FPR نیز برابر با صفر خواهد بود. (شکل پایین سمت چپ). رده بندی که

مساحت زیر منحنی ROC آن بیشتر باشد رده بندی بهتر است. این مساحت با AUC در شکل نشان داده شده است. یک رده بند با مقدار AUC بزرگتر از 0.5 بهتر از یک رده بند تصادفی است. چنانچه مقدار AUC کمتر از 0.5 باشد موردی در مدل شما اشتباه است. مقدار AUC یک رده بند کامل برابر با یک است. معمولاً اگر مدل شما رفتاری مناسب داشته باشد با انتخاب مقداری برای آستانه که باعث نزدیک شدن TPR به یک می شود و FPR را نزدیک به صفر نگه میدارد. به رده بندی خوبی دست پیدا می کنید.



د) یک مجموعه اعتبار سنجی جداگانه زمانی مفید است که شما بخواهید پارامترهای مدل خود را بهینه کنید و عملکرد مدل را بر روی داده های جدید ارزیابی کنید. مجموعه اعتبار سنجی به شما این امکان را می دهد که مدل خود را بر روی داده هایی که در آموزش دیده نشده اند، امتحان کنید و پارامترهایی را انتخاب کنید که بهترین نتیجه را بدهند. مجموعه آزمایشی فقط برای تخمین خطای نهایی مدل بر روی داده های جدید استفاده می شود و نباید در فرآیند بهینه سازی مدل تاثیر گذار باشد. اگر شما فقط از دو مجموعه آزمایشی و آموزشی استفاده کنید، ممکن است خطای مدل را در آزمایش بالاتر یا پایین تر از حقیقت برآورد کنید و عملکرد مدل را به درستی نشان ندهید.

$$E_{out} = E(h(x) - y)^2$$

سوال 2

$$E_{out} = \int (h(x) - y)^2 p(x, y) dx dy$$

تابع مایل احتمال مشترک x, y

$$\rightarrow E_{out} = \int (h(x) - E[y|x] + E[y|x] - y)^2 p(x, y) dx dy$$

با استفاده از ماکتور مربعی

$$E_{out} = \int (h(x) - E[y|x])^2 p(x, y) dx dy + \int (E[y|x] - y)^2 p(x, y) dx dy$$

$$+ 2 \int (h(x) - E[y|x]) (E[y|x] - y) p(x, y) dx dy$$

برای سه جمله می توانیم برابر با صفر است: زیرا

$$p(x, y) dx dy = \int (h(x) - E[y|x]) (E[y|x] - y) p(x, y) dx dy =$$

$$= \int (h(x) - E[y|x]) p(y|x) p(x) dy x = \int (h(x) - E[y|x]) p(y|x) dx$$

$$\int p(x) dx = 0$$

بنابراین تمامی خارج از قوس را می توان به صفر نزدیک کرد

$$E_{out} = \int (h(x) - E[y|x])^2 p(x, y) dx dy + \int (E[y|x] - y)^2 p(x, y) dx dy$$

این کمیت که در این عبارت داریم جمله اول را کمیت کنیم این جمله همواره مثبت است و زمانی که صفر شود

خطای خارج از قوس کمیت نه شود. جمله دوم نمایانگر خطای ناگزیر از تغییرات تصادفی y است و به دلیل

تکرار ندارد.



Page: ()

SUBJECT: Year: Month: Day:

فصلت بود: برای اثبات $E(x)$ مقدار صفر دارد. داریم:

$$h^*(x) = E[y|x]$$

$$E[h^*(x)|x] = E[E[y|x]|x] = E[y|x] = h^*(x)$$

که در آن $E(x)$ نشان دهنده مقدار مورد انتظار است. داریم:

$$y = h^*(x) + \varepsilon(x)$$

$$E[y|x] = E[h^*(x) + \varepsilon(x)|x] = E[h^*(x)|x] + E[\varepsilon(x)|x]$$

$$= h^*(x) + E[\varepsilon(x)|x] = h^*(x) + E[\varepsilon(x)|x]$$

ماضوف $h^*(x)$ از دو طرف معادله می‌توانیم بنویسیم:

$$E[\varepsilon(x)|x] = 0 = E[\varepsilon(x)|x] = 0$$

یعنی معادله $E(x)$ برابر با صفر است.

(3)

Page: ()

SUBJECT: Year: Month: Day:

مقدار هدف پیش‌بینی \hat{y}_i (مقدار واقعی y_i)

سوال (3)

در میزبان برای غایب مقدار (MSE) می‌توانیم محاسبه کنیم:

آن را برای این سوال محاسبه کنید:

(Y_i)	(\hat{Y}_i)
$(Y_1) = 10$	$(\hat{Y}_1) = 12$
$(Y_2) = 20$	$(\hat{Y}_2) = 18$
$(Y_3) = 15$	$(\hat{Y}_3) = 10$
$(Y_4) = 25$	$(\hat{Y}_4) = 28$
$(Y_5) = 30$	$(\hat{Y}_5) = 32$

(Mean Squared Error)

می‌توانیم مجموع مربعات خطا

$$MSE = \frac{1}{2} \sum_{i=1}^N (f(x_i) - y_i)^2$$

$$MSE = \frac{1}{2} \sum_{i=1}^5 ((\hat{Y}_i) - Y_i)^2 =$$

$$(\hat{Y}_1 - Y_1)^2 = (12 - 10)^2 = 4, (\hat{Y}_2 - Y_2)^2 = (18 - 20)^2 = 4$$

$$(\hat{Y}_3 - Y_3)^2 = (10 - 15)^2 = 25, (\hat{Y}_4 - Y_4)^2 = (28 - 25)^2 = 9$$

$$(\hat{Y}_5 - Y_5)^2 = (32 - 30)^2 = 4$$

$$MSE = \frac{1}{2} \sum_{i=1}^5 (4 + 4 + 25 + 9 + 4) = \frac{46}{2} = 23$$

Senobar

4) مدل ۲ بهتر است؛ زیرا حساسیت (sensitivity) بالاتری نسبت به بقیه مدل ها برخوردار است. بدلیل اینکه حساسیت نشان می دهد چند درصد از تومورهای بدخیم را مدل به درستی شناسایی کرده است و این برای ما خیلی مهم است زیرا هزینه ی FN (false negative) خیلی بالا است. به این معنا که فرض کنید شخصی تومور بدخیم دارد ولی مدل ما آن را سالم پیش بینی کرده است. در اینجور مسائل بدنبال دقت یا حساسیت (sensitivity) بالا هستیم.

$$Accuracy = \frac{TN+TP}{TN+TP+FN+FP}$$

$$sensitivity = \frac{TP}{TP+FN}$$

$$specificity = \frac{TN}{TN+FP}$$

با قراردادن مقادیر بالا برای هریک از مدل های زیر بدست می آوریم:
کل داده ها را ۱۰۰۰ نفر در نظر می گیریم.

مدل 1:

actual		
Predicted	TP=20	FP=0
	FN=80	TN=900

مدل 2:

actual		
Predicted	TP=70	FP=110
	FN=30	TN=790

مدل 3:

actual		
Predicted	TP=40	FP=27
	FN=60	TN=873

5) Confusion Matrix:

برای class 1 داریم:

Predicted	actual			
		class1	class2	class3
	class1	TP=7	FP=0	FP=4
	class2	FN=0	TN=6	TN=1
	class3	FN=1	TN=0	TN=1

Class1: TP=7, FP=0+4=4, FN=0+1=1, TN=6+1+1+0=8

sensitivity=TP/TP+FN=7/(7+1)=7/8
FNR=1-Sensitivity=1-7/8=1/8=12.5%

Accuracy = $TN+TP / TN+TP+FN+FP = (7+8)/(7+8+1+4) = 15/20$

Error Rate=1-Accuracy=1-15/20=1/4=0/25=25%

Precision=tp/tp+fp=7/(7+4)=7/11

recall=tp/tp+fn=7/8

F1-Score= 2 * recall * precision / recall + precision

F1-Score=(2*7/8*7/11) / (7/8+7/11)=0/73=73%

برای کلاس ۲ و ۳ هم به همین صورت داریم:

Class2: TP=6, FP=0+1=1, FN=0+0=0, TN=4+7+1+1=13

Class3: TP=1, FP=0+1=1, FN=4+1=5, TN=6+7+0+0=13

سوال (7) مجموعه داده های آزمون (در درج های A, B, C, و فصول: کلاس) همراه با نتایج آزمون

دسته بندی (فصول یعنی 0) داده شده است. دو نقطه روی منحنی ROC برای سایر آستانه 0.5 در

18	پیدا کرده و رسم کنید.	predicted	class	C	B	A
19	0.97	1	1	A	2	10
20	0.61	0	0	B	1	20
21	0.77	1	1	A	3	30
22	0.91	1	1	B	2	40
23	0.12	0	0	B	1	15

آستانه 0.5 است.
 FN: به معنای تعداد داده های که فصول 1 گرفته اند است و فصول 0 که بیشتر از آستانه است.
 TN: به معنای تعداد داده های که فصول 0 گرفته اند است و فصول 1 که بیشتر از آستانه است

$$FPR = \frac{FP}{FP + TN} \quad TPR = \frac{TP}{TP + FN}$$

برای آستانه 0.5 داریم:
 $TP = 3, FN = 0, FP = 1$

$$TN = 1 \quad TPR = \frac{3}{3 + 0} = 1 \quad FRR = \frac{1}{1 + 1} = 0.5$$

برای آستانه 0.8 داریم:
 $TP = 2, FN = 1, FP = 0, TN = 2$

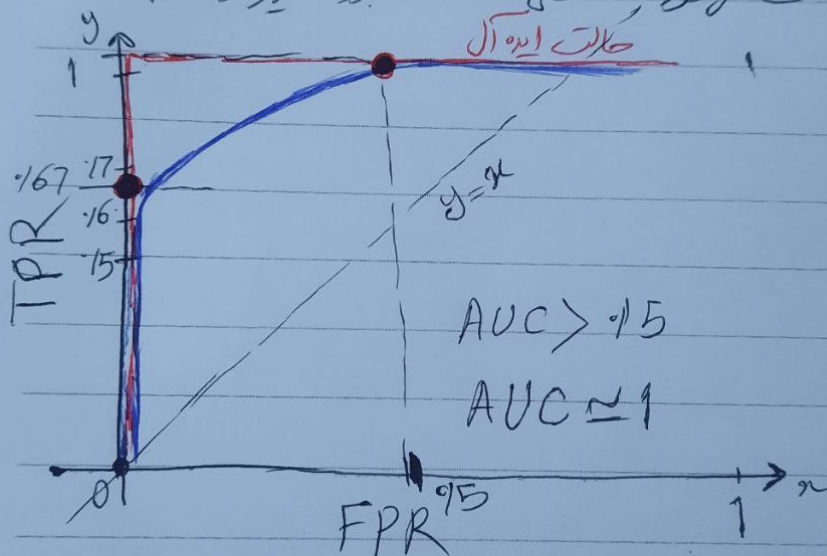
$$TPR = \frac{2}{2 + 1} = 0.67 \quad FPR = \frac{0}{0 + 2} = 0$$

حالا با رسم نقطه TPR و FPR برای آستانه 0.5 و 0.8 روی منحنی ROC را بدست می آوریم

Subject:

Year : Month : Day : ()

معنی ROC شامل نقاط (0,0) که به معنای همان سالم است و نقطه (1,1) که به معنای همان بیمار است می شود. معنی ROC نسبت زیر قوسه بود:



(یا انتخاب مقدار برای آن که باعث ترسیدن TPR به یک می شود و FPR را نزدیک به صفر نگه می دارد به روشی خوب دست پیدا می کنید)

(8

(الف

$$\text{Entropy}([3+, 7-]) = -3/10 \log_2 (3/10) + (7/10 \log_2 7/10) = 0.87$$

(ب

$$\text{Gain}(S, A) \equiv \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

values(Tid)=(1,2,3,4,5,6,7,8,9,10)

S=([3+, 7-])

27 Monday
2023 February
۰۶ شعبان ۱۴۴۴

ظریف مصور
Zarif Mosavar

دوشنبه
1401
اسف

(ب

$S_1, S_2, S_3, \dots, S_{10} = [1+, 9-]$

$\sum_{v \in \text{Values}(Tid)} \frac{|S_v|}{|S|} \text{Entropy}(S_v) =$

$-\frac{1}{10} \text{Entropy}(S_1) - \frac{1}{10} \text{Entropy}(S_2) - \dots$

$-\frac{1}{10} \text{Entropy}(S_{10}) = \boxed{0.467}$

$\text{Gain}(S, Tid) = \text{Entropy}(S) - \sum \frac{|S_v|}{|S|} \text{Entropy}(S_v)$

$= 0.87 - 0.467 = 0.403$

Values (Tid) = 1, 2, 3, 4, 5, 6, 7, 8, 9

$$S = [3+, 7-]$$

$$S_1 \leftarrow [1+, 0], S_2 \leftarrow [1-, 0], S_3 \leftarrow [1-, 0], S_4 \leftarrow [1-, 0], S_5 \leftarrow [1+, 0], S_6 \leftarrow [1-, 0]$$

$$S_7 \leftarrow [1-, 0], S_8 \leftarrow [0, 1+], S_9 \leftarrow [1-, 0], S_{10} \leftarrow [0, 1+]$$

$$\text{Entropy}(S_v) = \frac{1}{10} \log_2 \frac{1}{10} = 0.332$$

$$\sum_{n=1}^n \frac{|S_v|}{|S|} \text{Entropy}(S_v) = 10 \times 0.332 = 0.332$$

$$0.187 - \sum \frac{|S_v|}{|S|} \text{Entropy}(S_v) = 0.187 - 0.332 = 0.538$$

25

Values (Marital Status) = Single, Married, Divorced

$$S = [3+, 7-]$$

$$S_{\text{Single}} \leftarrow [2-, 2+] , S_{\text{married}} \leftarrow [4-, 0]$$

$$S_{\text{Divorced}} \leftarrow [1-, 1+]$$

$$\text{Gain}(S, \text{Marital Status}) = \text{Entropy}(S) - \sum \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= 0.187 - \frac{4}{10} \text{Entropy}(S_{\text{Single}}) - \frac{4}{10} \text{Entropy}(S_{\text{married}})$$

$$- \frac{2}{14} \text{Entropy}(S_{\text{Divorced}}) =$$

روز حمایت از بیماران نادر

Values (Marital Status) = Single, Married, Divorced

$$S = [3+, 7-]$$

$$S_{\text{Single}} \leftarrow [2-, 2+] \quad , \quad S_{\text{married}} \leftarrow [4-, 0]$$

$$S_{\text{Divorced}} \leftarrow [1-, 1+]$$

$$\begin{aligned} \text{Gain}(S, \text{Marital Status}) &= \text{Entropy}(S) - \sum \frac{|S_v|}{|S|} \text{Entropy}(S_v) \\ &= 1.87 - \frac{4}{10} \text{Entropy}(S_{\text{Single}}) - \frac{4}{10} \text{Entropy}(S_{\text{married}}) \\ &\quad - \frac{2}{14} \text{Entropy}(S_{\text{Divorced}}) = \end{aligned}$$

روز حمایت از بیماران نادر

اسفند

$$\text{Entropy}(S_{\text{Single}}) = -\frac{2}{10} \log_2 \frac{2}{10} - \frac{2}{10} \log_2 \frac{2}{10}$$
$$= 1.46 + 1.46 = 2.92$$

$$\text{Entropy}(S_{\text{married}}) = -\frac{4}{10} \log_2 \frac{4}{10} + 0 = 1.528$$

$$\text{Entropy}(S_{\text{Divorced}}) = -\frac{1}{10} \log_2 \frac{1}{10} + \frac{1}{10} \log_2 \frac{1}{10} = 1.664$$

$$\begin{aligned} \text{Gain}(S, \text{Marital Status}) &= 1.87 - \left(\frac{4}{10}\right) 2.92 \\ &\quad - \left(\frac{4}{10}\right) 1.528 - \left(\frac{2}{10}\right) 1.664 = 0.15 \end{aligned}$$

0.15

هرچه داده‌ها یکنواخت تر باید به دسته‌های بیشتر تقسیم کنیم. در اینجا ویژگی آیدی ۱۰ فرزند دارد که در هر فرزند ۱ نمونه قرار می‌گیرد. به همین دلیل برای شرط تست ویژگی مناسب نمی‌باشد.

طرح رسم درخت تصمیم: ویژگی را به عنوان گره ریشه انتخاب می‌کنیم که داده را به ۲ زیرمجموعه مجزا تقسیم کند.
به این منظور از معیار محاسبات (Information Gain) استفاده می‌کنیم و ویژگی را به دو دسته تقسیم می‌کنیم که کمترین آنتروپی را در هر دسته ایجاد کند.

Values (Refund) = Yes, No

$S = [3+, 7-]$ $S_{Yes} = [3-, 0]$ $S_{No} = [3+, 4-]$

$Gain(S, \text{Refund}) = Entropy(S) - \sum \frac{|S_v|}{|S|} Entropy(S_v) = .187 - \frac{3}{10} (.15) - \frac{7}{10} (.4) = 0$

$Entropy(S_{Yes}) = -\frac{3}{10} \log_2 \frac{3}{10} = .51$

$Entropy(S_{No}) = -\frac{3}{10} \log_2 \frac{3}{10} - \frac{4}{10} \log_2 \frac{4}{10} = .528$

PAYCO

Values (Tid) = 1, 2, 3, 4, 5, 6, 7, 8, 9

$S = [3+, 7-]$

$S_1 = [1+, 0]$, $S_2 = [1-, 0]$, $S_3 = [1-, 0]$, $S_4 = [1-, 0]$, $S_5 = [1+, 0]$, $S_6 = [1-, 0]$

$S_7 = [1-, 0]$, $S_8 = [0, 1+]$, $S_9 = [1-, 0]$, $S_{10} = [0, 1+]$

$Entropy(S_v) = \frac{1}{10} \log_2 \frac{1}{10} = .332$

$\sum_{n=1}^{10} \frac{|S_v|}{|S|} Entropy(S_v) = 10 \times .332 = .332$

$.187 - \sum \frac{|S_v|}{|S|} Entropy(S_v) = .187 - .332 = .1538$

اسم فرد

$Entropy(S_{single}) = -\frac{2}{10} \log_2 \frac{2}{10} - \frac{2}{10} \log_2 \frac{2}{10} = .464$

$.464 + .464 = .928$

$Entropy(S_{married}) = -\frac{4}{10} \log_2 \frac{4}{10} = .528$

$Entropy(S_{divorced}) = -\frac{1}{10} \log_2 \frac{1}{10} + \frac{1}{10} \log_2 \frac{1}{10} = .664$

$Gain(S, \text{Marital Status}) = .187 - (\frac{4}{10}) .928 - (\frac{4}{10}) .528 - (\frac{2}{10}) .664 = .15$

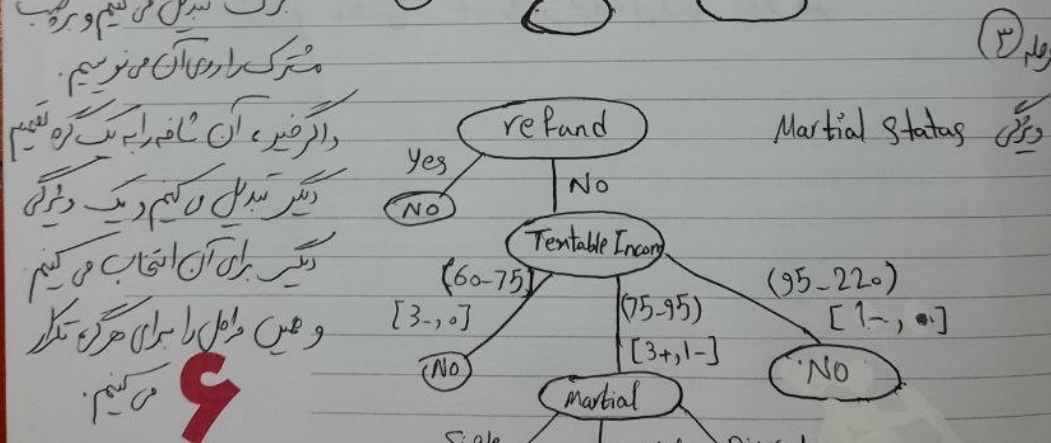
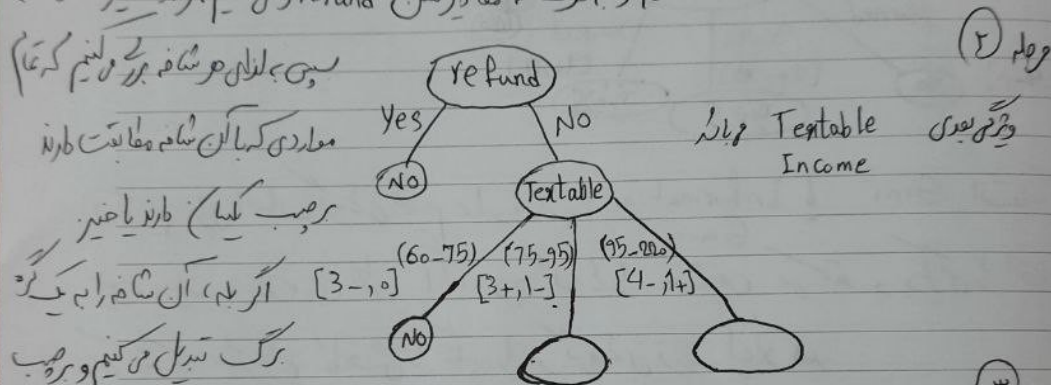
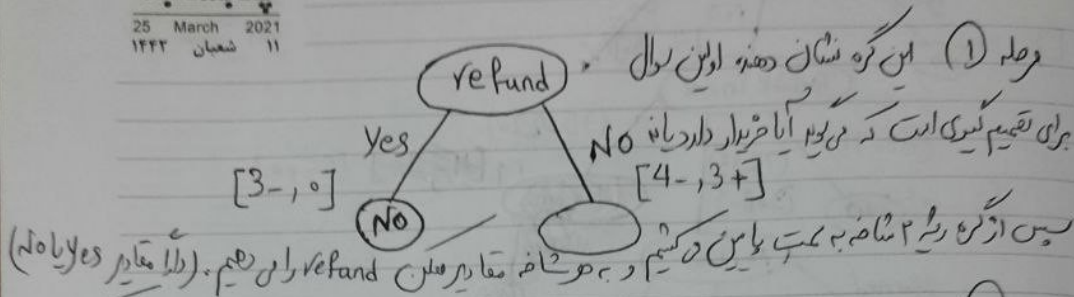


فروردین

پنجشنبه

25 March 2021
۱۴۴۲ شعبان ۱۱

گروه Refund با کمترین آنتروپی به عنوان گره انتخاب می‌کنیم



فروردین

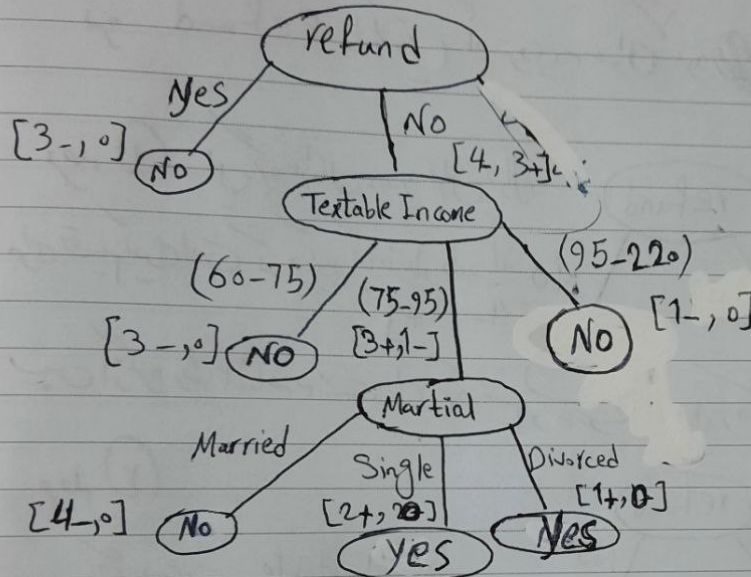
جمعه

26 March 2021
۱۴۴۲ شعبان ۱۲

Y

فروردین
شنبه

27 March 2021
۱۳ شعبان ۱۴۴۲



در اینجا برای انتخاب گره در تصمیم از معیار Information Gain یا Gini استفاده کنیم.
این معیار میزان بی نظری یا کنترس گره را مشخص می‌کند و گره را انتخاب
می‌کند که بیشترین کاهش ناظمی را در گره‌های فرزند ایجاد کند.

ه) کلاس مربوط به نمونه تست {NO, Married, 82k} طبق درخت تصمیم رسم شده در قسمت «ر» دارای لیبل NO می‌باشد.