

پاسخ تکلیف Decision Tree

درس یادگیری ماشین

امیرحسین ابوالحسنی

۴۰۰۴۰۵۰۰۳

- 1 Suppose there is an attribute, "A," that consists of random values, and these values do not have any correlation with the class labels. Additionally, assume that "A" has a sufficient number of distinct values such that no two instances in the training dataset share the same value for "A." What would be the outcome if a decision tree is built using this attribute? What challenges or issues might arise in this scenario?

پاسخ

قسمت اول

با توجه به الگوریتم ID3، در ابتدا information gain ناشی از هر ویژگی را سنجیده و آن ویژگی که بیشترین gain را دارد انتخاب می‌کنیم (تعداد کلاس‌های هدف = k):

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

$$Entropy(S_v) = - \sum_v p_v \log(p_v) = -(P(S_v = 0) \log P(S_v = 0) + P(S_v = 1) \log P(S_v = 1) + \dots + P(S_v = k) \log p(S_v = k))$$

به علت یکتایی این ویژگی (کلید اصلی بودن) برای هر نمونه، همه ترم‌های $P(S_v = l) \log P(S_v = l)$ برابر با صفر می‌شود. زیرا

$$P(S_v = l) = 0$$

یا

$$P(S_v = l) = 1$$

در نتیجه یکی از مضرب‌ها ۰ خواهد شد و کل ترم را ۰ خواهد کرد. بدین صورت است که نتیجه می‌گیریم:

$$Entropy(S_v) = 0$$

و این ویژگی برای ریشه انتخاب می‌گردد:

$$\arg \max \{Gain(S_v) | \forall A \in \text{Header}\} = A$$

در نتیجه این کار، ارتفاع درخت ۱ شده و به تعداد مقادیر ویژگی A، شاخه خواهیم داشت.

قسمت دوم

در صورتی که این ویژگی با ویژگی هدف هیچ رابطه‌ای نداشته باشد، استفاده از این ویژگی کاملاً اشتباه است و منجر به overfit می‌شود. زیرا عملاً هیچ جایی برای generalization باقی نمی‌ماند.

2 Answer the questions according to the following dataset:

Weekend	Weather	Parents	Money	Decision (Category)
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay in
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

2.1 Create a decision tree model using the given dataset to predict the value of the final column, using all other columns as input features except for the first one(weekend). Clearly explain each step of the process, including your calculations, reasoning, and decisions made while constructing the tree. What is the model's overall classification accuracy?

Root Node - \

Decision

Cinema	Tennis	Stay in	Shopping
6	2	1	1

$$Entropy(S) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S) = -(0.6 \times -0.73 + 0.2 \times -2.32 + 0.1 \times -3.32 + 0.1 \times -3.32) = 1.56$$

Money

Value	Cinema	Tennis	Stay in	Shopping
Rich	3	2	1	1
Poor	3	0	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Rich}) = -(0.42 \times -1.25 + 0.28 \times -1.83 + 0.14 \times -2.83 + 0.14 \times -2.83) = 1.82$$

$$Entropy(S_{Poor}) = -(0.42 \times -1.25) = 0.52$$

$$Gain(S, \text{Money}) = Entropy(S) - \sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{7}{10} \times 1.82 + \frac{3}{10} \times 0.52 = 1.43$$

$$Gain(S, \text{Money}) = 1.56 - 1.43 = 0.13$$

Parents

Value	Cinema	Tennis	Stay in	Shopping
Yes	5	0	0	0
No	1	2	1	1

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{\text{Yes}}) = -(\frac{5}{5} \times 0) = 0$$

$$Entropy(S_{\text{No}}) = -(0.2 \times -2.32 + 0.4 \times -1.32 + 0.2 \times -2.32 + 0.2 \times -2.32) = 1.92$$

$$Gain(S, \text{Parents}) = Entropy(S) - \sum_{v \in \text{Parents}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Parents}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{5}{10} \times 0 + \frac{5}{10} \times 1.92 = 0.96$$

$$Gain(S, \text{Parents}) = 1.56 - 0.96 = 0.6$$

Weather

Value	Cinema	Tennis	Stay in	Shopping
Sunny	1	2	0	0
Windy	3	0	0	1
Rainy	2	0	1	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{\text{Sunny}}) = -(\frac{1}{3} \times -1.59 + \frac{2}{3} \times -0.59) = 0.92$$

$$Entropy(S_{\text{Windy}}) = -(0.75 \times -0.41 + 0.25 \times -2) = 0.8$$

$$Entropy(S_{\text{Rainy}}) = -(\frac{2}{3} \times -0.59 + \frac{1}{3} \times -1.59) = 0.92$$

$$Gain(S, \text{Weather}) = Entropy(S) - \sum_{v \in \text{Weather}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Weather}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{3}{10} \times 0.92 + \frac{4}{10} \times 0.8 + \frac{3}{10} \times 0.92 = 0.87$$

$$Gain(S, \text{Weather}) = 1.56 - 0.87 = 0.69$$

Picking The Best Attribute

Attribute	Information Gain
Money	0.13
Parents	0.6
Weather	0.69

ویژگی انتخابی، Weather می باشد.

Sunny Node - ۲

Decision

Cinema	Tennis	Stay in	Shopping
1	2	0	0

For W_1, W_2, W_{10}

$$Entropy(S) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S) = -(0.33 \times -1.59 + 0.66 \times -0.59) = 0.91$$

Money

Value	Cinema	Tennis	Stay in	Shopping
Rich	1	2	0	0
Poor	0	0	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Rich}) = -(0.33 \times -1.59 + 0.66 \times -0.59) = 0.91$$

$$Entropy(S_{Poor}) = 0$$

$$Gain(S, Money) = Entropy(S) - \sum_{v \in Money} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in Money} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{3}{3} \times 0.91 + \frac{0}{3} \times 0 = 0.91$$

$$Gain(S, Money) = 0.91 - 0.91 = 0$$

Parents

Value	Cinema	Tennis	Stay in	Shopping
Yes	1	0	0	0
No	0	2	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Yes}) = -(\frac{1}{1} \times 0) = 0$$

$$Entropy(S_{No}) = -(\frac{2}{2} \times 0) = 0$$

$$Gain(S, Parents) = Entropy(S) - \sum_{v \in Parents} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in Parents} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{1}{3} \times 0 + \frac{2}{3} \times 0 = 0$$

$$Gain(S, Parents) = 0.91 - 0 = 0.91$$

Picking The Best Attribute

Attribute	Information Gain
Money	0
Parents	0.91

ویژگی انتخابی، Parents می‌باشد.

Rainy Node - ۳

Decision

Cinema	Tennis	Stay in	Shopping
2	0	1	0

For W_4, W_5, W_6

$$Entropy(S) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S) = -(0.33 \times -1.59 + 0.66 \times -0.59) = 0.91$$

Money

Value	Cinema	Tennis	Stay in	Shopping
Rich	0	0	1	0
Poor	2	0	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Rich}) = -(1 \times 0 + 1 \times 0) = 0$$

$$Entropy(S_{Poor}) = -(1 \times 0) = 0$$

$$Gain(S, Money) = Entropy(S) - \sum_{v \in Money} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in Money} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{1}{3} \times 0 + \frac{2}{3} \times 0 = 0$$

$$Gain(S, Money) = 0.91 - 0 = 0.91$$

Parents

Value	Cinema	Tennis	Stay in	Shopping
Yes	2	0	0	0
No	0	0	1	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Yes}) = -(\frac{2}{2} \times 0) = 0$$

$$Entropy(S_{No}) = -(\frac{1}{1} \times 0) = 0$$

$$Gain(S, Parents) = Entropy(S) - \sum_{v \in Parents} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in Parents} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{1}{3} \times 0 + \frac{2}{3} \times 0 = 0$$

$$Gain(S, Parents) = 0.91 - 0 = 0.91$$

Picking The Best Attribute

Attribute	Information Gain
Money	0.91
Parents	0.91

ویژگی انتخابی، Parents یا Money می باشد.

Windy Node - ۴

Decision

Cinema	Tennis	Stay in	Shopping
3	0	0	1

For W_3, W_7, W_8, W_9

$$Entropy(S) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S) = -(\frac{3}{4} \times -0.41 + \frac{1}{4} \times -2) = 0.8$$

Money

Value	Cinema	Tennis	Stay in	Shopping
Rich	2	0	0	1
Poor	1	0	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Rich}) = -(0.33 \times -1.59 + 0.66 \times -0.59) = 0.91$$

$$Entropy(S_{Poor}) = -(1 \times 0) = 0$$

$$Gain(S, \text{Money}) = Entropy(S) - \sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{3}{4} \times 0.91 + \frac{1}{4} \times 0 = 0.68$$

$$Gain(S, \text{Money}) = 0.8 - 0.68 = 0.12$$

Parents

Value	Cinema	Tennis	Stay in	Shopping
Yes	2	0	0	0
No	1	0	0	1

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{\text{Yes}}) = -(\frac{2}{2} \times 0) = 0$$

$$Entropy(S_{\text{No}}) = -(\frac{1}{2} \times -1 + \frac{1}{2} \times -1) = 1$$

$$Gain(S, \text{Parents}) = Entropy(S) - \sum_{v \in \text{Parents}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Parents}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{2}{4} \times 1 + \frac{2}{4} \times 0 = 0.5$$

$$Gain(S, \text{Parents}) = 0.8 - 0.5 = 0.3$$

Picking The Best Attribute

Attribute	Information Gain
Money	0.12
Parents	0.3

ویژگی انتخابی، Parents می باشد.

(Parents - No) Node - Δ

Decision

Cinema	Tennis	Stay in	Shopping
1	0	0	1

For W_7, W_8

$$Entropy(S) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S) = -(\frac{1}{2} \times -1 + \frac{1}{2} \times -1) = 1$$

Money

Value	Cinema	Tennis	Stay in	Shopping
Rich	0	0	0	1
Poor	1	0	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{\text{Rich}}) = -(1 \times 0) = 0$$

$$Entropy(S_{\text{Poor}}) = -(1 \times 0) = 0$$

$$Gain(S, \text{Money}) = Entropy(S) - \sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{1}{2} \times 0 + \frac{1}{2} \times 0 = 0$$

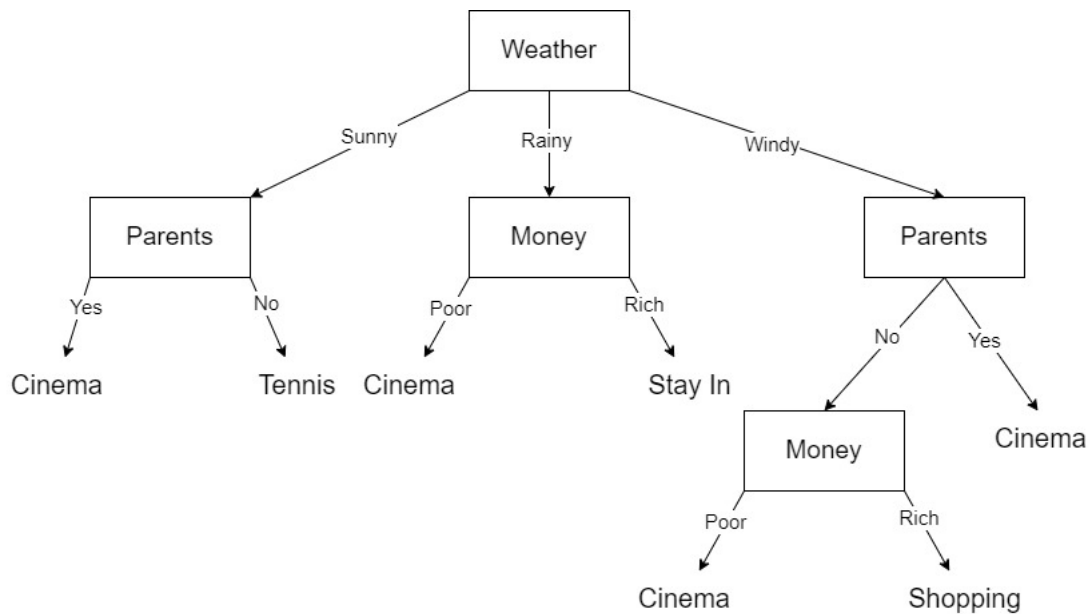
$$Gain(S, \text{Money}) = 1 - 0 = 1$$

Picking The Best Attribute

Attribute	Information Gain
Money	0.12

ویژگی انتخابی، Money می باشد.

شکل درخت تصمیم:



Accuracy: 100%

2.2 Construct a decision tree model using only the first 6 samples from the dataset(W1 - W6). Evaluate the model's classification performance on these initial 6 samples as the training set. Then, use the model to classify the remaining samples in the dataset. What is the classification accuracy for both the training and test datasets? Discuss your findings and explain the reasons behind the observed results.

یادگیری با داده های آموزشی:

Root Node - ۱

Decision

Cinema	Tennis	Stay in	Shopping
4	1	1	0

$$Entropy(S) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S) = -(\frac{4}{6} \times -0.59 + \frac{1}{6} \times -2.64 + \frac{1}{6} \times -2.64) = 1.23$$

Money

Value	Cinema	Tennis	Stay in	Shopping
Rich	2	1	1	0
Poor	2	0	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Rich}) = -(\frac{2}{4} \times -1 + \frac{1}{4} \times -2 + \frac{1}{4} \times -2) = 1.5$$

$$Entropy(S_{Poor}) = -(\frac{2}{2} \times 0) = 0$$

$$Gain(S, Money) = Entropy(S) - \sum_{v \in Money} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in Money} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{4}{6} \times 1.5 + \frac{2}{6} \times 0 = 1$$

$$Gain(S, Money) = 1.23 - 1 = 0.23$$

Parents

Value	Cinema	Tennis	Stay in	Shopping
Yes	4	0	0	0
No	0	1	1	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Yes}) = -(\frac{4}{4} \times 0) = 0$$

$$Entropy(S_{No}) = -(\frac{1}{2} \times -1 + \frac{1}{2} \times -1) = 1$$

$$Gain(S, Parents) = Entropy(S) - \sum_{v \in Parents} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in Parents} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{4}{6} \times 0 + \frac{2}{6} \times 1 = 0.33$$

$$Gain(S, Parents) = 1.23 - 0.33 = 0.9$$

Weather

Value	Cinema	Tennis	Stay in	Shopping
Sunny	1	1	0	0
Windy	1	0	0	0
Rainy	2	0	1	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{Sunny}) = -(\frac{1}{2} \times -1 + \frac{1}{2} \times -1) = 1$$

$$Entropy(S_{Windy}) = -(\frac{1}{1} \times 0) = 0$$

$$Entropy(S_{Rainy}) = -(\frac{2}{3} \times -0.59 + \frac{1}{3} \times -1.59) = 0.92$$

$$Gain(S, Weather) = Entropy(S) - \sum_{v \in Weather} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in Weather} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{2}{6} \times 1 + \frac{1}{6} \times 0 + \frac{3}{6} \times 0.92 = 0.79$$

$$Gain(S, Weather) = 1.23 - 0.79 = 0.44$$

Picking The Best Attribute

Attribute	Information Gain
Money	0.23
Parents	0.9
Weather	0.44

ویژگی انتخابی، Parents می‌باشد.
(Parents - No) Node - ۲

Decision

Cinema	Tennis	Stay in	Shopping
0	1	1	0

For W_2, W_5

$$Entropy(S) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S) = -(\frac{1}{2} \times -1 + \frac{1}{2} \times -1) = 1$$

Money

Value	Cinema	Tennis	Stay in	Shopping
Rich	0	1	1	0
Poor	0	0	0	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{\text{Rich}}) = -(\frac{1}{2} \times -1 + \frac{1}{2} \times -1) = 1$$

$$Entropy(S_{\text{Poor}}) = -(0 \times 0) = 0$$

$$Gain(S, \text{Money}) = Entropy(S) - \sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Money}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{2}{2} \times 1 = 1$$

$$Gain(S, \text{Money}) = 1 - 1 = 0$$

Weather

Value	Cinema	Tennis	Stay in	Shopping
Sunny	0	1	0	0
Windy	0	0	0	0
Rainy	0	0	1	0

$$Entropy(S_v) = - \sum_{v \in S} p_v \log(p_v)$$

$$Entropy(S_{\text{Sunny}}) = -(\frac{1}{1} \times 0) = 0$$

$$Entropy(S_{\text{Windy}}) = 0$$

$$Entropy(S_{\text{Rainy}}) = -(\frac{1}{1} \times 0) = 0$$

$$Gain(S, \text{Weather}) = Entropy(S) - \sum_{v \in \text{Weather}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$\sum_{v \in \text{Weather}} \frac{|S_v|}{|S|} Entropy(S_v) = \frac{1}{2} \times 0 + \frac{1}{2} \times 0 = 0$$

$$Gain(S, \text{Weather}) = 1 - 0 = 1$$

Picking The Best Attribute

Attribute	Information Gain
Money	0
Weather	1

ویژگی انتخابی، Weather می باشد.

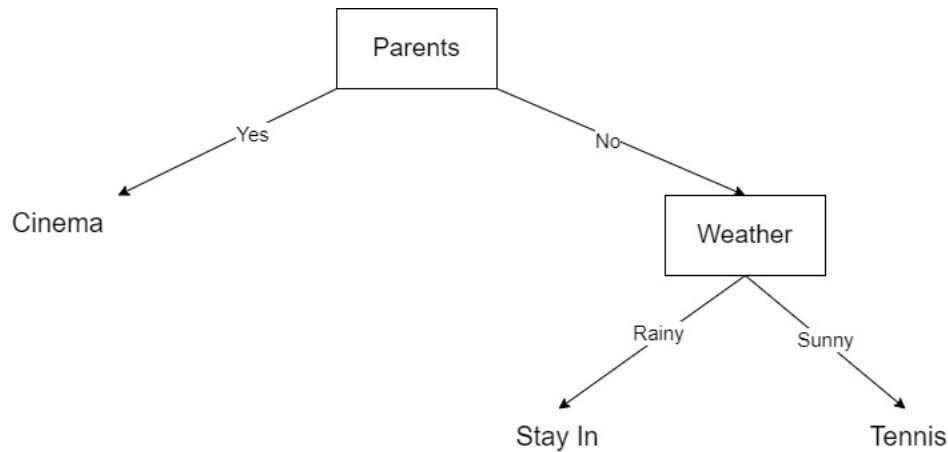
دقت مدل روی داده های آموزشی

$$Accuracy_{Train} = \frac{\text{Correctly Labeled}}{\text{All train samples}} = \frac{6}{6} = 100\%$$

دقت مدل روی داده های تست

$$Accuracy_{Test} = \frac{\text{Correctly Labeled}}{\text{All test samples}} = \frac{2}{4} = 50\%$$

شکل درخت تصمیم:



به نظر می رسد یکی از دلایل اصلی افت مقدار دقت در داده های تست، نبود نمونه هایی با مقادیر مختلف برای ویژگی ها در داده های آموزش بوده است، برای مثال برای مقدار Windy هیچ پیشبینی را نمی توان انجام داد زیرا مدل آن را تا الان بیرون از تسلط ویژگی Parents ندیده است.

2.3 In scenarios where only a limited number of labeled examples are available for training (and no extra data is available for testing or validation), propose a specific pruning technique that could be integrated into the decision tree algorithm to prevent overfitting. Justify why you believe this technique would be effective.

با توجه به اینکه داده تستی نداریم، از این روش که بگذاریم درخت overfit کند و سپس Post Pruning به طریقی دیگر استفاده می کنیم. می توان در الگوریتم اضافه کرد، پس از اتمام کار، از هر برگ شروع میکنیم و در صورتی که مقدار آنتروپی آن از آستانه ای کمتر بود، به راس والد رفته و اینکار را انجام می دهیم تا جایی که این اتفاق نیفتد، سپس از آن راس که آنتروپی آن دیگر کمتر از مقدار آستانه نیست، زیر درخت آن راس را حذف می کنیم تا از بیش برآزش جلوگیری به عمل آمده و قابلیت تعمیم پذیری مدل بالاتر رود.