



گزارش تمرین عملی دوم درس هوش مصنوعی
آریا جلالی
۹۸۱۰۵۶۶۵

۱ کتابخانه‌های استفاده شده:

```
1 import pandas as pd
2 import numpy as np
3 import gvgen
4 import os
5 import matplotlib.pyplot as plt
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import confusion_matrix,
  ConfusionMatrixDisplay
8 from graphviz import Source
```

از کتابخانه‌ی numpy برای ذخیره کردن ویژگی‌های متفاوت رستوران‌ها و افراد دیابتی استفاده شده است و از کتابخانه‌ی pandas برای خواندن فایل‌های csv استفاده شده است. از کتابخانه‌های gvgen و matplotlib و graphviz برای نشان دادن درخت نهایی استفاده شده است. در نهایت از کتابخانه‌ی sklearn برای انتخاب سمپل‌های رندوم و محاسبه خطای برنامه استفاده شده است که در ادامه کاربرد آن‌ها را به صورت مفصل توضیح خواهیم داد.

۲ توضیح بخش‌های مختلف کد:

```
1 # function for discretizing our continuous data.
2 def discretize(data: np.array, number_of_bins):
3     global bins
4     for i in range(8):
5         bins.append(pd.cut(data[i], number_of_bins, retbins=True)
6                        [1])
7         data[i] = np.digitize(data[i], bins[i], right=True)
```

این تابع وظیفه‌ی گسسته‌سازی اطلاعات برای افراد دیابتی را دارد و روش استفاده شده. این تابع به این صورت کار می‌کند که برای هر ستون ابتدا تابع pd.cut را صدا می‌زند و این تابع یک لیست را برمی‌گرداند که هر 2 عضو مجاور آن نشان دهنده‌ی یک بازه هستند و در نهایت می‌توان اطلاعات هر ستون را با توجه به مقدار آن‌ها با استفاده از تابع np.digitize که دو ورودی داده‌ها و لیست bin ها را می‌گیرد در بازه‌ی متناظر خود قرار داد.

```
1 def test(row: np.array, root):
2     global Correct_Guess
3     if len(root.children) == 0:
4         y_pred.append(int(root.plus > root.negative))
5         if int(root.plus > root.negative) == int(row[-1]):
6             Correct_Guess += 1
7         return
8     if row[root.attribute] not in root.children.keys():
9         if np.random.rand() > 0.5:
10            y_pred.append(int(row[-1]))
11            Correct_Guess += 1
12     else:
```

```

13         y_pred.append(1 - int(row[-1]))
14     return
15     test(row, root.children[row[root.attribute]])

```

تابع test برای تست نهایی برنامه استفاده می‌شود و برای تمامی افراد تابع test را به صورت بازگشتی تا به رسیدن یک برگ صدا می‌زند و نتیجه را بررسی می‌کند. دقت کنید اگر مقداری را برای اولین بار بر روی درخت ببینیم با استفاده از تابع random بررسی می‌کنیم آیا جواب ما درست بوده است یا خیر. با دیدن هر جواب درست متغیر Correct_Guess global را یکی افزایش می‌دهیم و مقدار پیشبینی شده را در لیستی به نام y_pred قرار می‌دهیم تا در نهایت نتایج را با استفاده از Confusion Matrix بررسی کنیم.

```

1 # Decision tree
2 class Tree:
3     def __init__(self, data, parent_attribute=None, depth=0):
4         if parent_attribute is None:
5             parent_attribute = []
6         self.children = dict()
7         self.attribute = -1
8         self.depth = depth
9         self.attribute_name = 'Leaf'
10        self.information_gain = -1
11        self.entropy = -1
12        self.plus = 0
13        self.negative = 0
14        self.parent_attribute = parent_attribute
15        self.data = np.array(data)
16
17    def __str__(self):
18        if len(self.children) == 0:
19            return self.attribute_name + "\n" + "Entropy: " +
20            str(self.entropy) + "\n" + "[" + str(
21                self.plus) + "/" + str(
22                self.negative) + "]"
23            return self.attribute_name + "\n" + "Gain: " + str(self.
24            information_gain) + "\n" + "[" + str(
25                self.plus) + "/" + str(
26                self.negative) + "]"

```

در این بخش و ادامه کلاس اصلی درخت تصمیم که Tree نام دارد را بررسی می‌کنیم. این کلاس در constructor خود 3 ورودی می‌گیرد که این ورودی‌ها به ترتیب برابر با data, parent_attributes, depth هستند که در ادامه به توضیح تک تک آن‌ها خواهیم پرداخت.

ورودی data همان افراد مورد بررسی هر Node را نشان می‌دهد که برای تمامی اعضا را شامل می‌شود و در مراحل بعد با انتخاب ویژگی با توجه به معیارهای متفاوت (information gain یا gini) افراد فیلتر می‌شوند و به مراحل بعدی داده می‌شوند.

parent_attribute برای هر Node نشان دهنده ویژگی‌هایی است که پدر و اجداد آن تا

الان بررسی کرده‌اند و در این مرحله باید skip شوند و نیازی به بررسی آن‌ها نیست. دقت کنید برای root این ورودی برابر با یک لیست خالی است، زیرا هنوز ویژگی‌ای انتخاب نشده است.

depth نشان دهنده‌ی عمق Node مورد بررسی است که در ابتدا برای root برابر با صفر قرار داده شده است و فرزندان هر Node دارای عمق $depth + 1$ هستند. این ورودی depth برای محدود کردن سائز درخت و جلوگیری از overfit شدن استفاده شده است که در ادامه در رابطه با آن توضیحاتی خواهیم داد.

children فرزندان یک Node را نشان می‌دهد که آن را با استفاده از یک dictionary نشان می‌دهیم که key نشان دهنده‌ی ویژگی مورد بررسی در Node فرزند و value متناظر با آن کلاس فرزند است.

plus & negative به ترتیب نشان دهنده‌ی این هستند که در data فعلی چند نفر از افراد دیابت دارند و چند نفر ندارند. از این ویژگی‌ها برای محاسبه‌ی معیارهایی مانند information gain یا gini و نشان دادن روی درخت نهایی به صورت گرافیکی استفاده خواهیم کرد.

```
1 def train(self):
2     self.plus = np.sum(self.data.T[-1])
3     self.negative = self.data.shape[0] - self.plus
4     self.attribute_name = str(int(self.plus > self.negative))
5     self.entropy = calculate_entropy(self.plus, self.negative)
6
7     if self.depth == max_depth or self.plus == 0 or self.
negative == 0:
8         return
9
10    difference = dict()
11
12    for i in range(len(self.data[0]) - 1):
13        if i in self.parent_attribute:
14            continue
15        for row in self.data:
16            if row[i] in difference:
17                difference[row[i]].append(row)
18            else:
19                difference[row[i]] = []
20                difference[row[i]].append(row)
21        self.choose_attribute(difference, i)
22        if self.attribute == i:
23            self.children.clear()
24            next_parent_attribute = self.parent_attribute.copy()
25            next_parent_attribute.append(self.attribute)
26            for key, value in difference.items():
27                self.children[key] = Tree(value, depth=self.
depth + 1, parent_attribute=next_parent_attribute)
```

```

28         difference.clear()
29
30     for key, value in self.children.items():
31         value.train()
32
33     if len(self.children) == 0:
34         self.attribute_name = str(int(self.plus > self.negative)

```

تابع train تابع اصلی ما را برای ساخت درخت تصمیم تشکیل می‌دهد. نحوه‌ی کار به این صورت است که ابتدا یک dictionary به نام difference می‌گیریم و در ادامه روی تمام 8 ویژگی داده شده loop می‌زنیم و هر مقدار جدیدی را که می‌بینیم آن را به عنوان یک key در difference dictionary قرار می‌دهیم و برای هر key یک لیست می‌گیریم که نقطه مشترک این افراد این است که ویژگی i ام آن‌ها برابر با key متناظر با لیستی هستند که در آن قرار دارند.

در ادامه تابع choose_attribute با ورودی های i و difference صدا می‌زنیم که نحوه‌ی کار آن را در ادامه توضیح خواهیم داد در ادامه بررسی می‌کنیم که آیا این ویژگی که آن را مورد بررسی قرار دادیم بهتر از ویژگی‌های پیشین است؟ و در آن صورت dictionary فرزندان را clear می‌کنیم و فرزندان جدید را از difference می‌خوانیم و به children اضافه می‌کنیم. در نهایت نیز در صورت داشتن فرزندان و برگ نبودن تابع train را بر روی تک تک فرزندان صدا می‌زنیم.

```

1 def choose_attribute(self, difference, attribute):
2     remainder_A = 0
3     for key, value in difference.items():
4         plus = 0
5         negative = 0
6         for each in value:
7             if each[-1] == 1:
8                 plus += 1
9             else:
10                negative += 1
11        remainder_A += ((plus + negative) / (self.plus + self.
12        negative)) * calculate_entropy(plus, negative)
13    if self.information_gain < self.entropy - remainder_A:
14        self.information_gain = self.entropy - remainder_A
15        self.attribute = attribute
16        self.attribute_name = HEADERS[attribute]

```

در تابع choose_attribute مقدار remainder(A) به صورت مقابل با استفاده از یک for loop محاسبه می‌شود.

$$Remainder(A) = \sum_{k=1}^d \frac{p_k + n_k}{p + n} B\left(\frac{p_k}{p_k + n_k}\right)$$

که در معادله بالا d نشان دهنده‌ی مقادیر ممکن و p_k و n_k به ترتیب نشان دهنده‌ی افراد دیابتی و غیر دیابتی‌ای هستند که در زیرمجموعه k قرار دارند. و در نهایت تابع Entropy B

را از طریق فرمول مقابل محاسبه می‌کند.

$$B(q) = -(q * \log_2 * q + (1 - q) * \log_2 * (1 - q))$$

در نهایت مقدار information gain را از طریق فرمول مقابل محاسبه می‌کنیم.

$$Gain(A) = B(\frac{p}{p+n}) - Remainder(A)$$

و در صورتی که $Gain$ محاسبه شده بیشتر از مقدار فعلی باشد ویژگی attribute کلاس به ویژگی مورد بررسی تغییر پیدا می‌کند.

```
1 def calculate_entropy(plus, negative):
2     q = plus / (plus + negative)
3     if q == 0 or q == 1:
4         return 0
5     return -(q * np.log2(q) + (1 - q)
6             * np.log2(1 - q))
```

در تابع بالا نیز تعداد افراد دیابتی و سالم به تبدیل به متغیر q می‌شوند و entropy طبق فرمول بالا محاسبه می‌شود.

```
1 def get_edge_name(key, val: Tree):
2     attribute_bin = bins[val.attribute]
3     key = int(key)
4     return str(round(attribute_bin[key - 1], 2)) + " <= x <= " +
5           str(round(attribute_bin[key], 2))
6
7 def graphMaker(g, my_tree: Tree):
8     if len(my_tree.children) == 0:
9         myItem = g.newItem(my_tree.__str__())
10        return myItem
11    else:
12        myItem = g.newItem(my_tree.__str__())
13        for key, val in my_tree.children.items():
14            newTree = graphMaker(g, val)
15            link = g.newLink(myItem, newTree)
16            g.propertyAppend(link, "color", "darkblue")
17            g.propertyAppend(link, "label", get_edge_name(key,
18my_tree))
19        return myItem
20
21 def makeVisualGraph(my_tree):
22     g = gvgen.GvGen()
23     graphMaker(g, my_tree)
24     string = ""
25     my_file = open("output_graphviz.txt", 'w')
26     g.dot(my_file)
```

```

27 my_file.close()
28 my_file = open("output_graphviz.txt", 'r')
29 lines = my_file.readlines()[1:]
30 for line in lines:
31     string = string + line
32     src = Source(string)
33     src.render(view=True)

```

از 3 تابع بالا برای نمایش گرافیکی درخت تصمیم استفاده کرده‌ام به نحوی که ابتدا تابع `makeVisualGraph` صدا زده می‌شود و آن نیز تابع `graphMaker` را صدا می‌زند که آن با اجرای یک `dfs` و استفاده از کتابخانه‌ی `gvgen` یک درخت به زبان `dot` که زبان رسمی موتور گرافیکی `graphviz` است را در یک فایل `txt` می‌نویسد و در ادامه آن را تبدیل به یک فایل `dot` می‌کنیم و با استفاده از موتور گرافیکی `graphviz` درخت تصمیم را نمایش می‌دهیم.

```

1 if __name__ == '__main__':
2     # read attributes from file.
3     df = pd.read_csv('diabetes.csv')
4     HEADERS = df.columns
5
6     # Put all attributes into a numpy matrix.
7     attributes = df.to_numpy().T
8     discretize(data=attributes, number_of_bins=5)
9     attributes_train, attributes_test = train_test_split(
10         attributes.T, test_size=0.2)
11
12     root = Tree(attributes_train)
13     root.train()
14     makeVisualGraph(root)
15
16     for row in attributes_train:
17         test(row, root)
18     print("Accuracy for train data: " + str(round(Correct_Guess /
19         len(attributes_train) * 100, 2)) + "%")
20     Correct_Guess = 0
21
22     cm = confusion_matrix(attributes_train.T[-1], y_pred)
23
24     disp = ConfusionMatrixDisplay(confusion_matrix=cm)
25     disp.plot()
26     plt.show()
27
28     y_pred = []
29
30     for row in attributes_test:
31         test(row, root)
32     print("Accuracy for test data: " + str(round(Correct_Guess /
33         len(attributes_test) * 100, 2)) + "%")
34
35     cm = confusion_matrix(attributes_test.T[-1], y_pred)

```

```

33
34     disp = ConfusionMatrixDisplay(confusion_matrix=cm)
35     disp.plot()
36     plt.show()
37
38     average = 0
39     for i in range(100):
40         attributes_train, attributes_test = train_test_split(
41             attributes.T, test_size=0.2)
42
43         root = Tree(attributes_train)
44         root.train()
45
46         for row in attributes_test:
47             test(row, root)
48         average += round(Correct_Guess / len(attributes_test) *
49             100, 2)
50         Correct_Guess = 0
51
52     print("Average Accuracy of tree on test data is : " + str(
53         round(average / 100, 2)))

```

تکه کد بالا که بخش اصلی کد را نشان می‌دهد ابتدا data را از طریق فایل داده شده می‌خواند و در ادامه با استفاده از تابع‌های گفته شده درخت تصمیم را تشکیل می‌دهد و در ادامه با استفاده از تابع test درخت تصمیم را بر روی attributes_test اجرا می‌کند و دقت برنامه را به کاربر نشان می‌دهد و در نهایت این کار را برای 100 بار تکرار می‌کند تا مقدار متوسطی از عملکرد درخت داشته باشیم.

برای جدا کردن data به 2 بخش train و test از تابع test_train_split کتابخانهی sklearn استفاده می‌کنیم.

۳ بررسی معیارهای استفاده شده:

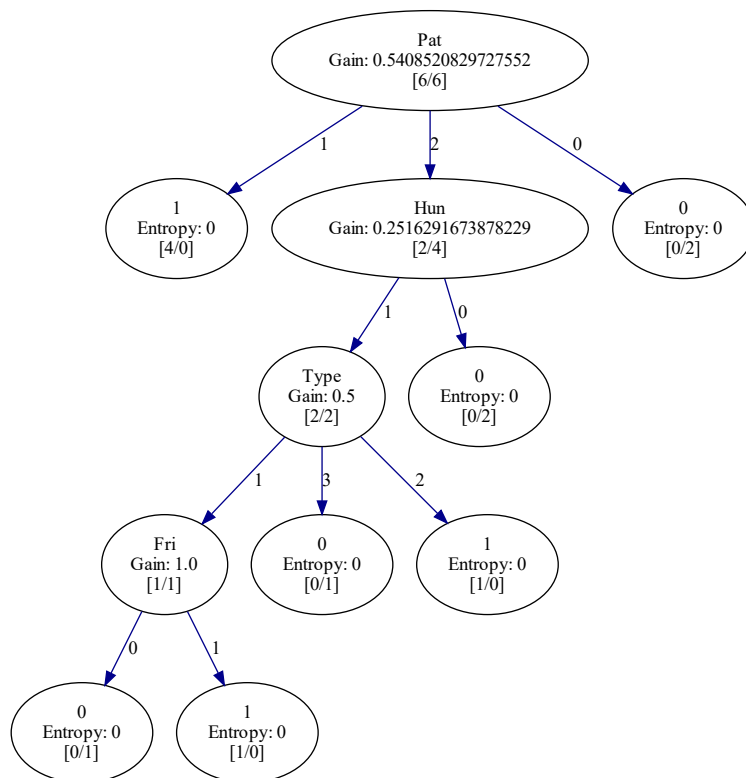
همانطور که در بخش‌های قبل‌تر نیز گفته شد علاوه بر معیار information gain از معیار gini index نیز استفاده شده است که محاسبه‌ی آن به صورت مقابل است:

$$Gini\ Impurity(k) = 1 - \left(\frac{p_k}{p_k + n_k}\right)^2 - \left(\frac{n_k}{p_k + n_k}\right)^2$$

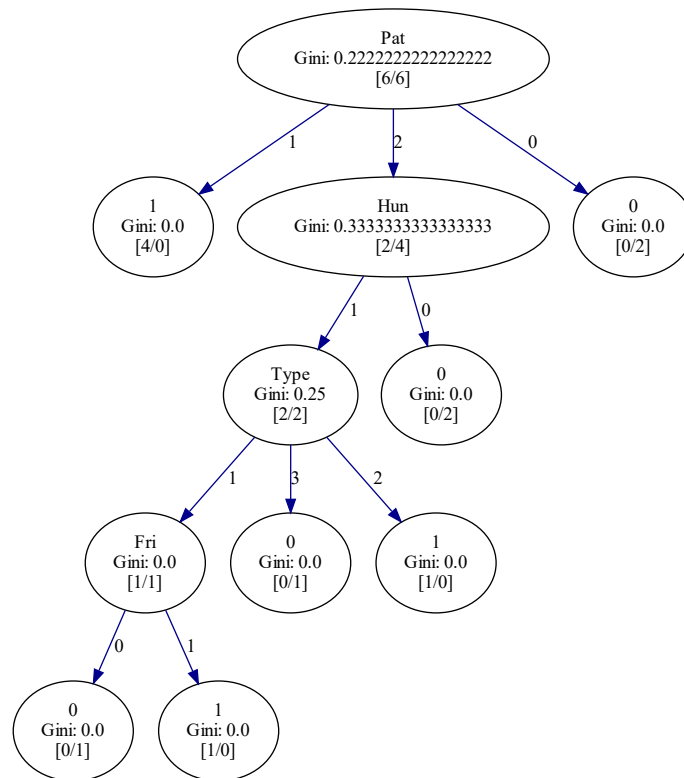
$$Gini\ Index = \sum_{k=1}^d Gini\ Impurity(k) \times \frac{p_k + n_k}{p + n}$$

دقت کنید معیار gini index برخلاف معیار information gain در مقادیر پایین نشان دهنده‌ی بهتر بودن ویژگی مورد بررسی است و از طرفی محاسبه‌ی آن نسبت به information gain از لحاظ محاسباتی آسان‌تر است و در ادامه نتایج این 2 معیار را مقایسه خواهیم کرد.

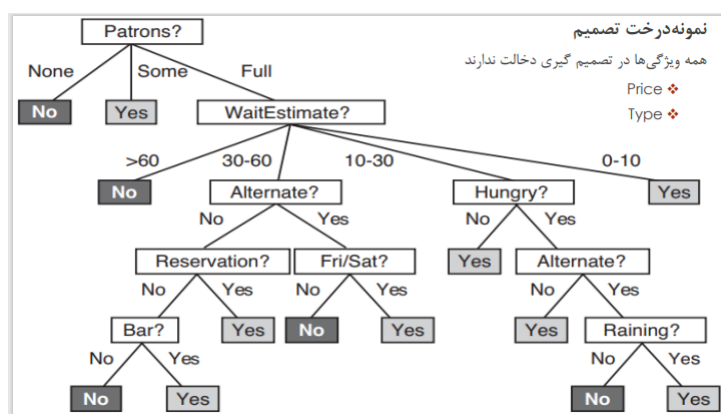
۴ درخت تولید شده برای داده‌های رستوران:



شکل ۱: درخت تصمیم تولید شده با معیار information gain



شکل ۲: درخت تصمیم تولید شده با معیار gini index



شکل ۳: نمونه درخت تولید شده در اسلاید درس

همانطور که از تصاویر بالا نیز مشخص است هر 3 درخت تولید شده یکسان شده‌اند و دلیل این عمل برای درخت اول و سوم واضح است زیرا هر 2 یک الگوریتم را پیاده‌سازی کرده‌اند ولی جالب است که معیار gini index نیز همان درخت را تشکیل داده است. دلیل اینکار را می‌توان به این حقیقت که تمام داده‌ها برای ساخت درخت استفاده شده‌اند و باید دقت آن‌ها 1 شود نسبت داد.

دقت کنید در معیار index gini روی هر Node مقدار gini impurity آن نود روی آن نوشته شده است و در معیار information gain روی Node های میانی مقدار Gain و بر روی برگ‌ها مقدار entropy نوشته شده است و در هر دو معیار ویژگی مورد بررسی در هر Node روی آن نوشته شده است و در نهایت به ترتیب تعداد افراد دیابتی و سالم در هر Node به فرم [Plus,Negative] نوشته شده است.

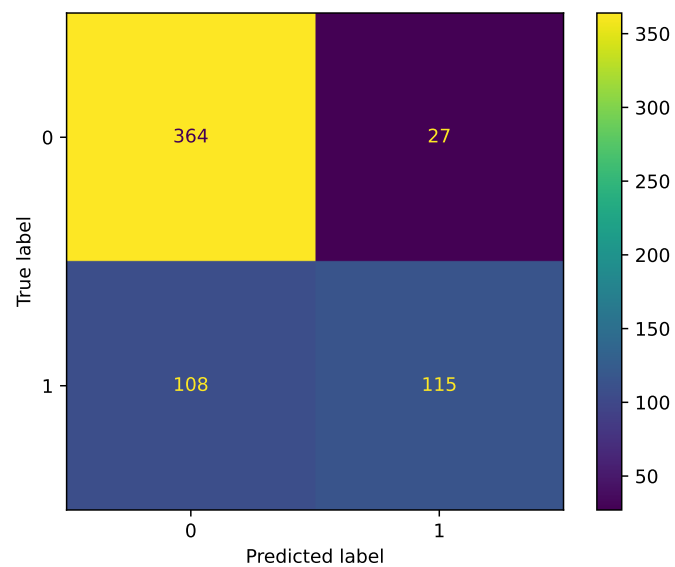
۵ بررسی دقت درخت‌های تولید شده بر روی داده‌های تست:



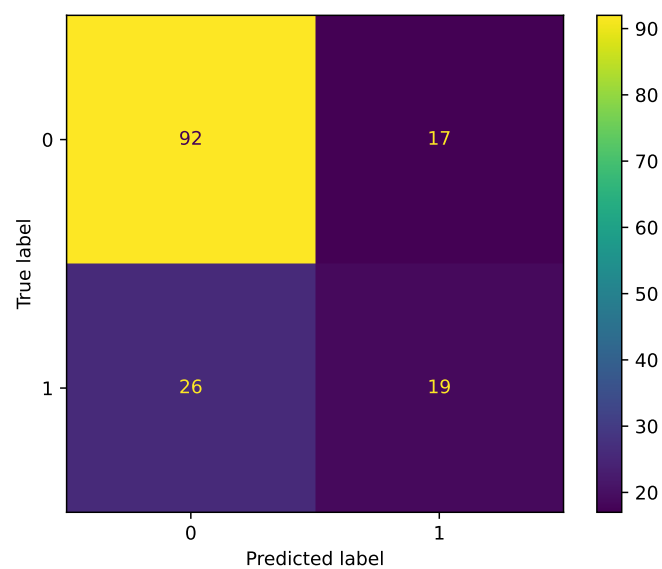
شکل ۴: درخت تولید شده با استفاده از معیار gini index

دقت کنید به دلیل بزرگ بودن درخت تصویر بالا کوچک نشان داده شده است ولی به دلیل بودن فرمت عکس می‌توانید هر مقدار که بخواهید روی آن zoom کنید از طرفی تمامی عکس‌ها در فایل‌های فرستاده شده قرار داده شده‌اند.

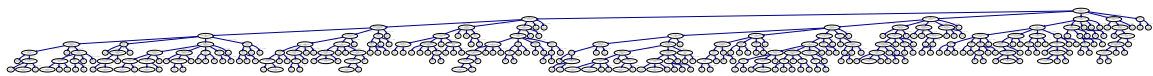
درخت بالا با حداکثر عمق 3 و تعداد 5 bin طراحی شده است و دقت آن برای داده‌های train برابر با 78.01% و برای داده‌های test برابر با 72.08% است و مقدار متوسط آن پس از 100 بار اجرا روی داده‌های تست برابر با 72.7% است.



شکل ۵: ماتریکس confusion برای داده‌های train

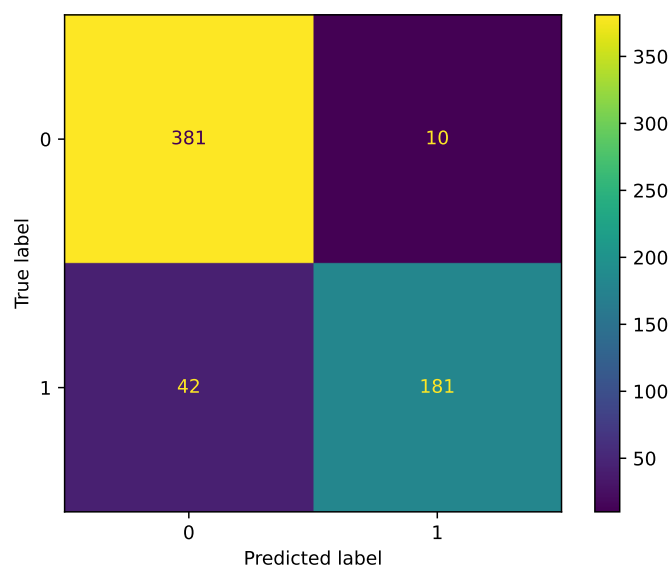


شکل ۶: ماتریکس confusion برای داده‌های test

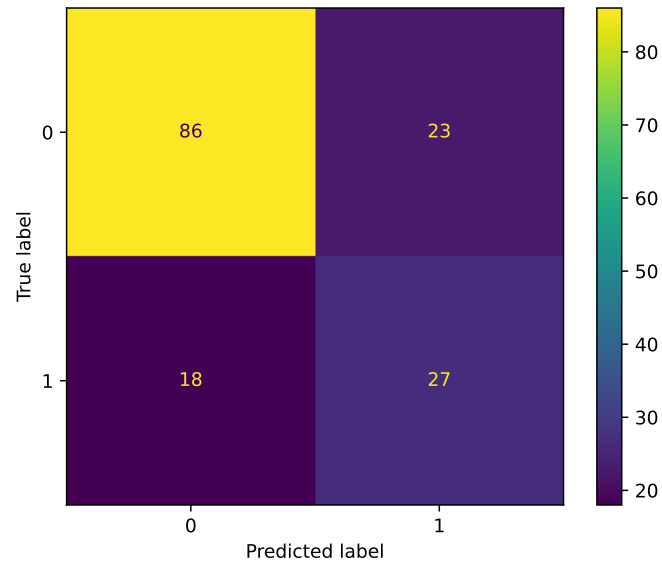


شکل ۷: درخت تولید شده با معیار gini index

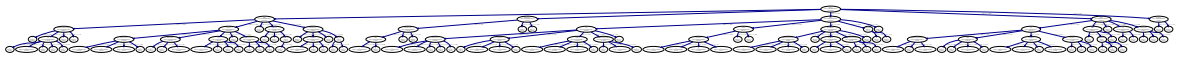
درخت بالا با حداکثر عمق 10 و تعداد 5 bin طراحی شده است و دقت آن برای داده‌های train برابر با 91.53% و برای داده‌های test برابر با 73.38% است و مقدار متوسط آن پس از 100 بار اجرا روی داده‌های تست برابر با 69.28% است.



شکل ۸: ماتریکس confusion برای داده‌های train

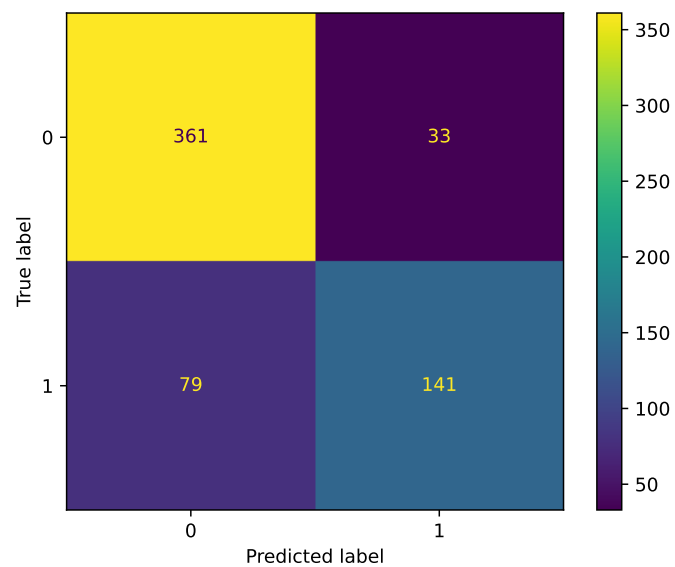


شکل ۹: ماتریکس confusion برای داده‌های test

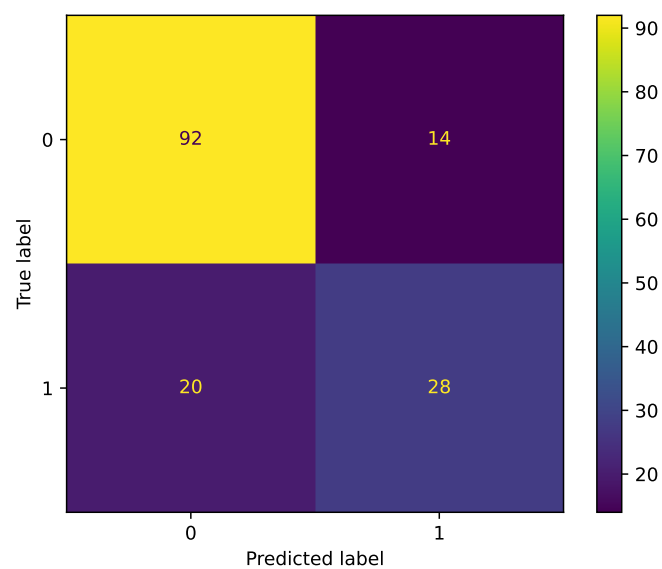


شکل ۱۰: درخت تولید شده با استفاده از معیار information gain

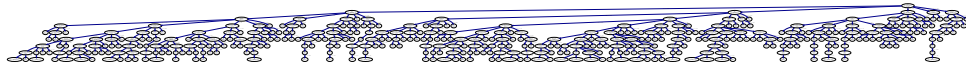
درخت بالا با حداکثر عمق 4 و تعداد 5 bin طراحی شده است و دقت آن برای داده‌های train برابر با 81.76% و برای داده‌های test برابر با 77.92% است و مقدار متوسط آن پس از 100 بار اجرا روی داده‌های تست برابر با 71.5% است.



شکل ۱۱: ماتریکس confusion برای داده‌های train

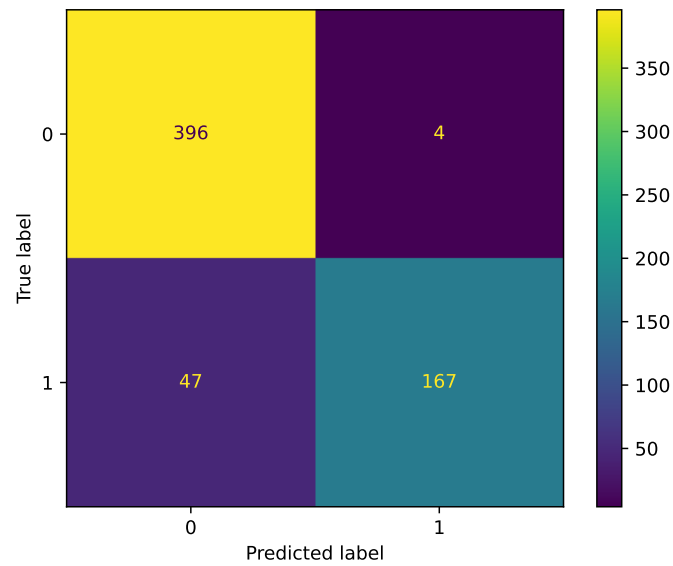


شکل ۱۲: ماتریکس confusion برای داده‌های test

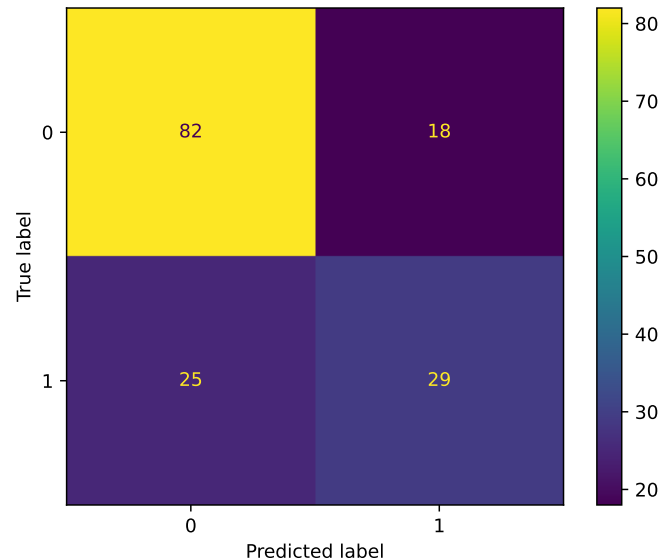


شکل ۱۳: درخت تولید شده با استفاده از معیار information gain

درخت بالا با حداکثر عمق 10 و تعداد 5 bin طراحی شده است و دقت آن برای داده‌های train برابر با 91.69% و برای داده‌های test برابر با 72.08% است و مقدار متوسط آن پس از 100 بار اجرا روی داده‌های تست برابر با 69.92% است.



شکل ۱۴: ماتریکس confusion برای داده‌های train



شکل ۱۵: ماتریکس confusion برای داده‌های test

۶ نتایج نهایی:

با چندین بار اجرای برنامه برای هر 2 معیار gini index و information gain متوجه شدم گرفتن دقت بالای 80% برای داده‌های test تقریباً با روش decision tree ناممکن است و حتی کتابخانه‌ی معتبر و قدیمی sklearn نیز در بهترین حالت دقتی برابر با 79% درصد داشت و به نظر می‌رسد دقت موردانتظار این الگوریتم با معیارهایی مانند gini index و information gain باید حدود 70% باشد. با این حال با نگاهی به داده‌ی دیابت متوجه می‌شویم بسیاری از داده‌ها ثبت نشده‌اند و به عنوان مثال افراد زیادی دارای فشار خونی و ضخامت پوست و BMI صفر هستند و می‌توانستم این مقادیر را در نظر بگیرم و یا مقادیر 0 را با میانگین آن‌ها جایگزین کنم ولی در آن صورت درخت من general نخواهد بود و برای همین dataset خاص طراحی شده و باید برای هر dataset یک کد جدید یا تغییراتی روی کد قبلی داد.

در نهایت راجع به موضوع overfit متوجه شدم هرچه تعداد bin ها و حداکثر عمق ممکن بیشتر باشد پدیده‌ی overfit بیشتر رخ می‌دهد که این اتفاق از نتیجه‌های بالا مشخص است که تنها افزایش حداکثر عمق ممکن باعث شد دقت بر روی داده‌های train افزایش چشمگیری داشته باشد و در طرف دیگر دقت بر روی داده‌های test کاهش یابد. برای کاهش مشکل overfit حداکثر عمق ممکن را برای هر دو معیار را برابر با 3 قرار دادم و تعداد bin ها را نیز برابر با 5 قرار دادم. این 2 hyper parameter نسبت به بقیه نتایج بهتری کسب کردند.

در ادامه با مقایسه بین 2 معیار gini index و information gain متوجه شدم با اینکه ممکن

است دقت information gain تا حد کمی بیشتر باشد ولی هزینه‌ی محاسبات آن نیز بیشتر از معیار gini index است و باید به نحوی بین دقت بالاتر و محاسبات بیشتر trade-off انجام شود. البته با توجه به نتایج بالا متوجه می‌شویم انجام دادن محاسبات بیشتر برای افزایش دقت به اندازه‌ی 1 یا 2 درصد ارزش ندارد.

۷ فایل‌های ارسال شده:

تمام کدهای اجرایی همراه با عکس‌ها و فایل README برای اجرای کد در فایل زیپ فرستاده شده قرار داده شده‌اند.