

## کاربرد درخت رده بندی برای پیشگویی رتینوپاتی دیابتیک و مقایسه با مدل تابع تشخیص در بیماران دیابتی نوع ۲

دکتر سید محسن حسینی<sup>۱</sup>، دکتر مهدی تذهیبی<sup>۱</sup>، دکتر مسعود امینی<sup>۲</sup>، اصغر زارع<sup>۳</sup>، دکتر حسن جهانی هاشمی<sup>۴</sup>

### خلاصه

**مقدمه:** بیماری دیابت یکی از بیماری‌های شایع در جهان به شمار می‌رود و رتینوپاتی دیابتی به عنوان یک بیماری چشمی، این گروه از بیماران را به شدت درگیر می‌نماید و یکی از عوامل کوری آنان در سنین بالا به شمار می‌رود. با توجه به شیوع بالای دیابت نوع II در جامعه و همچنین خطر بالای رتینوپاتی در این گروه از بیماران، در این مطالعه سعی گردید تا با توجه به مدل رده بندی درختی CART (Classification and Regression Tree) یک الگوی پیش‌گویی و رده‌بندی برای این بیماران معرفی شود.

**روش‌ها:** این بررسی با توجه به اطلاعات ۳۷۳۴ بیمار دیابتی نوع II انجام شد. بیماران به صورت تصادفی به دو گروه نمونه‌ی آزمون و یادگیری تقسیم شدند. با توجه به نمونه‌ی یادگیری، مدل رده بندی درخت CART معرفی و دقت آن با توجه به نمونه‌ی آزمون مشخص گردید. نتایج به دست آمده از الگوی درختی CART با نتایج به دست آمده از مدل تابع تشخیص برای کل بیماران مقایسه شد.

**یافته‌ها:** در این مطالعه، مدل رده بندی درختی با دقت ۰/۶۷ با توجه به دوره بیماری، فشار خون سیستولیک، سن، جنسیت، تری‌گلیسرید، کلسترول و قند خون ناشتا به دست آمد. این الگوی رده بندی از حساسیت ۰/۷۱ و ویژگی ۰/۶۲ برخوردار بود. در این بررسی میزان دقت رده بندی درختی معرفی شده با میزان دقت الگوی رده بندی بر اساس تابع تشخیص تقریباً یکسان بود.

**نتیجه‌گیری:** در این بررسی مشخص گردید که دوره‌ی بیماری دیابت یکی از مهمترین عوامل تعیین کننده‌ی رتینوپاتی می‌باشد؛ به طوری که خطر رتینوپاتی در بیماران دیابت نوع II که از بیماری دیابت آن‌ها بیشتر از ۷/۵ سال می‌گذرد نسبت به سایر بیماران بیشتر است.

**واژگان کلیدی:** درخت رده بندی، تحلیل ممیزی (تابع تشخیصی)، رتینوپاتی، دیابت نوع II.

<sup>۱</sup> استادیار آمار زیستی، گروه آمار و اپیدمیولوژی، دانشکده‌ی بهداشت، دانشگاه علوم پزشکی اصفهان، اصفهان، ایران.

<sup>۲</sup> استاد، مرکز تحقیقات غدد و متابولیسم، دانشگاه علوم پزشکی اصفهان، اصفهان، ایران.

<sup>۳</sup> کارشناس ارشد، گروه آمار و اپیدمیولوژی، دانشکده‌ی بهداشت، دانشگاه علوم پزشکی اصفهان، اصفهان، ایران.

<sup>۴</sup> دانشیار آمار زیستی، گروه پزشکی اجتماعی، دانشکده‌ی پزشکی، دانشگاه علوم پزشکی قزوین، قزوین، ایران.

**نویسنده‌ی مسؤؤل:** دکتر سید محسن حسینی، استادیار آمار زیستی، گروه آمار و اپیدمیولوژی، دانشکده‌ی بهداشت، دانشگاه علوم پزشکی اصفهان، اصفهان، ایران.

E-mail: hosseini@hlth.mui.ac.ir

## Using Classification Tree for prediction of Diabetic Retinopathy on Type II Diabetes

Sayed Mohsen Hosseini MSc PhD<sup>\*</sup>, Mehdi Tazhibi MSc PhD<sup>\*</sup>, Massoud Amini MD<sup>\*\*</sup>,  
Asghar Zaree MSc<sup>\*\*\*</sup>, Hassan Jahani Hashemi MD<sup>\*\*\*\*</sup>

### Abstract

**Background:** Diabetes disease is a common disease in the world and diabetic retinopathy that is an eye disease, involve these patients so hardly that leads patients to blindness in elderly. Because of high incidence of type II diabetes in the society and the danger of retinopathy in this group of patients, in this study we attempted to introduce a classification and predictive model according to Classification and Regression Tree (CART) model for this disease.

**Methods:** This study was performed according to the information of about 3734 patients with type II diabetes, consulted to Isfahan Metabolic and Endocrine Research Center from 1991 to 2006. According to the CART model, a classification pattern was introduced for predicting of retinopathy in these patients.

**Findings:** In this study, classification tree model (CART) obtained with accuracy of 67 percents according to duration of disease, age, sex, systolic blood pressure, triglyceride, cholesterol and fasting blood sugar. This classification model had the sensitivity of 71% and specificity of 62%.

**Conclusion:** By this study founds that the duration of diabetes is one of the most important element of retinopathy, in such a manner that the danger of retinopathy in type II diabetes patients with more than 7.5 years is more than other patients.

**Keywords:** Diabetic retinopathy, Type II diabetes, Classification tree.

<sup>\*</sup> Assistant Professor, Department of Biostatistics and Epidemiology, School of Health, Isfahan University of Medical Sciences, Isfahan, Iran.

<sup>\*\*</sup> Professor, Endocrine and Metabolism Research Center, Isfahan University of Medical Sciences, Isfahan, Iran.

<sup>\*\*\*</sup> Resident MSc in Biostatistics, Department of Biostatistics and Epidemiology, School of Health, Isfahan University of Medical Sciences, Isfahan, Iran.

<sup>\*\*\*\*</sup> Associate Professor, Department of Community Medicine, School of Medicine, Isfahan University of Medical Sciences, Isfahan, Iran.

Corresponding Author: Sayed Mohsen Hosseini MSc, PhD, E-mail: hosseini@hlth.mui.ac.ir

## مقدمه

دیابت یکی از بیماری‌های شایع و مزمن در جهان به شمار می‌رود (۱) و رتینوپاتی به عنوان یکی از شایع‌ترین عوارض میکروواسکولر این گروه از بیماران را به شدت گرفتار می‌نماید (۲)؛ به طوری که خطر کوری در دیابتی‌ها ۲۵ برابر نسبت به سایر بیماران می‌باشد (۳). رتینوپاتی دیابتی یکی از مهمترین عواملی است که باعث ایجاد اختلالات بینایی یا کوری در سنین کار و فعالیت در جوامع پیشرفته می‌شود (۴-۵)؛ به طوری که این بیماری یکی از سه عامل مهم کوری در آمریکا به شمار می‌رود و منجر به کوری در سنین ۲۰ تا ۷۴ سالگی می‌گردد (۶-۷). عوامل مختلفی در بروز رتینوپاتی نقش دارند که می‌توان از طول دوره‌ی بیماری دیابت، سن، وجود پروتئین در ادرار، فشار خون، حاملگی، نوع کنترل قند خون و چربی‌های خون نام برد (۸،۹). با توجه به شیوع بالای دیابت نوع II در جامعه و همچنین خطر بالای رتینوپاتی دیابتی در بیماران مبتلا، در این بررسی سعی گردید تا با توجه به عوامل مؤثر در رتینوپاتی یک الگوی رده بندی درختی برای این بیماران معرفی گردد.

درخت رده بندی (Classification and Regression Tree) در زمینه‌ی تشخیص و پیش‌گویی پزشکی به طور گسترده‌تری در مقایسه با انواع دیگر درخت رده بندی مورد استفاده قرار می‌گیرد (۱۰-۱۲). روش‌های معمول آماری به سه دلیل عمده نمی‌توانند در این زمینه مورد استفاده قرار گیرند؛ اول این که به طور معمول، در اطلاعات پزشکی با مقادیر بسیاری از متغیرهای کمکی سر و کار داریم که این حجم زیاد متغیرها استفاده از روش‌های معمول را با مشکل و پیچیدگی‌هایی روبه‌رو می‌سازد؛ به عنوان نمونه اثرات متقابل در این داده‌ها به قدری پیچیده است که نمی‌توان

با روش‌های معمول آماری آن‌ها را بررسی کرد. دوم این که روش‌های معمول آماری نیازمند وجود بعضی از فرضیات نظیر نرمال بودن توزیع و واریانس‌های برابر می‌باشند که در داده‌های پزشکی اغلب چنین فرضیاتی برقرار نخواهد بود. درخت تصمیم با توجه به استفاده از شیوه‌های ناپارامتری نیازی به داشتن توزیع خاصی برای مشاهدات ندارد (۱۳-۱۵، ۱۰). سوم این که نتایج به دست آمده از روش‌های معمول آماری به سادگی قابل استفاده نیست. درخت تصمیم یکی از شیوه‌های مناسب رده‌بندی می‌باشد که با وجود مشکلات ذکر شده می‌تواند به خوبی عمل نماید (۱۶). به علاوه درخت تصمیم نسبت به داده‌های پرت حساس نمی‌باشد و برای اطلاعات با حجم زیاد بسیار مناسب است (۱۷-۱۸، ۱۳-۱۵).

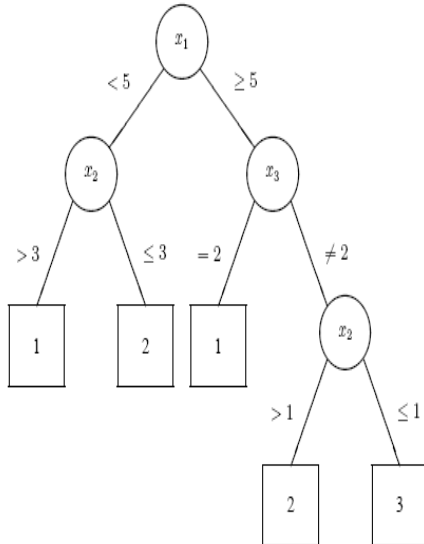
## روش‌ها

تحقیقات دیابت اصفهان مراجعه کرده بودند، انجام گرفت. از بین این بیماران، کسانی که وضعیت رتینوپاتی مشخص شده داشتند، برای مطالعه انتخاب شدند؛ تعداد این بیماران ۳۷۳۴ نفر بود. برای این بیماران اطلاعاتی نظیر سن، جنسیت، مدت بیماری، شاخص توده‌ی بدنی (BMI)، وضعیت سیگاری بودن، میزان قند خون ناشتا، سطح کلسترول، تری‌گلیسیرید خون، فشار خون سیستولیک و دیاستولیک، وجود پروتئین در ادرار و هموگلوبین A<sub>1c</sub> در نظر گرفته شد.

در این مطالعه برای معرفی یک الگوی رده بندی از بررسی آماری CART به کمک نرم‌افزار STATISTICA<sup>6</sup> استفاده گردید. CART یکی از انواع درخت رده بندی می‌باشد که توسط Breiman و همکاران در سال ۱۹۸۴ میلادی معرفی گردید. این مدل

$$i(t) = gini(t) = 1 - \sum_{j=1}^k p^2 [c = c_j | T]$$

$$= \sum_{k \neq l} p(c = c_k | T) \cdot p(c = c_l | T)$$



نمودار ۱. نمونه‌ای از درخت رده بندی CART

این رابطه هنگامی که مشاهدات فقط متعلق به یک رده باشند، برابر صفر است و زمانی که احتمال هر رده برابر باشد، بیشترین مقدار ممکن را اختیار می‌نماید. با در نظر گرفتن متغیر کمکی  $x$  که بر اساس آن گره‌ی  $t$  به  $n$  زیر شاخه تقسیم می‌گردد (هر زیر شاخه با  $T_j$  برای  $j$  از ۱ تا  $n$  نشان داده می‌شود)، یک کاهش در تابع ناچوری خواهیم داشت که بر اساس شاخص جینی به صورت زیر تعریف می‌شود:

$$gini\ gain = GG(T, X) = gini(t) - \sum_{j=1}^n p(T_j | T) \cdot gini(T_j)$$

از بین چندین متغیر، متغیری مناسب است که مقدار بیشتری را برای  $GG(T, X)$  اختیار نماید؛ این ملاکی برای انتخاب بهترین متغیر از بین چندین متغیر می‌باشد (۱۳، ۲۵، ۱۱-۱۰). از این رو با توجه به تابع ناچوری و شاخص جینی ابتدا مقدار تابع ناچوری در حالت کلی برای متغیر پاسخ محاسبه می‌گردد. در مرحله‌ی بعد برای

یک گراف غیر چرخشی شبیه درخت با تقسیمات دوتایی بر اساس متغیرهای کمکی را برای معرفی یک الگوی رده بندی و تشخیصی معرفی می‌نماید (۱۱، ۱۳، ۱۹-۲۲).

یک درخت تصمیم از سه جز اصلی شامل ریشه، گره‌ی داخلی و برگ تشکیل شده و روند بدین گونه است که ابتدا یک متغیر کمکی به عنوان ریشه انتخاب و با توجه به اهداف مطالعه به چندین گره‌ی داخلی تقسیم می‌شود. هر گره‌ی داخلی نیز مانند ریشه به گره‌های دیگری تقسیم می‌شود تا در نهایت به هر گره یک رده از متغیر پاسخ منتسب گردد؛ این گره‌ها برگ نامیده می‌شود (۱۰، ۱۷، ۲۳-۲۴). نمودار ۱ نمونه‌ای از درخت رده بندی CART می‌باشد.

در این بررسی، مشاهدات به صورت تصادفی به دو گروه نمونه‌ی آزمون و یادگیری تقسیم شدند. نمونه‌ی یادگیری (فراگیر) زیرمجموعه‌ای از مشاهدات می‌باشد که درخت رده بندی بر اساس آن طراحی می‌گردد؛ نمونه‌ی آزمون نیز زیرمجموعه‌ای از داده‌هاست که به منظور تعیین دقت درخت طراحی شده با نمونه‌ی یادگیری، مورد استفاده قرار می‌گیرد (۱۶-۱۷، ۲۴).

به منظور انتخاب متغیرهای مهم در الگوی رده بندی درختی، در این بررسی از تابعی تحت عنوان تابع ناچوری (Impurity Function) و شاخصی به نام جینی (Gini) استفاده گردید. تابع ناچوری برای گره‌ای مانند  $t$  و متغیر وابسته با  $k$  رده  $(c_1, \dots, c_k)$  به صورت زیر تعریف می‌شود:

$$i(t) = \Phi[p(c = c_1 | t), \dots, p(c = c_k | t)]$$

شاخص جینی اغلب در مدل‌های درختی با تقسیمات دوتایی در هر گره مورد استفاده قرار می‌گیرد و به صورت زیر تعریف می‌شود:

بندی اشتباه، حساسیت و ویژگی الگوی معرفی بر مبنای تشخیص بالینی برای دو نمونه‌ی آزمون و یادگیری محاسبه و با استفاده از آزمون  $\chi^2$  مقایسه گردید. به علاوه با روشی مشابه، نتایج حاصل از الگوی درختی معرفی شده با نتایج به دست آمده از روش معمول آماری تابع تشخیص نیز مقایسه گردید.

#### یافته‌ها

در بین ۳۷۳۴ بیمار تحت بررسی، ۲۰۰۱ بیمار (۵۳/۶ درصد) دارای رتینوپاتی و ۱۷۳۳ بیمار دیگر (۴۶/۴ درصد) بدون رتینوپاتی بودند. میانگین سنی این بیماران برابر  $10/44 \pm 52/15$  سال و میانگین دوره‌ی بیماری دیابت  $6/01 \pm 7/04$  سال بود. ۶۴/۳ درصد بیماران را زنان تشکیل می‌دادند و از بین این بیماران تنها ۱۲/۳ درصد سیگاری بودند. جدول ۱ وضعیت بیماران را برای سن، دوره‌ی بیماری، شاخص توده‌ی بدنی، میزان قند خون، کلسترول و تری‌گلیسیرید خون، فشار خون سیستولی و دیاستولی و هموگلوبین A<sub>1c</sub> را در کل بیماران و در دو گروه دارای رتینوپاتی و بدون رتینوپاتی نشان می‌دهد. با توجه به جدول ۱، به جز تری‌گلیسیرید و کلسترول خون، از نظر آماری بین میانگین خصوصیات کمی بررسی شده در دو گروه دارای رتینوپاتی و بدون رتینوپاتی تفاوت آماری معنی‌دار وجود داشت. جدول ۲ فراوانی و درصد فراوانی بیماران را در دو گروه دارای رتینوپاتی و بدون رتینوپاتی با توجه به جنسیت، سیگاری بودن و وجود پروتئین در ادرار نشان می‌دهد. با توجه به جدول ۲، دو گروه دارای رتینوپاتی و بدون رتینوپاتی از نظر درصد فراوانی جنسیت، سیگاری بودن و وجود پروتئین در ادرار تفاوت معنی‌دار داشتند.

تمام متغیرهای کمکی، با توجه به بهترین تقسیمات دوتایی برای متغیر پاسخ، مقدار تابع ناجوری در هر یک از دو زیر مجموعه‌ی ایجاد شده محاسبه و میانگین وزنی آن‌ها از مقدار تابع ناجوری کل کم می‌گردد.

از بین متغیرهای کمکی، متغیری که دارای بیشترین مقدار برای این رابطه باشد، در گام اول رده بندی درختی انتخاب می‌شود. در برخورد با متغیرهای کمی، از تقسیمات دوتایی استفاده نموده، نقطه‌ای مانند a (نقطه‌ی برش) را تعیین می‌کنیم. لازم به ذکر است که نقطه‌ی برش در بسیاری از الگوهای رده بندی درختی توسط خود شاخص به کار برده شده (در این جا شاخص جینی)، مشخص می‌شود. در برخورد با متغیر کیفی، هر سطح متغیر به عنوان یک زیر شاخه‌ی درخت رده بندی در نظر گرفته می‌شود (۲۹-۲۶، ۱۲).

در مدل CART برای انتخاب اندازه‌ی مناسب از درخت رده بندی، از روشی تحت عنوان ارزش پیچیدگی (Cost complexity) استفاده می‌شود. یک الگوی درختی زمانی مناسب است که علاوه بر این که برای مشاهدات موجود (نمونه‌ی فراگیری) خوب عمل می‌نماید، برای مشاهدات جدید (نمونه‌ی آزمون) نیز مناسب باشد. با وجود این که با بزرگ‌تر شدن اندازه‌ی درخت رده بندی، دقت رده بندی در مورد نمونه‌ی یادگیری افزایش می‌یابد، این دقت برای نمونه‌ی آزمون از اندازه‌ای به بعد کاهش می‌یابد. روش ذکر شده در واقع تعادلی بین دقت درخت رده بندی و اندازه‌ی آن برقرار می‌نماید و با توجه به اندازه‌ی خطا و تعداد گره‌های درخت، تصمیم‌گیری می‌نماید که کدام گره از درخت را باید حذف نمود (۲۲-۱۹، ۱۳، ۱۱).

در این مطالعه بعد از تعیین الگوی درختی بر اساس مدل CART، برای تعیین مناسب‌ترین مدل، میزان گروه

با توجه به الگوی رده بندی درختی CART، چهار الگوی درختی با معیارهای مختلف برای ساخت درخت رده بندی معرفی گردید. به کمک میزان گروه‌بندی اشتباه، حساسیت و ویژگی مدل‌های معرفی شده، مناسب‌ترین مدل با توجه به آزمون  $\chi^2$  انتخاب گردید. مدل رده بندی درختی به دست آمده مطابق نمودار ۲ بود.

جدول ۱. مقایسه‌ی مقادیر کمی خصوصیات بررسی شده در کل بیماران و دو گروه دارای رتینوپاتی و بدون رتینوپاتی

| P-value | بدون رتینوپاتی |         | دارای رتینوپاتی |         | کل بیماران |         | متغیر                     |
|---------|----------------|---------|-----------------|---------|------------|---------|---------------------------|
|         | SD             | میانگین | SD              | میانگین | SD         | میانگین |                           |
| < ۰/۰۰۱ | ۱۰/۵۱          | ۵۰/۱۱   | ۱۰/۰۵           | ۵۳/۹۱   | ۱۰/۴۴      | ۵۲/۱۵   | سن                        |
| < ۰/۰۰۱ | ۴/۷            | ۴/۹۷    | ۶/۴۷            | ۸/۸۳    | ۶/۰۱       | ۷/۰۴    | دوره‌ی بیماری             |
| < ۰/۰۰۱ | ۴/۶۶           | ۲۷/۹۶   | ۴/۲۳            | ۲۶/۶۷   | ۴/۴۸       | ۲۷/۲۷   | شاخص توده‌ی بدنی          |
| < ۰/۰۰۱ | ۴۲/۶۵          | ۱۶۹/۳۲  | ۴۶/۵۳           | ۱۸۱/۱۴  | ۴۵/۱۳      | ۱۷۵/۵۹  | قند خون ناشتا             |
| ۰/۹۴۲   | ۳۸/۱۷          | ۲۱۸/۶۲  | ۳۸/۷۷           | ۲۱۸/۵۲  | ۳۸/۴۸      | ۲۱۸/۵۷  | کلسترول                   |
| ۰/۸۳۳   | ۱۱۲/۰۵         | ۲۱۳/۵۹  | ۱۰۴/۸۴          | ۲۱۲/۳۵  | ۱۰۸/۲۸     | ۲۱۲/۹۳  | تری‌گلیسرید               |
| < ۰/۰۰۱ | ۱۵/۴۶          | ۱۲۶/۶۳  | ۱۶/۰۷           | ۱۳۱/۸۶  | ۱۶         | ۱۲۹/۴۵  | فشار خون سیستولی          |
| ۰/۰۵۱   | ۷/۵۸           | ۸۰/۲۴   | ۷/۵۴            | ۸۰/۷۴   | ۷/۵۶       | ۸۰/۵۱   | فشار خون دیاستولی         |
| < ۰/۰۰۱ | ۱/۸۸           | ۸/۶۹    | ۱/۹۸            | ۹/۰۷    | ۱/۹۴       | ۸/۸۹    | هموگلوبین A <sub>1c</sub> |

جدول ۲. فراوانی و درصد فراوانی در دو گروه رتینوپاتی و بدون رتینوپاتی برای جنسیت، سیگاری بودن و وجود پروتئین در ادرار

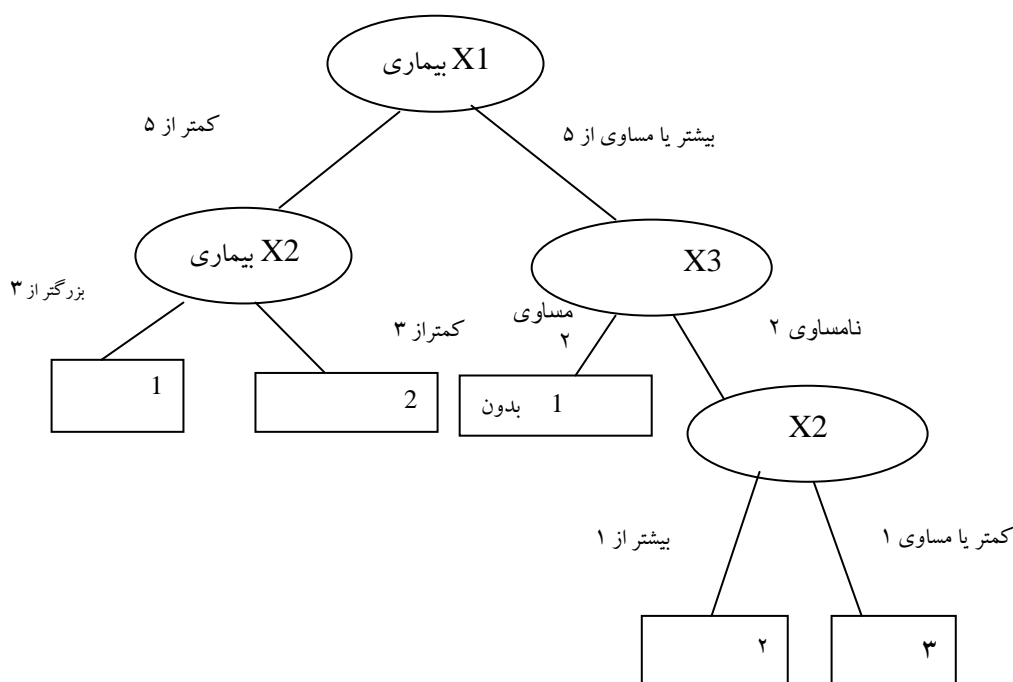
| P-value | بدون رتینوپاتی |      | دارای رتینوپاتی |      | کل بیماران |      |            |                       |
|---------|----------------|------|-----------------|------|------------|------|------------|-----------------------|
|         | فراوانی        | درصد | فراوانی         | درصد | فراوانی    | درصد |            |                       |
| < ۰/۰۰۱ | ۴۱/۴           | ۵۵۱  | ۵۸/۶            | ۷۷۹  | ۳۵/۷       | ۱۳۳۰ | مرد        | جنسیت                 |
|         | ۴۹/۲           | ۱۱۸۰ | ۵۰/۸            | ۱۲۱۹ | ۶۴/۳       | ۲۳۹۹ | زن         |                       |
| ۰/۰۲۳   | ۴۰/۷           | ۱۶۸  | ۵۹/۳            | ۲۴۵  | ۱۲/۲       | ۴۱۳  | سیگاری     | وضعیت سیگاری          |
|         | ۴۶/۶           | ۱۳۸۳ | ۵۳/۴            | ۱۵۸۵ | ۸۷/۸       | ۲۹۶۷ | غیر سیگاری |                       |
| < ۰/۰۰۱ | ۲۹/۵           | ۲۶۳  | ۷۰/۵            | ۶۲۹  | ۳۰/۹       | ۸۹۲  | دارد       | وجود پروتئین در ادرار |
|         | ۵۰/۸           | ۱۰۱۲ | ۴۹/۲            | ۹۸۱  | ۶۹/۱       | ۱۹۹۳ | ندارد      |                       |

جدول ۳. مقایسه‌ی گروه‌بندی اشتباه، حساسیت و ویژگی در نمونه‌ی آزمون و یادگیری

| P-value | نمونه‌ی آزمون | نمونه‌ی فراگیر |                  |
|---------|---------------|----------------|------------------|
| ۰/۳۴۲   | ۰/۳۱          | ۰/۳۳           | گروه‌بندی اشتباه |
| ۰/۷۲۹   | ۰/۷۲          | ۰/۷۱           | حساسیت           |
| ۰/۲۷۶   | ۰/۶۵          | ۰/۶۲           | ویژگی            |

حساسیت و ویژگی در دو نمونه‌ی آزمون و یادگیری تفاوت معنی‌داری از نظر آماری وجود نداشت. از این رو می‌توان نتیجه گرفت که الگوی درختی معرفی شده برای مشاهدات جدید می‌تواند با دقتی در حدود ۶۷ درصد تصمیم‌گیری نماید که این تصمیم‌گیری از حساسیت ۷۱ و ویژگی ۶۲ درصد برخوردار می‌باشد.

بر اساس الگوی رده‌بندی درختی معرفی شده، بیماران به نه رده‌ی متمایز دسته بندی می‌شوند. با توجه به این مدل درختی، دوره‌ی بیماری دیابت مهمترین نقش را برای تعیین رتینوپاتی دیابتی ایفا می‌کند. جدول ۳ معیارهای صحت الگوی رده بندی درختی را برای نمونه‌ی فراگیر و آزمون نشان می‌دهد. این جدول نشان می‌دهد که بین مقادیر گروه‌بندی اشتباه،



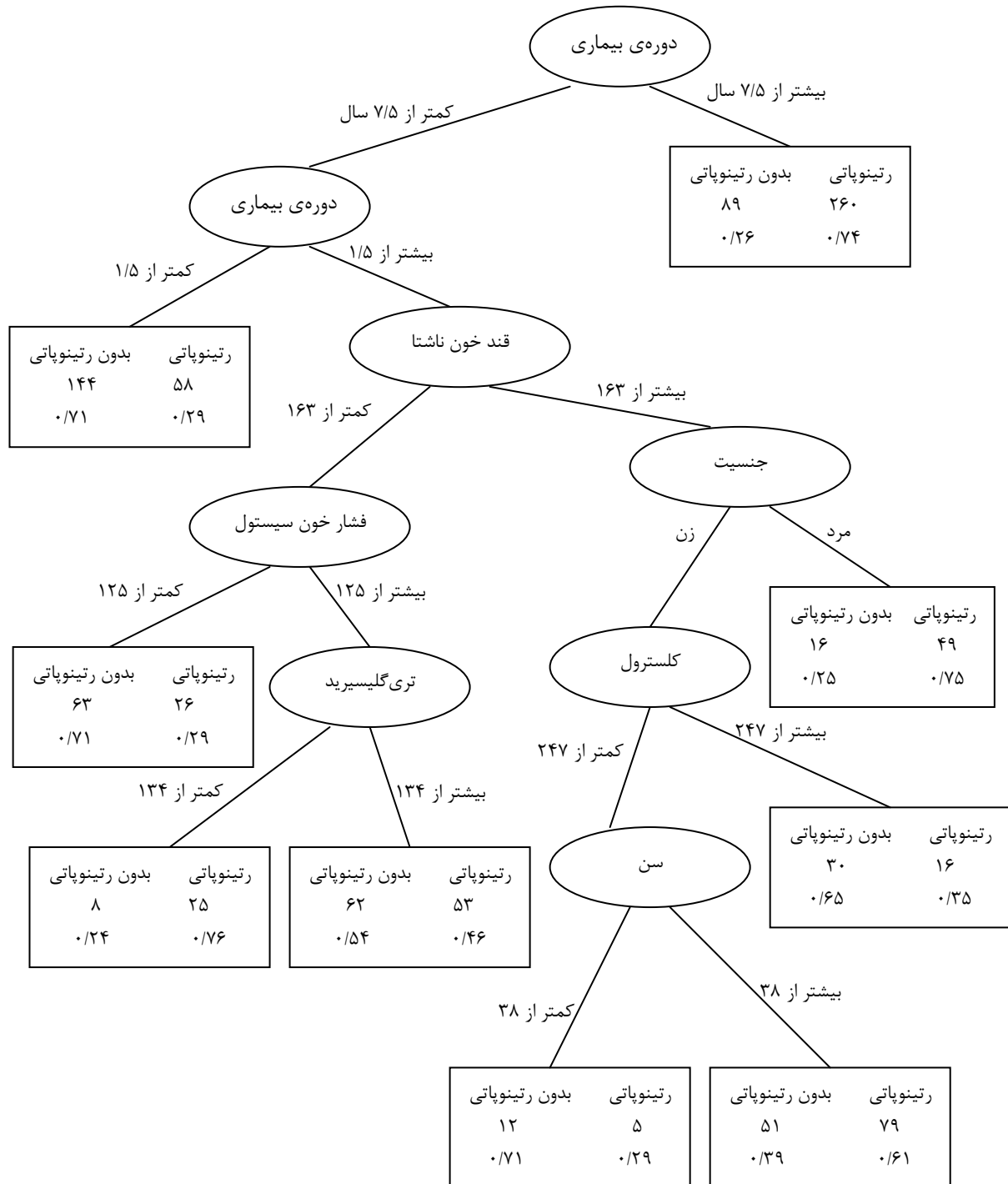
نمودار ۱: نمونه ای از درخت رده بندی CART

جدول ۴. ضرایب استاندارد شده‌ی متغیرهای مهم در تابع تشخیص

| متغیر                 | ضریب استاندارد شده |
|-----------------------|--------------------|
| جنسیت                 | -۰/۱۴۴             |
| دوره‌ی بیماری         | ۰/۶۰۸              |
| شاخص توده‌ی بدنی      | -۰/۲۲۲             |
| قند خون ناشتا         | ۰/۲۴۷              |
| فشار خون سیستولی      | ۰/۶۱۵              |
| فشار خون دیاستولی     | -۰/۳۸۴             |
| وجود پروتئین در ادرار | ۰/۳۱۹              |

بدنی، میزان قند خون ناشتا، فشار خون سیستولی و دیاستولی و وجود پروتئین در ادرار، یک الگوی رده‌بندی خطی با ضرایب استاندارد شده مطابق جدول ۴ را معرفی نمود.

در این بررسی برای اطمینان از کارایی مدل درختی، الگوی رده بندی بر اساس روش کلاسیک آماری تابع تشخیص نیز به دست آمد. تابع تشخیص با توجه به متغیرهای جنسیت، دوره‌ی بیماری، شاخص توده‌ی



نمودار ۲. درخت رده‌بندی با اندازه‌ی نه بر اساس نمونه‌ی فراگیر و آزمون

**بحث**

با توجه به یکسان بودن میزان گروه‌بندی اشتباه در نمونه‌ی آزمون و یادگیری می‌توان نتیجه گرفت که مدل درختی معرفی شده می‌تواند برای مشاهدات جدید با دقتی معادل ۶۷ درصد تصمیم‌گیری نماید. بر اساس مدل معرفی شده، دوره‌ی بیماری دیابت به عنوان مهم‌ترین عامل برای رده بندی شناخته شد؛ به طوری که بیماران دیابتی نوع II که از سابقه‌ی بیماری آن‌ها بیشتر از ۷/۵ سال می‌گذرد، نسبت به گروه‌های دیگر، در معرض خطر بیشتری برای رتینوپاتی هستند. بسیاری از مطالعات دیگر نیز رابطه‌ی مستقیمی بین دوره‌ی بیماری دیابت و شیوه‌ی رتینوپاتی نشان داده و دوره‌ی بیماری به عنوان یکی از مهم‌ترین عوامل تأثیرگذار در رتینوپاتی شناخته شده است (۳۳-۳۰، ۶). با این حال، الگوی معرفی شده توسط درخت تصمیم، یک نمودار قابل درک و ساده برای رده بندی معرفی می‌نماید (۳۵-۳۴). به علاوه در الگوهای درخت تصمیم، نیازی به در نظر گرفتن اثرات متقابل نیست. در این الگوها اگر اثر متقابل وجود داشته باشد، در درخت تصمیم بدون در نظر گرفتن ظاهر می‌گردد. همچنین شدت عوامل مؤثر بر ابتلا به بیماری را تعیین می‌کند (۱۶، ۱۰).

با توجه به مطالعه‌ی انجام شده در آمریکا، با گروه‌بندی بر روی دوره‌ی بیماری تفاوتی بین شیوع رتینوپاتی در دیابت نوع II و نوع I دیده نشد (۳۶).

الگوی درختی معرفی شده در این بررسی، فشار خون سیستمیک را نیز به عنوان یک عامل برای رده بندی معرفی می‌نماید که تأثیرات فشار خون، به خصوص فشار خون سیستمیک، بر روی رتینوپاتی در چندین مطالعه قید شده است و فشار خون به عنوان یک ریسک فاکتور رتینوپاتی شناخته می‌شود

(۳۹-۳۷، ۹، ۶). بر اساس نتایج به دست آمده از مدل درختی، بیماران با سطح پایین تری‌گلیسرید و کلسترول نسبت به گروه دیگر بیماران در خطر بیشتری قرار دارند؛ این نکته در بعضی از مطالعات ذکر شده ولی اختلاف معنی‌داری بین دو گروه دیده نشده است. با توجه به ارتباط مستقیم شاخص توده‌ی بدنی با این دو عامل، نتایج به دست آمده در این مطالعه با نتایج دیگر مطالعات (۳۸-۳۷، ۹، ۵) همخوانی دارد.

مدل درختی معرفی شده در این بررسی، حاکی از اختلاف شیوع رتینوپاتی در مدل تابع تشخیص معرفی شده است. این مدل به منظور رده بندی رتینوپاتی از میزان دقت ۰/۶۷ با حساسیت ۰/۷۱ و ویژگی ۰/۶۴ برخوردار می‌باشد.

باید توجه داشت که در مقایسه‌ی مدل درخت رده بندی و تابع تشخیص، هدف مقایسه‌ی تک تک متغیرها نمی‌باشد؛ بلکه همان طور که اشاره شد، برای اطمینان از کارایی مدل درختی، میزان دقت درخت رده‌بندی و تابع تشخیص مقایسه گردید که از نظر آماری اختلاف معنی‌داری بین آن‌ها دیده نشد.

مدل درختی معرفی شده در این بررسی، حاکی از اختلاف شیوع رتینوپاتی بین مردان و زنان بود. با این که در مطالعات مختلف شیوع بالاتر رتینوپاتی در مردان نسبت به زنان عنوان شده ولی در اکثر مطالعات جنسیت به عنوان یک عامل تأثیرگذار شناخته نشده است (۳۷-۳۶، ۹، ۶، ۳).

**نتیجه‌گیری**

در این بررسی مشخص گردید که دوره‌ی بیماری دیابت یکی از مهم‌ترین عوامل تعیین کننده‌ی رتینوپاتی می‌باشد؛ به طوری که خطر رتینوپاتی در بیماران دیابت



داشت که تعریف دوره‌ی بیماری، مدت زمان تشخیص دیابت بیمار تا شروع مطالعه بوده و در آن، دوره‌ی کمون لحاظ نشده است. همچنین یکی دیگر از محدودیت‌های مطالعه‌ی حاضر آن بود که درخت‌های رده‌بندی اغلب در هر مرحله از تقسیم، فقط از یک متغیر بهره می‌برند (۱۴-۱۳).

نوع II، که از بیماری آن‌ها بیشتر از ۷/۵ سال می‌گذرد، نسبت به سایر بیماران بیشتر می‌باشد. همچنین بر اساس مدل CART می‌توان به خوبی مدل تابع تشخیص، در مورد بیماری رتینوپاتی در بیماران دیابت نوع II تصمیم‌گیری نمود.

هر چند در این مطالعه دوره‌ی بیماری به عنوان مهم‌ترین عامل رتینوپاتی شناخته شد اما بایستی توجه

## References

- Vrijhoef HJ, Diederiks JP, Spreuwenberg C. Effects on quality of care for patients with NIDDM or COPD when the specialised nurse has a central role: a literature review. *Patient Educ Couns* 2000; 41(3): 243-50.
- Rodriguez J, Sanchez R, Munoz B, West SK, Broman A, Snyder RW, et al. Causes of blindness and visual impairment in a population-based sample of U.S. Hispanics. *Ophthalmology* 2002; 109(4): 737-43.
- Jamal-u-Din, Qureshi MB, Khan AJ, Khan MD, Ahmad K. Prevalence of diabetic retinopathy among individuals screened positive for diabetes in five community-based eye camps in northern Karachi, Pakistan. *J Ayub Med Coll Abbottabad* 2006; 18(3): 40-3.
- Gupta S, Ambade A. Prevalence of diabetic retinopathy and influencing factors amongst type 2 diabetics from central India. *Int J Diab Dev Ctries* 2004; 24:75-8.
- Rema M, Deepa R, Mohan V. Prevalence of retinopathy at diagnosis among type 2 diabetic patients attending a diabetic centre in South India. *Br J Ophthalmol* 2000; 84(9): 1058-60.
- Brown JB, Pedula KL, Summers KH. Diabetic retinopathy: contemporary prevalence in a well-controlled population. *Diabetes Care* 2003; 26(9): 2637-42.
- Fonseca V, Munshi M, Merin LM, Bradford JD. Diabetic retinopathy: a review for the primary care physician. *South Med J* 1996; 89(9): 839-50.
- Marshall G, Garg SK, Jackson WE, Holmes DL, Chase HP. Factors influencing the onset and progression of diabetic retinopathy in subjects with insulin-dependent diabetes mellitus. *Ophthalmology* 1993; 100(8): 1133-9.
- Harney F. Diabetic retinopathy. *Medicine-Abingdon* 2006; 34(3): 95-8.
- Podgorelec V, Kokol P, Stiglic B, Rozman I. Decision trees: an overview and their use in medicine. *J Med Syst* 2002; 26(5): 445-63.
- Breiman L, Friedman J, Stone CJ, Olshen RA. Classification and regression trees. 1<sup>st</sup> ed. London: Chapman and Hall/CRC; 1984.
- Marshall RJ. The use of classification and regression trees in clinical epidemiology. *J Clin Epidemiol* 2001; 54(6): 603-9.
- Timofeev R. Classification and patients attending a diabetic center Regression Trees (CART) Theory and application. [Thesis]. Berlin, Humboldt Uni; 2004
- Tan PN, Steinbach M, Kumar V. Introduction to data mining. Canada: Pearson Addison Wesley; 2006. p. 145-206.
- Yohannes Y, Webb P. Classification and regression trees, CART: a user manual for identifying indicators of vulnerability to famine and chronic food insecurity. Washington, DC: International Food Policy Research Institute; 1999.
- Lewis RJ. An Introduction to Classification and Regression Tree (CART) Analysis, presented at Annual Meeting of the Society for Academic Emergency Medicine. Annual Meeting of the Society of Academic Emergency Medicine. 2000. Available from: URL: [http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.95.4103]
- Rema M, Premkumar S, Anitha B, Deepa R, Pradeepa R, Mohan V. Prevalence of diabetic retinopathy in urban India: the Chennai Urban Rural Epidemiology Study (CURES) eye study, I. *Invest Ophthalmol Vis Sci* 2005; 46(7): 2328-33.
- Safavian SR, Landgrebe D. A survey of decision tree classifier methodology. *IEEE Trans Systems, Man, & Cybernetics* 1991; 21(3): 660-7.
- Mitra S, Acharya T. Data Mining: multimedia, soft computing, and bioinformatics. London: John Wiley and Sons; 2003.
- Wilkinson L. Tree structured data analysis: AID, CHAID, and CART. *Proceedings of Sawtooth*

- Software Conference; 1992; Sun Valley.
21. Kim H, Loh WY. Classification trees with unbiased multiway split. *J Amer Statist Assoc* 2001; 98: 598-604.
22. Scoh C, Willett R, Nowak R. *CORT: Classification and regression Trees*. Proceedings of the 6<sup>th</sup> International Conference of IEEE; April 2003.
23. Utgoff PE. Decision trees. In: Wilson RA, Keil FC, Editors. *The MIT encyclopedia of the cognitive sciences*. Bradford: MIT Press; 1998.
24. Quinlan JR. *C4.5: programs for machine learning*. San Francisco: Morgan Kaufmann Publishers Inc; 1994. p. 235-40.
25. Buntine W, Niblett T. A further comparison of splitting rules for decision-tree induction. *Machine Learning* 1992; 8(1): 75-85.
26. Quinlan JR. Induction of decision trees. *Machine Learning* 1986; 1(1): 81-106.
27. Shih YS. Selecting the best splits for classification trees with categorical variables. *Statistics & Probability Letters* 2001; 54(4): 341-5.
28. Loh WY, Shih YS. Split selection methods for classification trees. *Statistica Sinica* 1997; 7: 815-40.
29. Hoare R. Using CHAID for classification problems. *Proceedings of New Zealand Statistical Association Conference on machine learning*; 1997. New Zealand.
30. Dobra AV. *Scalable classification and regression tree construction*. New York: Cornell University; 2003.
31. Kull CE, Abrahasson M. Prevalence of retinopathy difference in a population with type 1 diabetes. *Diabet Med* 2002; 19(11): 923-4.
32. Porta M, Bandello F. Diabetic retinopathy: A clinical update. *Diabetologia* 2002; 45(12): 1617-34.
33. Dandona L, Dandona R, Naduvilath TJ, McCarty CA, Rao GN. Population based assessment of diabetic retinopathy in an urban population in southern India. *Br J Ophthalmol* 1999; 83(8): 937-40.
34. Frank E, Wang Y, Inglis S, Holmes G, Witten IH. Using model trees for classification. *Machine Learning* 1998; 32(1): 63-76.
35. Frank E. *Pruning decision trees and lists*. Newzealand: Waikato Univ; 2000.
36. Varma R, Torres M, Pena F, Klein R, Azen SP. Prevalence of diabetic retinopathy in adult Latinos: the Los Angeles Latino eye study. *Ophthalmology* 2004; 111(7): 1298-306.
37. Janghorbani M, Amini M, Ghanbari H, Safaiee H. Incidence of and risk factors for diabetic retinopathy in Isfahan, Iran. *Ophthalmic Epidemiol* 2003; 10(2): 81-95.
38. Abdollahi A, Malekmadani MH, Mansoori MR, Bostak A, Abbaszadeh MR, Mirshahi A. Inflammatory markers and diabetic retinopathy in type 1 diabetes. *Acta Medica Iranica*, 2006; 44(6): 415-9.
39. Izuora KE, Chase HP, Jackson WE, Coll JR, Osberg IM, Gottlieb PA, et al. Inflammatory markers and diabetic retinopathy in type 1 diabetes. *Diabetes Care* 2005; 28(3): 714-5.