

# پیش‌بینی و تعیین عوامل مؤثر بر بقای پنج ساله کلیه پیوندی در داده‌های نامتوازن با رویکرد فراابتکاری و یادگیری ماشین

نسیبه امامی\* و زینب حسنی

گروه علوم کامپیوتر، دانشکده علوم پایه، دانشگاه کوثر بجنورد، ایران

## چکیده

در مرحله نهایی نارسایی کلیه، پیوند کلیه می‌تواند عمر بیماران را طولانی کند و کیفیت زندگی بیمار را بسیار بهبود بخشد. بعد از عمل پیوند کلیه، بورسی میزان یا پیش‌بینی بقای کلیه پیوندی اهمیت زیادی دارد. این مطالعه بر روی بیماران کلیه پیوندی بیمارستان‌های امام رضا(ع) و چهارمین شهید محراب کرمانشاه در سال‌های ۲۰۱۲-۲۰۰۱ انجام شده است. از آن جایی که داده‌های نامتوازن باعث ناکارامدی مدل‌های یادگیری ماشین می‌شوند، ابتدا داده‌های نامتوازن با دو روش پیش‌نمونه‌برداری و زیرنمونه‌برداری متوازن شدند؛ سپس عوامل اثرگذار بر بقای پیوند کلیه به کمک الگوریتم فراابتکاری ژنتیک شناسایی شده و مدل یادگیر طبقه‌بند نزدیک ترین همسایه برای پیش‌بینی بقای پنج ساله کلیه پیوندی به کار گرفته شد. بقای کلیه پیوندی در روش پیش‌نمونه‌برداری با دقیقت ۹۶/۸ درصد و زیرنمونه‌برداری با دقیقت ۸۹/۲ درصد پیش‌بینی شد. هم‌چنین، ویژگی‌های وزن، سن دهنده و غیرنده، اوره قبل پیوند، کراتین قبل پیوند، هموگلوبین قبل و بعد پیوند، جنسیت دهنده، RH دهنده و غیرنده، بیماری اولیه، سن دهنده بالای سی و سن گیرنده بالای چهل، به عنوان ویژگی‌های تأثیرگذار در بقای کلیه پیوندی شناسایی شد. مقایسه نتایج به دست آمده از این پژوهش با مطالعات پیشین، برتری مدل پیشنهادی را از نقطه نظر دقیقت دقت مدل نشان می‌دهد. به عبارتی متوازن‌سازی داده‌ها همراه با انتخاب ویژگی بهینه منجر به ارائه مدل پیش‌بینی دقیق تری می‌شود.

وازگان کلیدی: پیوند کلیه، داده‌های نامتوازن، الگوریتم ژنتیک، نزدیک ترین همسایگی.

## Prediction and determining the effective factors on the survival transplanted kidney for five-year in imbalanced data by the meta-heuristic approach and machine learning

Nasibeh Emami\* & Zeinab Hassani

Department of Basic Science, Kosar University of Bojnord, Iran

### Abstract

Chronic kidney failure is one of the most widespread diseases in Iran and the world. In general, the disease is common in high health indexes societies due to increased longevity. Treatment for chronic kidney failure is dialysis and kidney transplantation. Kidney transplantation is an appropriate and effective strategy for patients with End-Stage Renal Disease (ESRD), and it provides a better life and reduces mortality risk for patients. In contrast to many benefits that kidney transplantation has in terms of improving physical and mental health and the life's quality in kidney transplantation patients, it may be rejected because of host's immune response to the received kidney, and it consequences the need for another transplantation, or even death will have to. In fact, a patient that can survive for years with dialysis, he may lose his life with an inappropriate transplantation or be forced into high-risk surgical procedures.

According to the above, the study of predicting the survival of kidney transplantation, its effective factors and providing a model for purposing of high prediction accuracy is essential. Studies in the field of survival

\*Corresponding author • تاریخ ارسال مقاله: ۱۳۹۶/۱۰/۴ • تاریخ آخرین بازنگری: ۱۳۹۷/۵/۱۶ • تاریخ پذیرش: ۱۳۹۷/۷/۱۴ <sup>\*</sup>نویسنده عهده‌دار مکاتبات

of kidney transplantation include statistical studies, artificial intelligence and machine learning. In all of the studies in this field, researchers have sought to identify a more effective set of features in survival of transplantation and the design of predictive models with higher accuracy and lower error rate.

This study carried out on 756 kidney transplant patients with 21 features of Imam Reza and Fourth Shahid Merab hospital in Kermanshah from 2001 to 2012. Some features set to binary value and other features have real continuous values. Due to data are unbalance, which led to convergence of classification model to majority class, so over sampling and under sampling techniques has been used for achieving higher accuracy.

To identify the more effective features on the survival of the kidney transplantation, the genetic meta-heuristic algorithm is used. For this purpose binary coding for each chromosome has been used; it is combining three single-point, two-point, and uniform operators to make better generations, better convergence and achieve higher accuracy rate. The genetic search algorithm plays a vital role in searching for such a space in a reasonable time because data search space is exponential. In fact, in balanced data, genetic algorithm determines the effective factors and the K-nearest neighbor model with precision of classification as the evaluator function was used to predict the five-year survival of the kidney transplantation. Based on the results of this study, in comparison to similar studies for prediction of survival transplanted kidney, the five-year survival rate of transplanted kidney was appropriate in these models. Also the effective factors in over sampling and under sampling methods with a precision of 96.8% and 89.2% are obtained respectively. in addition weight, donor and recipient age, pre-transplantation urea, pre-transplantation creatinine, hemoglobin before and after transplantation, donor gender, donor and recipient RH, primary illness, donor age up 30 and receipt age up 40 were identified as the effective features on kidney transplantation survival. Comparing the results of this study with previous studies shows the superiority of the proposed model from the point of view of the models' precision. In particular, balancing the data along the selection of optimal features leads to a high precision predictive model.

**Keywords:** Kidney Transplantation, imbalance data, Genetic Algorithm, K- nearest neighbors.

عوامل مؤثر بر آن و ارائه یک مدل بهمنظور دقّت پیش‌بینی بالاتر ضرورت می‌یابد و در مراکز متعدد مطالعات مرتبط با بقای کلیه پیوندی انجام گرفته است [3,1]. مطالعات انجام شده جهت بررسی بقای پیوند از نوع مطالعات آماری و هوش مصنوعی است [1]. تجزیه و تحلیل بقا را با استفاده از استانداردهای آماری، می‌توان مدل مبتنی بر جمعیت در نظر گرفت که پیش‌بینی‌ها براساس احتمال و یا فاصله از تخمین‌های جمعیت مشتق می‌شود. یادگیری ماشین و داده‌کاوی، ابزارهای تصمیم‌گیری را به جای یک جمعیت از بیماران برای هر بیمار ارائه می‌دهد که این وجه تمایز آن است. در اینجا تأکید بر یافتن دانش مفید، معترض و کاملی از داده‌های است [5]. هدف از این مطالعه، شناسایی عوامل مؤثر و برآورد میزان بقای پنج ساله کلیه پیوندی به کمک الگوریتم ژنتیک و روش‌های یادگیری ماشین و داده‌کاوی است که بر روی داده‌های بیمارستان امام رضا(ع) و چهارمین شهید محرب کرمانشاه در سال‌های ۱۴۰۰-۱۴۰۲ انجام شده است. با توجه به این‌که داده‌های مورد استفاده در این مطالعه نامتوازن است و داده‌های نامتوازن باعث ناکارامدی نتایج می‌شوند؛ از این‌رو برای بهبود نتایج متوازن‌سازی داده‌ها ضرورت دارد.

این مطالعه در ۶ بخش سازماندهی شده است. در بخش ۲ به پیشینهٔ پژوهش و در بخش ۳ به مواد و روش‌ها خواهیم پرداخت؛ در بخش ۴ به بیان الگوریتم پیشنهادی می‌پردازیم؛ سپس یافته‌ها و بحث در بخش ۵ مطرح می‌شود و در بخش ۶ نتیجه‌گیری آمده است.

## ۱- مقدمه

یکی از مهم‌ترین بیماری‌های کلیوی، نارسایی کلیه است که با کاهش شدید عملکرد کلیه همراه است. مهم‌ترین اقدام برای درمان بیماری، برطرف کردن عامل بیماری است تا باقی‌مانده بافت سالم کلیه حفظ شود و کلیه‌ها بتوانند کار خود را به خوبی انجام دهند. در مراحل اولیه نارسایی کلیه بیمار با مصرف دارو ممکن است، بیماری اش را کنترل کند، اما اگر نارسایی کلیه در مراحل پیشرفتی باشد، کلیه قادر نخواهد بود عمل تصفیه خون را به خوبی انجام دهد و سموم را از بدن خارج کند. بنابراین افراد مبتلا باید با روش مصنوعی دیالیز و درنهایت، پیوند کلیه درمان شوند. در مرحله نهایی نارسایی کلیه، پیوند کلیه، زندگی مطلوب‌تر و کاهش خطر مرگ و میر را برای بیماران فراهم می‌کند [1]. سابقهٔ پیوند کلیه به سال ۱۹۵۴ میلادی (۱۳۳۳ شمسی) در آلمان بر می‌گردد. نخستین پیوند کلیه در ایران در سال ۱۹۶۷ در شیراز انجام گرفت. منابع تأیین عضو پیوندی شامل افراد زندهٔ فامیل، افراد زندهٔ غیرفamilی و جسد است. یکی از اهداف اصلی برنامه‌های پیوند عضو، فراهم‌کردن عضو پیوندی مناسب برای هر بیمار نیازمند به آن است که امروزه کمبود عضو، بزرگ‌ترین مانع در دست‌یابی به این هدف به‌شمار می‌رود [2]. بررسی میزان یا پیش‌بینی بقای کلیه پیوندی بعد از عمل پیوند اهمیت زیادی دارد؛ زیرا بیماری که می‌تواند سال‌ها با دیالیز زندگی کند با پیوندی نامناسب ممکن است، زندگی خود را از دست دهد؛ بنابراین انجام مطالعه بر روی پیش‌بینی بقای کلیه پیوندی و

فصلنامه



به نام طبقه کمینه هستند. تعداد نمونه داده‌های بیماران بیشتر از تعداد نمونه‌های افراد سالم است. الگوریتم‌های طبقه‌بندی با داده‌های نامتوازن باعث هم‌گراشدن طبقه کمینه به طبقه بیشینه می‌شود؛ به طوری که نمونه‌های طبقه کمینه در طبقه بیشینه طبقه‌بندی می‌شوند؛ درنهایت طبقه‌بندی ناکارا خواهد بود. یکی از روش‌های حل این مشکل، استفاده از روش‌های نمونه‌سازی از جمله بیش‌نمونه‌برداری و زیرنمونه‌برداری است. اسکوبی و همکارانش، در سال ۲۰۱۵ دو روش نمونه‌برداری (بیش‌نمونه‌برداری و زیرنمونه‌برداری) را با سیزده مجموعه داده‌های نامتوازن با الگوریتم‌های طبقه‌بندی بیزین ساده و J48 مطالعه کردند و نشان داده‌اند که روش بیش‌نمونه‌برداری در الگوریتم طبقه‌بندی دارای دقیق‌تر است [10]. در این مطالعه سعی بر ارائه مدلی است که دقیق‌تر بالاتری نسبت به مدل‌های پیشین داشته باشد.

### ۳- مواد و روش‌ها

در این بخش به بیان مقدماتی راجع به مواد و روش‌های به کاررفته در مدل پیشنهادی می‌پردازیم.

#### ۳-۱- نمونه‌برداری

در مجموعه‌داده‌های متوازن، تعداد نمونه‌ها در هر طبقه مجموعه‌داده برابر است؛ درصد متوازن‌سازی مجموعه داده‌های متوازن، در حدود پنجاه درصد است. در حالی که در مجموعه داده‌های نامتوازن، تعداد نمونه‌های یک طبقه از طبقه دیگر بیشتر است؛ درصد متوازن‌سازی مجموعه داده‌های نامتوازن، خیلی کمتر از پنجاه درصد است. از طرفی طبقه‌بندی داده‌های نامتوازن ناکارا است. بنابراین طبقه‌بندی مجموعه داده‌های نامتوازن یک چالش مهم در یادگیری ماشین است [10].

در الگوریتم‌های استاندارد حوزه یادگیری ماشین و داده‌کاوی، فرض بر این است که توزیع طبقه‌ها متوازن باشد. از این‌رو در صورت استفاده از این الگوریتم‌ها در طبقه‌بندی داده‌های نامتوازن، مدل به دست آمده به سمت نمونه‌های آموزشی طبقه بزرگ‌تر متمایل می‌شود که سبب کاهش دقیق‌تر مدل حاصل در پیش‌بینی طبقه کمینه می‌شود [12,11]؛ لذا در این مطالعه جهت رسیدن به مدلی با دقیق‌تر در مرحله پیش‌پردازش داده، از روش نمونه‌برداری برای متوازن‌سازی داده‌ها استفاده می‌شود. نمونه‌برداری به دو روش بیش‌نمونه‌برداری از طبقه کمینه و زیرنمونه‌برداری از طبقه بیشینه است.

## ۲- پیشینهٔ پژوهش

در سال ۲۰۰۷ سنتوری و همکاران به ارائه مدل یادگیری شبکه عصبی برای پیش‌بینی نتیجه بقای کلیه پیوندی در مورد پیوند کلیه اطفال پرداخته‌اند که دقیق‌تر مدل حاصل ۸۷/۱۴ درصد است [6]. در سال ۲۰۰۹، اشرفی و همکاران جهت بررسی بقای کلیه پیوندی در اصفهان با استفاده از شبکه عصبی مدلی با دقیق‌تر درصد ارائه داده‌اند. همچنان نشان دادند بر اساس نتایج این مطالعه، میزان بقای کلیه پیوندی در پیوندهای زنده، بیش‌تر از پیوندهای جسدی است [3].

در مطالعه انجام‌شده در سال ۲۰۱۰ در ایران توسط حشیانی و همکاران، عامل سن دهنده و سن گیرنده عامل معناداری بر بقای کلیه پیوندی شناسایی شدند [4]. در سال ۲۰۱۲ در آمریکا، برای مدل سازی ۵۱۴۴ مورد پیوند کلیه با ۴۸ متغیر از روش شبکه بیزی استفاده شد، شاخص توده بدنی، نژاد، جنسیت گیرنده و سن دهنده از عوامل مؤثر در بقای پیوند شناسایی شدند. در مطالعه‌ای در سال ۲۰۱۳ در آلمان برای پیش‌بینی میزان فیلتراسیون گلومار گیرنده یک سال پس از پیوند، روش ماشین‌بردار پشتیبان استفاده شد و متغیرهای مؤثر سن و سطح کراتینین دهنده و جنسیت و وزن گیرنده شناسایی شدند [8]. در سال ۲۰۱۳ بیرونی و همکاران با ارائه مدل شبکه عصبی پرسپترون چندلایه بر روی ۷۵۶ داده بیمارستان امام رضا(ع) و چهارمین شهید محرب کرمانشاه که در سال‌های ۲۰۰۲-۲۰۱۲ جمع‌آوری شده است به پیش‌بینی بقای پنج ساله کلیه پیوندی پرداخته و مدلی با دقیق‌تر درصد به دست آورده‌اند [9]. میرزائی و همکاران در سال ۲۰۱۶، مطالعه‌ای بر روی اطلاعات پرونده ۴۲۳ بیمار پیوند کلیه مرکز آموزشی-درمانی افضلی بور شهر کرمان انجام داده‌اند. آنها الگوریتم ژنتیک را جهت شناسایی متغیرهای تأثیرگذار به کار برده‌اند. ویژگی‌های شاخص توده بدنی، جنسیت گیرنده، سن دهنده، همسانی گروه خونی دهنده و گیرنده و سابقه پیوند کلیه به عنوان متغیرهای بهینه شناسایی شدند که صحت پیش‌بینی مدل پیشنهادی با این تعداد متغیر ۹۱,۶٪ درصد است [1].

در تمامی مطالعات انجام‌شده در این زمینه، پژوهش‌گران به دنبال شناسایی مجموعه ویژگی‌های عوامل مؤثر در بقای پیوند و همچنین طراحی مدل‌های پیش‌بینی با دقیق‌تر بالاتر و خطای پایین‌تر بوده‌اند.

به طور عمومی داده‌های پزشکی، داده‌های نامتوازن هستند؛ به عبارتی در مجموعه داده‌های پزشکی، تعداد نمونه‌های یک طبقه به نام طبقه بیشینه بیش‌تر از طبقه دیگر

### ۱-۱-۳- بیش نمونه برداری

این روش برای تعادل رساندن توزیع طبقه، از روش جایگزینی نمونه های طبقه کمینه استفاده می کند، تا توازن در مجموعه آموزشی برقرار شود.

### ۱-۲- زیر نمونه برداری

در این روش به طور تصادفی نمونه هایی از طبقه بیشینه حذف می شود؛ تا زمانی که طبقه کمینه درصدی از طبقه بیشینه شود. به این ترتیب توازن در مجموعه آموزشی برقرار می شود [14,13].

## ۲- الگوریتم ژنتیک

الگوریتم ژنتیک در سال ۱۹۷۵ توسط هلند بر طبق نظریه تکاملی داروین ارائه شده است. الگوریتم با یک مجموعه از جمعیت (جواب های اولیه) که از طریق کروموزوم نشان داده می شوند، شروع می شود. در الگوریتم ژنتیک هر کروموزوم به عنوان راه حل مسئله در نظر گرفته می شود. به طوری که هر کروموزوم شامل تعدادی زن است که به صورت عددی برطبق مسئله کدگذاری شده است. تعداد جمعیت در الگوریتم ژنتیک براساس پیچیدگی مسئله تعیین می شود. ابتدا کروموزوم ها براساس مسئله کدگذاری و در ادامه پارامترهای الگوریتم از جمله تعداد جمعیت تعیین می شود؛ سپس در گام بعدی با استفاده از تابع ارزیاب هر کروموزوم در نسل نخست ارزیابی می شود. بر اساس مقادیر تابع ارزیاب، کروموزوم ها برای ایجاد نسل بعدی انتخاب می شوند. چند روش برای انتخاب کروموزوم های والدین وجود دارد. یکی از رایج ترین روش های انتخاب کروموزوم انتخاب چرخ رول است که تابع ارزیاب کروموزوم را بر طبق بالاترین مقدار برای نسل های بعدی انتخاب می کند. نسل های بعدی براساس روش های الهام از تکامل طبیعی، مانند ترکیب و جهش تولید می شوند. در مرحله ترکیب، ویژگی های والدین برای ایجاد فرزندان ترکیب می شوند تا این که کروموزوم های بهتری ایجاد شوند. در ادامه در مرحله جهش در یک تعداد از مقادیر کروموزوم های مرحله ترکیب، تغییر ایجاد می شود تا از هم گرایی به بهینه محلی اجتناب شود و همچنین تنوع در نسل بعدی ایجاد شود. امید است که جمعیت جدید نسبت به جمعیت قبلی بهتر باشد. این فرآیند تا برقراری شرطی که تعیین شده است، ادامه می یابد [16,15].

مجموعه داده های پیشکی دارای تعداد ویژگی های مختلفی است که سوابق بیمار را توصیف می کنند. برخی از

این ویژگی ها ممکن است، در پیش بینی و تشخیص بیماری مؤثر نباشند یا تعدادی از ویژگی ها تأثیر مشابه با هم در پیش بینی و تشخیص بیماری داشته باشند. از طرفی مجموعه داده با تمام ویژگی های بیمار نیازمند حافظه و زمان بیشتر محاسباتی است. همچنین ممکن است، دقیق طبقه بندی را کاهش دهن. از این رو تجزیه و تحلیل مجموعه داده با ابعاد بالا و تعداد نمونه کم در هوش مصنوعی اهمیت خاصی دارد. یکی از روش های مورد استفاده انتخاب ویژگی مؤثر با استفاده از الگوریتم ژنتیک است.

انتخاب ویژگی و تعیین عوامل مؤثر بیشتر به انتخاب یک زیر مجموعه از ویژگی های اولیه مربوط می شود که برای هدف مورد نظر اطلاعات لازم و کافی را دربر داشته باشد، به طوری که ویژگی های انتخاب شده معیار طبقه بندی را بهبود می دهد و پیش بینی و تشخیص بیماری با دقیق بالاتری انجام می شود. بنابراین انتخاب ویژگی و تعیین عوامل مؤثر کمک می کند تا ویژگی های افزونه و نامرتب کاهش یابند و همچنین کمک به بهبود عملکرد طبقه بند می کند [17].

### ۳- طبقه بند نزدیک ترین همسایه

الگوریتم نزدیک ترین همسایه برای یک داده آزمایشی به دنبال K نمونه از نزدیک ترین نمونه ها می شود. نزدیکی دو نمونه با به دست آوردن فاصله میان این دو نمونه محاسبه می شود. پس از یافتن این k داده مشابه با نمونه آزمایشی، با رأی اکثریت بر چسب طبقه داده آزمایشی انتخاب می شود. چنان چه مقدار یک برای k تنظیم شود، در این صورت طبقه نزدیک ترین داده به نمونه آزمایشی به عنوان طبقه تخمینی ارائه می شود [18]. در پیش بینی بقای پنج ساله کلیه پیوندی، الگوریتم طبقه بند نخستین نزدیک ترین همسایه با معیار فاصله اقلیدسی به کار رفته است. مهم ترین معیار برای تعیین کارایی یک الگوریتم طبقه بندی، دقیق یا نرخ طبقه بندی است که این معیار، دقیق کل یک طبقه بند را محاسبه می کند. معیار دقیق براساس رابطه (1) محاسبه می شود:

$$precision = \frac{TP}{TP+FP} \quad (1)$$

علاوه بر معیار دقیق، معیار صحت طبقه بندی براساس رابطه (2) برای مقایسه روش پیشنهادی با کارهای پیشین لحاظ شده است.

$$accuracy = \frac{TP+TN}{TP+T+FP+FN} \quad (2)$$

صفر و یک مقداردهی شده و مقادیر دیگر ویژگی‌ها مقادیر پیوسته حقیقی است که در مرحله پیش‌پردازش هنجارسازی شده است.

(جدول-۱): توصیف ویژگی‌های بیماران کلیه پیوندی  
(Table-1): description features of transplanted kidney patients

نام ویژگی	ردیف	نام ویژگی	ردیف
جنسیت دهنده	۱۲	وزن	۱
جنسیت گیرنده	۱۳	سن دهنده	۲
نسبت فامیلی	۱۴	سن گیرنده	۳
دهنه RH	۱۵	اورة قبل	۴
RH گیرنده	۱۶	اورة بعد	۵
سن دهنده بالا ۳۰	۱۷	کراتین بعد پیوند	۶
بیماری اولیه	۱۸	کراتین قبل پیوند	۷
سن گیرنده بالا ۴۰	۱۹	هموگلوبین بعد پیوند	۸
کراتین بعد از پیوند	۲۰	هموگلوبین قبل پیوند	۹
نتیجه بقای پنج ساله پیوند	۲۱	هموتین بعد پیوند	۱۰
		هموتین قبل پیوند	۱۱

#### ۴-۲- پیش‌پردازش داده

مرحله دوم، پیش‌پردازش داده است که یکی از بحرانی‌ترین مراحل موجود در فرآیند یادگیری ماشین و داده‌کاوی است. این مرحله شامل دو مرحله هنجارسازی و متوازن‌سازی داده‌های جمع‌آوری‌شده بیماران کلیه پیوندی است.

#### ۴-۲-۱- نرمال‌سازی

تولید نتایج کارا به داده‌های مناسب در مدل وابسته است. پیش‌پردازش داده‌ها با پاک‌سازی داده‌ها همراه است که برای کامل کردن داده‌های ناقص از میانگین و مذ استفاده شده است. بعضی از مقادیر ویژگی‌های داده در محدوده یا دامنه متفاوتی قرار دارند؛ از این‌رو برای یک‌پارچه‌سازی داده‌ها از رابطه نرمال‌سازی زیر استفاده شده است.

$$Ni = \frac{(xi - \mu_i)}{\sigma_i} \quad (3)$$

که  $xi$  مقدار داده اولیه به‌ازای هر رکورد ،  $\mu_i$  و  $\sigma_i$  به ترتیب میانگین و انحراف معیار مقادیر هر ویژگی و  $Ni$  مقادیر نرمال شده است.

#### ۴-۲-۲- متوازن‌سازی

مجموعه داده اولیه جمع‌آوری شده از بیماران کلیه پیوندی نامتوازن هستند. به طوری که از ۷۵۶ داده بیماران کلیه پیوندی داده‌های متعلق به طبقه قبول پیوند-طبقه کمینه-۹۲- ۶۶۴ رکورد و داده‌های متعلق به طبقه رد پیوند-طبقه بیشینه-۲۱ رکورد

TN: بیان گر تعداد رکوردهایی است که دسته واقعی آنها منفی بوده و الگوریتم دسته‌بندی نیز دسته آنها را به درستی منفی تشخیص داده است.

TP: بیان گر تعداد رکوردهایی است که دسته واقعی آنها مثبت بوده و الگوریتم دسته‌بندی نیز دسته آنها را به درستی مثبت تشخیص داده است.

FP: بیان گر تعداد رکوردهایی است که دسته واقعی آنها منفی بوده و الگوریتم دسته‌بندی، دسته آنها را به اشتباه مثبت تشخیص داده است.

FN: بیان گر تعداد رکوردهایی است که دسته واقعی آنها مثبت بوده و الگوریتم دسته‌بندی، دسته آنها را به اشتباه منفی تشخیص داده است.

در این مطالعه معیار ارزیاب، دقّت طبقه‌بندی الگوریتم نخستین نزدیک‌ترین همسایه با اعتبار سنجی ضرب‌دری یکی بیرون به عنوانتابع شایستگی هر کروموزوم در نظر گرفته شده است. اعتبار سنجی ضرب‌دری یک روش ارزیابی است که مشخص می‌کند نتایج یک مدل تا چه اندازه قابل تعمیم و مستقل از داده‌های آموزشی است.

## ۴- الگوریتم پیشنهادی

در این مطالعه پیش‌بینی بقای پنج ساله کلیه پیوندی با استفاده از الگوریتم ژنتیک و الگوریتم نخستین نزدیک‌ترین همسایه انجام شده است. فرآیند کلی الگوریتم پیشنهادی در سه مرحله انجام شده است:

۱. جمع‌آوری داده
۲. پیش‌پردازش داده
۳. نرمال‌سازی
۴. متوزن‌سازی (روش پیش‌نمونه‌برداری و زیرنمونه‌برداری)
۵. مدل‌سازی (الگوریتم ژنتیک و الگوریتم اولین نزدیک‌ترین همسایه)

در ادامه هر مرحله به‌طور کامل شرح داده می‌شود.

#### ۱- جمع‌آوری داده

در این مرحله اطلاعات ۷۵۶ بیمار کلیه پیوندی با ۲۱ ویژگی از دو بیمارستان امام رضا(ع) و چهارمین شهید محرب کرمانشاه در سال‌های ۲۰۱۲-۲۰۰۱ جمع‌آوری شده است.

جدول (۱) ویژگی‌های بیماران کلیه پیوندی را نشان می‌دهد. جنسیت دهنده، جنسیت گیرنده، نسبت فامیلی، همبستگی گروه خونی و نتیجه بقای پنج ساله پیوند با مقادیر



است. در این صورت نرخ عدم توازن که برابر نسبت تعداد نمونه‌های طبقه کمینه به تعداد نمونه‌های طبقه بیشینه است در مجموعه‌داده مورد مطالعه دارای نرخ عدم توازن ۱۴ درصد است.

با توجه به این که داده‌های جمع‌آوری شده نامتوارن هستند و عدم توازن سبب می‌شود هم‌گرایی مدل طبقه‌بندی به طرف طبقه بیشینه صورت گیرد، که خود باعث کاهش کارایی طبقه‌بند می‌شود، لذا در مرحله پیش‌پردازشی برای متوازن‌سازی داده‌ها از روش‌های نمونه‌برداری استفاده شده است که دو روش پیش‌نمونه‌برداری و زیرنمونه‌برداری را شامل می‌شود.

#### ۴-۳- مدل‌سازی

برای شناسایی عوامل مؤثر بقای پنج‌ساله پیوند، مدلی براساس پیش‌بینی میزان بقای پنج‌ساله پیوند کلیه ارائه شده است که دست‌یابی به نتایج تحلیلی مشابه با اصل و تمام داده‌ها تضمین می‌شود. از طرفی پیداکردن بهترین عوامل مؤثر از میان تمام زیر مجموعه‌های ممکن، کاری بغرنج است؛ لذا استفاده از الگوریتم‌های مناسب برای جستجو در فضای ویژگی‌ها و انتخاب عوامل مؤثر امری ضروری است. بنابراین در مدل پیشنهادی، الگوریتم ژنتیک برای جستجو در فضای ویژگی‌ها به کار رفته است.

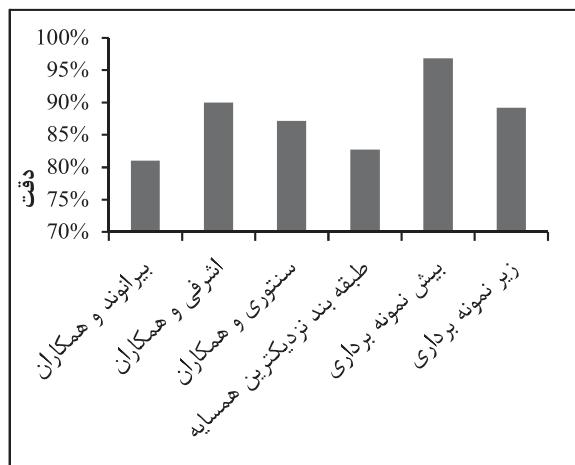
در الگوریتم ژنتیک باید ابتدا مقادیر پارامترهای الگوریتم مقداردهی شوند. مقداردهی پارامترها شامل تعیین تعداد اعضای جمعیت، نرخ ترکیب، نرخ جهش و تعداد تکرارهای الگوریتم است که در این مطالعه به ترتیب از مقادیر ۰/۹۵، ۰/۰۱، ۰/۰۲۰ استفاده شده است.

در این مدل به‌منظور شناسایی عوامل مؤثر برای بازنمایی هر کروموزوم از رشته‌ای از صفر و یک‌ها استفاده شده است که یک به معنی متغیر مؤثر در پیش‌بینی بقای پنج‌ساله کلیه پیوندی و صفر به معنی متغیر غیر مؤثر است. برای مثال رشته (۱۱۱۱۱۱۰۱۱۱۱۱۱۰۱۱۰۱۱۱۱۱۱۰) به عنوان یک کروموزوم در الگوریتم ژنتیک است که تمامی ویژگی‌ها به جز ویژگی‌های ۴، ۱۲ و ۲۰ انتخاب شده‌اند. با تولید جمعیت اولیه (کروموزوم‌ها) به صورت رشته دودویی تصادفی، عوامل انتخاب شده توسط الگوریتم ژنتیک به طبقه‌بند نخستین نزدیک‌ترین همسایه داده می‌شود.

این طبقه‌بند با ویژگی‌های انتخابی موجود در هر کروموزوم، و اعتبارسنجی ضرب‌دری یکی بیرون میزان

شایستگی هر کروموزوم را بر اساس دقت پیش‌بینی بقای پنج‌ساله کلیه پیوندی تعیین می‌کند. به عبارتی مقدار تابع شایستگی (برازندگی)، میزان دقت طبقه‌بند نخستین نزدیک‌ترین همسایه برای هر کروموزوم است که نشان‌دهنده میزان مؤثربودن ویژگی‌های انتخابی در پیش‌بینی داده بیماران پیوند کلیه در طبقه‌بند نخستین نزدیک‌ترین همسایه است. برای ساختن نسل بعد از عمل‌گرهای وراشتی، انتخاب، ترکیب و جهش استفاده می‌شود. ابتدا با عمل گر انتخاب چرخ رولت، کروموزوم‌ها با شایستگی بیشتر برای تولید مثل انتخاب می‌شوند؛ سپس عمل گر ترکیب دو کروموزوم والد را برای ساختن فرزندان ترکیب می‌کند تا کروموزوم‌های بهتری ایجاد شوند. در این مدل از تلفیق سه عمل گر ترکیب تک نقطه‌ای، دونقطه‌ای و یکنواخت استفاده شده است که در هر مرحله یکی از این سه عمل گر به صورت تصادفی برای ترکیب انتخاب می‌شود و دو فرزند برای نسل جدید تولید می‌کند [۱۵]. در ترکیب تک نقطه‌ای یک نقطه به صورت تصادفی در کروموزوم‌های والد انتخاب شده و با یکدیگر ترکیب می‌شوند (شکل ۱). در ترکیب دونقطه‌ای، دو نقطه به صورت تصادفی در کروموزوم‌های والد انتخاب شده و سپس ژن‌های بین نقاط انتخاب شده به صورت یک‌درمیان با یکدیگر عوض می‌شوند (شکل ۲). در ترکیب یکنواخت، تولید فرزندان از روی ماسکی که به صورت تصادفی تولید می‌شود، انجام می‌گیرد. به این صورت که اگر در ماسک بیتی برابر یک باشد، والد نخست در تولید فرزند نخست و والد دوم در تولید فرزند دوم به کار می‌رود. اگر در ماسک بیتی برابر صفر باشد، والد نخست در تولید فرزند دوم و والد دوم در تولید فرزند نخست به کار می‌رود (شکل ۳).

به کمک عمل گر جهش یک کروموزوم از نسل قبلی انتخاب و یک ژن از آن را به صورت تصادفی انتخاب و معکوس می‌کنیم. بعد از اعمال الگوریتم‌های ترکیب و جهش و با داشتن جمعیت قبلی، به دست آوردن کروموزوم‌های جدید و تعیین میزان شایستگی آن‌ها توسط طبقه‌بند نخستین نزدیک‌ترین همسایه کروموزوم‌ها با تابع شایستگی بیشتر، نسل بعد جمعیت را می‌سازند؛ این روال ادامه می‌باید تا به تعداد مراحل مشخصی از تکرار الگوریتم برسد [۱۹-۲۱]؛ در انتها با رسیدن به آخرین تکرار و توقف الگوریتم ژنتیک ویژگی‌های مؤثر با بیشترین دقت پیش‌بینی شناسایی می‌شوند.



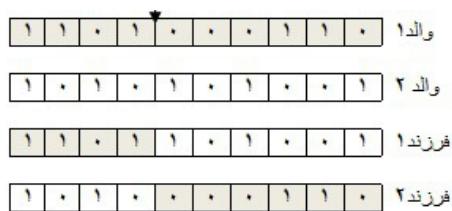
(شکل-۴): مقایسه روش‌های پیشنهادی با سایر روش‌ها  
(Figure-4): Comparison of proposed methods with other methods

هموتین بعد پیوند، هموتین قبل پیوند، کراتین بعد از پیوند با دقّت ۸۹/۲ درصد در روش زیرنمونه‌برداری شناسایی شدند؛ که در مقایسه با سایر روش‌های کارشده در حوزه برآورد میزان بقای کلیه پیوندی، از دقّت پیش‌بینی بالایی برخوردارند که می‌توان با درنظرگرفتن این عوامل به افزایش میزان بقای پنج ساله کلیه پیوندی کمک کرد.

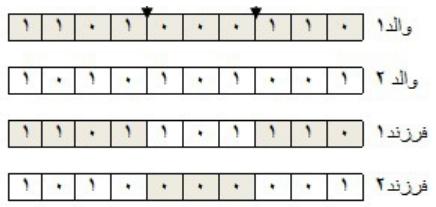
ستوری و همکاران در سال ۲۰۰۷ با استفاده از مدل یادگیری شبکه عصبی برای پیش‌بینی نتیجه بقای کلیه پیوندی اطفال را با دقّت ۸۷، ۱۴ درصد پیش‌بینی کردند [6]. در سال ۲۰۰۹، اشرفی و همکاران با استفاده از شبکه عصبی مدلی با دقّت نود درصد جهت بررسی بقای کلیه پیوندی در اصفهان ارائه دادند. همچنین بر اساس نتایج این مطالعه نشان دادند، میزان بقای کلیه پیوندی در پیوندی‌های زنده بیشتر از پیوندی‌های جسدی است [3]. در سال ۲۰۱۳ بیرونوند و همکاران با شبکه عصبی پرسپترون چندلایه بر روی ۷۵۶ داده بیمارستان امام رضا(ع) و چهارمین شهید محرب کرمانشاه در سال‌های ۲۰۱۲-۲۰۰۲ بررسی کردند و به مدلی با دقّت ۸۱ درصد دست یافتند [9].

شکل (۴) مقایسه‌ای از کارهای اخیر با روش پیشنهادی در این مقاله (در دو حالت بیش‌نمونه‌برداری و زیرنمونه‌برداری داده‌ها) را نشان می‌دهد.

نتایج مقایسه حاکی از آن است که روش نمونه‌برداری در مقایسه با مطالعات پیشین، از نقطه‌نظر دقّت برتر است. با توجه به این‌که داده‌های بیرونوند و همکارانش با داده‌های استفاده شده در این مطالعه یکسان است، الگوریتم پیشنهادی در این مقاله با متوازن‌سازی داده‌ها، دقّت قابل توجه‌ای را برآورد می‌کند و با اطمینان خاطر بیشتری می‌توان از مدل پیشنهادی برای پیش‌بینی بقای کلیه پیوندی استفاده کرد.



(شکل-۱): ترکیب تک نقطه‌ای  
(Figure-1): One-point crossover



(شکل-۲): ترکیب دو نقطه‌ای  
(Figure-2): Two-point crossover



(شکل-۳): ترکیب یک-نواخت  
(Figure-3): Uniform crossover

## ۵- یافته‌ها و بحث

مدل پیشنهادی در نرم‌افزار MATLAB روی لپ‌تاپی با پردازندۀ اینتل هفت‌هسته‌ای پیاده‌سازی شده است. مجموعه داده‌های جمع‌آوری شده که در بخش ۱-۴ بحث شد، برای ارزیابی عملکرد مدل پیشنهادی مورد استفاده قرار گرفت. با توجه به این‌که فضای جستجوی داده‌ها نمایی است، لذا استفاده از الگوریتم جستجوی ژنتیک برای جستجوی چنین فضایی بهمنظور تعیین عوامل اثرگذار بر میزان بقای پنج ساله کلیه پیوندی در زمان معقول نقش حیاتی دارد. براین اساس با جستجوی فضای زیرمجموعه‌ها در داده‌های کلیه پیوندی، مهم‌ترین عوامل مؤثر بر بقای پنج ساله کلیه در روش پیش‌نمونه‌برداری و زیرنمونه‌برداری شناسایی شدند.

در مدل به دست آمده بعد از رسیدن الگوریتم ژنتیک به تعداد تکرار معین عوامل وزن، سن دهنده، سن گیرنده، اوره، قبل پیوند، کراتین قبل پیوند، هموگلوبین بعد پیوند، هموگلوبین قبل پیوند، جنسیت دهنده، RH دهنده و گیرنده، سن دهنده بالای سی، بیماری اولیه، سن گیرنده بالای چهل به عنوان ویژگی‌های تأثیرگذار در بقای پنج ساله کلیه پیوندی با دقّت ۹۶/۸ درصد در روش پیش‌نمونه‌برداری و عوامل کراتین قبل پیوند، هموگلوبین بعد پیوند، هموگلوبین قبل پیوند،

## 7- References

## ۷- مراجع

- [۱] میرزایی، محترم، و فیروز آبادی، سید محمد، "کاربرد داده‌کاوی در پیش‌بینی بقای پیوند کلیه و شناسایی متغیرهای تأثیرگذار در بقای کلیه پیوندی"، مجله انفورماتیک سلامت و زیست پزشکی، مرکز تحقیقات انفورماتیک پزشکی، دوره ۳، شماره ۱، صفحات ۹-۱، ۱۳۹۵.
- [۱] M. Mirzaei, and M. Firooz Abadi, "The impact of data mining on prediction of renal transplantation survival and identifying the effective factors on the transplanted kidney," *Journal Of Health and Biomedical Informatics, Medical Informatics Research Center*, vol. 3, pp. 1-9, 2016.
- [۲] جوانروح گیوی، نیلوفر، علیمی، رسول، اسماعیلی، حبیب ا...، شاکری، محمدتقی، و شمسا، علی، "عوامل مؤثر بر بقای پیوند کلیه و برآورد خطر رد پیوند برای پیوند شدگان مراجعه کننده به بیمارستان قائم مشهد"، مجله دانشگاه علوم پزشکی خراسان شمالی، دوره ۵، شماره ۲، صفحات ۳۲۱-۳۱۵. ۱۳۹۲
- [۲] N. Javanroh Givi, R. Alimi, H. Esmaily, M. T. Shakeri, and A. Shamsa, "Assessment of effective factors on renal transplantation estimation of rejection hazard for transplanted in Mashhad Qaem hospital," *Journal of North Khorasan University of Medical Sciences*, vol. 5, pp. 315-321, 2013.
- [۳] اشرفی، مهدی، و همکاران، "پیش‌بینی بقای پنج ساله پیوند کلیه با استفاده از مدل شبکه عصبی مصنوعی: گزارش ۲۲ سال پی‌گیری از ۳۱۶ بیمار در اصفهان"، مجله دانشکده پزشکی، دانشگاه علوم پزشکی تهران، دوره ۶۷، شماره ۵، صفحات ۳۵۹-۳۵۳. ۱۳۸۸
- [۳] M. Ashrafi, and et. al, "Application of artificial neural network to predict graft survival after kidney transplantation: reports of 22 years follow up of 316 patients in Isfahan," *Tehran University Medical Journal*, vol. 67, pp. 353-359, 2009.
- [۴] الماسی حشیانی، امیر، رجایی‌فرد، عبدالرضا، حسن‌زاده، جعفر، و صلاحی، حشمت...، "تحلیل بقا پیوند کلیه و ارتباط آن با سن و جنس دهنده و گیرنده عضو بین بیماران پیوند شده"، مجله علمی دانشگاه علوم پزشکی سمنان، جلد ۱۱، شماره ۴، صفحات ۳۰۷-۳۰۲. ۱۳۸۹
- [۴] A. Almasi Hashiani, A. Rajaeefard, J. Hassanzade, and H. Salahi, "Survival analysis of renal transplantation and its relationship with age and sex," *Koomesh*, vol. 11, pp. 302-307, 2010.

میرزائی و همکاران در سال ۲۰۱۶، به کمک الگوریتم رژتیک و روش هم‌جوشی اطلاعات متغیرهای تأثیرگذار را شناسایی کردند که این مدل دارای صحت پیش‌بینی ۹۱,۶٪ درصد است. در الگوریتم پیشنهادی با شناسایی متغیرهای تأثیرگذار توسط الگوریتم رژتیک و طبقه‌بند نزدیک‌ترین همسایه صحت پیش‌بینی در روش بیش‌نمونه‌برداری و زیرنمونه‌برداری به ترتیب ۹۱/۸۹ و ۹۱/۸۵ درصد به دست آمده است [۱]. نتایج مقایسه نشان می‌دهد که روش پیشنهادی با متوازن‌سازی داده‌ها براساس بیش‌نمونه‌برداری، سبب ایجاد مدل پیش‌بینی با صحت بالاتر شده است.

## ۶- نتیجه‌گیری

امروزه روش‌های آماری با افزایش تعداد مشاهدات و تعداد ویژگی‌ها یا عوامل مربوط به یک مشاهده کارایی خود را از دست داده‌اند. ابزارهای یادگیری ماشین و داده‌کاوی کوششی برای به دست آوردن اطلاعات مفید از میان این داده‌های است. این مطالعه با رویکرد شناسایی عوامل مؤثر و پیش‌بینی بقای پنج ساله کلیه پیوندی در داده‌های نامتوازن به کمک الگوریتم رژتیک و روش یادگیری طبقه‌بند نزدیک‌ترین همسایه انجام شده است. با توجه به این که داده‌های آموزشی توزیع طبقه نامتوازنی دارند و طبقه‌بند نیازگرفته شده اغلب به طبقه بیشینه متمایل می‌شود، این موضوع به پیش‌بینی بسیار ضعیفی از طبقه کمینه منجر می‌شود؛ لذا برای بهبود کارایی از دو روش نامتوازن‌سازی بیش‌نمونه‌برداری و زیرنمونه‌برداری استفاده شد. از طرف دیگر افزایش تعداد ویژگی‌ها و عوامل، چالش‌های محاسباتی بالایی را به دنبال دارد و منجر به یافتن الگوهای نامعتبر می‌شود. بنابراین انتخاب عوامل مؤثر نقش مهمی را در تعیین عملکرد پیش‌بینی الگوریتم نزدیک‌ترین همسایه به عنوان یک مدل یادگیری بازی می‌کنند. با انتخاب عوامل مؤثر و کاهش ابعاد، داده‌ها علاوه‌بر کاهش زمان محاسبات، دقّت پیش‌بینی بر میزان بقای پنج ساله کلیه پیوندی افزایش یافته است.

همان‌طور که مشاهده شد در هر دو روش بیش‌نمونه‌برداری و زیرنمونه‌برداری دقّت پیش‌بینی میزان بقای پنج ساله کلیه پیوندی افزایش یافته است که نشان دهنده بهینه‌بودن روش پیشنهادی است. علاوه‌بر آن، روش بیش‌نمونه‌برداری نسبت به روش زیرنمونه‌برداری دقّت و صحت بهتری را داشته و روش مناسبی در جهت شناسایی عوامل مؤثر در پیش‌بینی بقای پنج ساله کلیه پیوندی است.

فصلنامه



- [16] H. Hoglund, " Tax payment default prediction using genetic algorithm-based variable selection," *Expert Syst Appl*, vol. 88, pp. 368-375, 2017.
- [17] S. Nagpal, S. Arora, S. Dey and Shreya, " Feature selection using gravitational search algorithm for biomedical data," *Procedia Comput Sci*, vol. 115, pp. 258-265, 2017.
- [18] X. We and et. al, "Top 10 algorithms in data mining," *knowl Inf Syst*, vol. 14, pp. 1-37, 2008.
- [19] J. H. Holland, " Adaptation in natural and artificial systems," *University of Michigan Press*, 1975.
- [20] D. E. Goldberg, "Genetic algorithms in search optimization and machine learning," *Addison-Wesley Publishing, INC. Reading. Mass*, 1989.
- [21] S. Olariu and A. Y. Zomaya, "Handbook of bioinspired algorithms and applications," *Taylor & Francis Group, LLC Press*, 2006.
- نسیبه امامی مدارک کارشناسی و کارشناسی ارشد خود را به ترتیب از دانشگاه یزد و شهید باهنر کرمان در رشته علوم کامپیوتر گرایش سامانه‌های هوشمند اخذ کرده است. وی هم‌اکنون عضو هیأت علمی دانشگاه کوثر بجنورد است و در زمینه‌های یادگیری ماشین، هوش مصنوعی، بهینه‌سازی و داده‌کاوی فعالیت دارد.**
- نشانی رایانامه ایشان عبارت است از:  
nasibeh.emami@kub.ac.ir
- زینب حسنی مدارک کارشناسی و کارشناسی ارشد خود را به ترتیب از دانشگاه یزد و مرکز تحصیلات تكمیلی زنجان در رشته علوم کامپیوتر گرایش سامانه‌های هوشمند اخذ کرده است. وی هم‌اکنون عضو هیأت علمی دانشگاه کوثر بجنورد است و در زمینه‌های یادگیری ماشین، هوش مصنوعی، بهینه‌سازی و هندسه محاسباتی فعالیت دارد.**
- نشانی رایانامه ایشان عبارت است از:  
Hassani@kub.ac.ir
- [5] T. D. Noia, and et al, "An end stage kidney disease predictor based on an artificial neural networks ensemble," *Expert Systems with Applications*, vol. 40, pp. 4438-4445, 2013.
- [6] G. Santori, I. Fontana, and U. Valente, "Application of an artificial neural network model to predict delayed decrease of serum creatinine in pediatric patients after kidney transplantation," *Transplant Proc*, vol. 39, pp. 1813-1819, 2007.
- [7] T. S. Brown, and et al, "Bayesian modeling of pre transplant variables accurately predicts kidney graft survival," *Am J Nephrol*, vol. 6, pp. 561-569, 2012.
- [8] J. Lasserre, S. Arnold, M. Vingron, P. Reinke, and C. P. Hinrichs, "Predicting the outcome of renal transplantation," *J Am Med Inform Assoc*, vol. 19, pp. 255-262, 2012.
- [9] A. H. Hashemian, B. beiranvand, M. Rezaei, A. Bardideh, and E. Zand-Karimi, "Comparison of artificial neural network of kidney transplant survival," *International Journal of Advanced Biological and biomedical Research*, vol. 1, pp. 1204-1212, 2013.
- [10] R. J. Oskouei and B. S. Bigham, "Over-sampling via under-sampling in strongly Imbalanced data," *International Journal of Advanced Intelligence Paradigms*, 2015.
- [11] M. M. Rahman, and D. N Davis, "Addressing the class imbalance problem in medical datasets," *Int J Machine Learning and Compute*, vol. 2, pp. 224-228, 2013.
- [12] N. V Chawla, "Data mining for imbalanced datasets: an overview," *Data Mining Knowledge Discovery Handbook*, 2005.
- [13] Y. Sun, A. K. C. Wong, and M. S Kamel, "Classification of imbalanced data: a review," *Int j patt Recogn Artif Intell*, vol. 4, pp. 687-719, 2009.
- [14] D. C. Li, C. W. Liu, and S. C. Hu, "A learning method for the class imbalance problem with medical datasets," *J comput Bio Medi*, vol. 5, pp. 509-518, 2010.
- [15] حسین خانی، فاطمه و ناصرشريف، بابک، "دو روش تبدیل ویژگی مبتنی بر الگوریتم‌های ژنتیک برای کاهش خطای دسته‌بندی ماشین بردار پشتیبان"، مجله پردازش علائم و داده‌ها، جلد ۲۴، شماره ۲، صفحات ۲۳-۴۹
- [15] F. Hoseinkhani, and B. Naser Sharif, " Two methods of converting feature based on genetic algorithms to reduce the classification error of support vector machine," *Journal Of signs and data Processing*, vol. 24, pp. 23-39, 2015.

