

## کاربرد داده کاوی در پیش بینی وضعیت بیماران نادر با استفاده از درخت تصمیم گیری؛ مطالعه موردی سامانه بنیاد بیماری های نادر ایران (سبنا)

محسن سجودی<sup>۱</sup>

فریبا ابراهیم بابایی<sup>۲</sup>

تاریخ دریافت: ۱۳۹۹/۰۴/۲۵ تاریخ چاپ: ۱۳۹۹/۰۴/۲۸

### چکیده

پایگاه داده ها در حوزه ی سلامت حاوی میزان وسیعی از داده های بالینی است که کشف ارتباطات و الگوها در آن می تواند به دانش جدید پزشکی بیانجامد. امروزه با توجه به پیدایش نظام های اطلاعات یکپارچه و رشد فن آوری اطلاعات، این مهم بیش از پیش نمایان شده است. این مطالعه با هدف شناسایی مزایای بالقوه ای که داده کاوی می تواند به بخش بهداشت و درمان، با استفاده از داده های بنیاد بیماری های نادر ایران به عنوان مطالعه موردی ارائه دهد، انجام پذیرفته است. معمول ترین روش داده کاوی که درخت تصمیم گیری می باشد برای تولید مدل پیش بینی با مصورسازی درخت برای انجام تحلیل پیش بینی بیماری نادر مورد استفاده قرار گرفته است. تمام مراحل فرایند داده کاوی با ابزاری به نام وکا (WEKA) انجام شده است. علاوه بر این، از وکا برای ارزیابی عملکرد پیش بینی از طریق اندازه گیری دقت، ویژگی و تحلیل حساسیت استفاده شده است. از جمله نتایج پژوهش حاضر، برخی از عواملی است که مراکز حمایتی بیمار می توانند در هنگام پیش بینی هزینه های درمان بیمار مورد توجه قرار دهند. شاخص جنس یا سن، به شدت تحت تاثیر مدت زمان بستری بیمار قرار می گیرد. بیمار سالخورده با بیماری نادر باید در مراکز درمانی تحت مراقبت و مدت زمان بستری طولانی تری نسبت به افراد جوان تر قرار گیرد. در نتیجه، درخت تصمیم گیری یک روش مفید و آموزنده برای انجام داده کاوی پیش بینی شده است.

### واژگان کلیدی

داده کاوی، مدل پیش بینی، بیماری نادر، درخت تصمیم گیری، تجزیه و تحلیل حساسیت، وکا

۱. دانشجوی دکتری مدیریت تحقیق در عملیات، دانشگاه فردوسی مشهد، ایران ([mohsened@gmail.com](mailto:mohsened@gmail.com))

۲. دانشجوی دکتری روانشناسی، دانشکده ادبیات و علوم انسانی، دانشگاه لرستان، خرم آباد، ایران ([faribaebrahimbabaie@yahoo.com](mailto:faribaebrahimbabaie@yahoo.com))

## ۱. مقدمه

سازمان‌های مراقبت بهداشتی-درمانی به دلیل پیچیدگی و تکرار فعالیت‌ها، سرشار از داده‌های بالینی هستند اما نیازمند تکنیک‌ها و ابزارهایی از قبیل داده‌کاوی می‌باشند که این داده‌ها را به اطلاعات مفید تبدیل نمایند (گیوسیدی<sup>۱</sup>، ۲۰۰۳). داده‌کاوی یا استخراج دانش از پایگاه داده‌ها<sup>۲</sup> فرایند مهم شناسایی الگوهای معتبر، جدید و قابل فهم در میان انبوهی از داده‌ها است. مفهوم داده‌کاوی شامل الگوریتم‌ها و متدهایی است که باعث استخراج اطلاعات از داده‌ها می‌شود (کلیرلی<sup>۳</sup>، ۲۰۰۶). داده‌ها در عصر امروزی یعنی عصر اطلاعات، عمده‌ترین دارایی برای سازمان‌های سلامت بوده است (پیاتسکی<sup>۴</sup>، ۲۰۰۰) و موفقیت سازمان‌های سلامت در گروهی جمع‌آوری، ذخیره و تحلیل آن‌ها است (بالی<sup>۵</sup>، ۲۰۰۵). با این وجود، جمع‌آوری و ذخیره‌ی میزان زیادی از داده‌ها می‌تواند یک نوع اتلاف محسوب شود؛ مگر این که داده‌ها به شکل سودمند استفاده شده و تبدیل به یک منبع مالی برای سازمان گردد. برای تبدیل این ارزش بالقوه به اطلاعات استراتژیک، بسیاری از سازمان‌ها به داده‌کاوی روی آورده‌اند (یوسفی و همکاران، ۲۰۱۶)؛ چرا که به واسطه‌ی داده‌کاوی امکان کشف روابط، روندها و الگوهای مخفی بین داده‌ها و دستیابی به دانش نوین در زمینه‌ی چالش‌های آشکار و نهان سازمان میسر خواهد شد. هدف این مطالعه شناسایی مزایای بالقوه‌ای است که داده‌کاوی می‌تواند به بخش سلامت در ایران ارائه نماید. این مطالعه براساس اطلاعات بیماران نادر متعلق به بنیاد بیماری‌های نادر ایران برای پیش‌بینی عوامل موثر بر بیماری‌های نادر با استفاده از درخت تصمیم‌گیری به عنوان تکنیک داده‌کاوی بهره گرفته است.

## ۲. مبانی نظری و ادبیات پژوهش

داده‌کاوی در اواخر دهه‌ی ۱۹۸۰ پدیدار گشت، در دهه‌ی ۱۹۹۰ گام‌های بلندی در این شاخه از علم برداشته شد و انتظار می‌رود در قرن حاضر نیز به رشد و توسعه‌ی خود ادامه دهد و پیش‌بینی‌ها حاکی از آن است که در دهه‌های آتی با توسعه‌ی انقلابی مواجه شود. مؤسسه‌ی فن‌آوری ماساچوست داده‌کاوی<sup>۶</sup> را یکی از ده فن‌آوری برتری می‌داند که نقش چشم‌گیری در تحول جهان خواهد داشت (راجرز<sup>۷</sup> و جویئر، ۲۰۱۱).

در ایران سیستم‌های اطلاعاتی، به پزشکان در انجام وظایف خود در رابطه با تصمیم‌گیری (سیستم پشتیبان تصمیم‌گیری) کمک می‌کند. بخش قابل توجهی از تحقیقات حوزه انفورماتیک پزشکی به صورت مستقیم و غیرمستقیم به این مقوله اختصاص دارد (ادل<sup>۸</sup> و همکاران، ۲۰۱۰). مزایای فناوری اطلاعات نه تنها در تصمیم‌گیری بلکه همچنین در پیاده‌سازی بهداشت الکترونیکی، مراقبت‌های پرستاری و غیره نیز دیده شود که می‌تواند خدمات بهداشت عمومی ایران را بهبود بخشد؛ پرایز<sup>۹</sup> و همکاران (۲۰۱۳) اظهار داشتند که کارکنان مدیریت اطلاعات مراقبت‌های بهداشتی برای کنترل هزینه‌های مراقبت بهداشتی و مدیریت استفاده از خدمات، از تلاش‌های مختلف مانند پیاده‌سازی پرونده‌های مدیریتی، بررسی بهره‌وری و مدیریت بیماری‌ها استفاده می‌کنند. با این حال، به نظر نمی‌رسد تمام این برنامه‌ها در کنترل هزینه‌ها درست عمل نمایند. آنها روش‌های مختلفی را برای شناسایی بیماران مبتلا به بیماری نادر (به دلیل ریسک بالا در پذیرش) و کنترل هزینه

<sup>1</sup> Giudici

<sup>2</sup> KDD (knowledge discovery in database)

<sup>3</sup> cleary

<sup>4</sup> Piatetsky-Shapiro

<sup>5</sup> Bali

<sup>6</sup> Massachusetts institute of technology

<sup>7</sup> Rogers

<sup>8</sup> Adel

<sup>9</sup> Price

ها به روش حرفه‌ای پزشکی و مدل‌های پیش‌بینی طراحی شده، پیشنهاد می‌کنند که از داده کاوی به عنوان یکی از انواع این روش‌ها یاد می‌کنند (پرایز و همکاران، ۲۰۱۳)؛ بنابراین، بنیاد بیماری‌های نادر ایران (رادو-آی-آر) باید از مدیریت داده خوبی بهره‌مند باشد. با این حال، بنیاد با برخی از چالش‌ها در فرایند مدیریت داده‌ها، به ویژه در بخش سلامت، جایی که حجم زیادی اطلاعات باید سازمان یافته و ذخیره شود، روبرو است. بنیاد بیماری‌های نادر ایران (رادو-آی-آر)، مرکز حمایتی، تشخیصی، درمانی و خدمات اجتماعی است و تنها حامی نهاد بیماران نادر در ایران می‌باشد. این نهاد کلیه خدمات را برای انواع بیماری‌های نادر مدیریت می‌کند. یک نمونه از آن، مدیریت بیماری‌های نادر که از مرگبارترین بیماری‌ها در سراسر جهان است می‌باشد (جولی<sup>۱۰</sup> و همکاران، ۲۰۱۱). بنیاد بیماری‌های نادر ایران سالهاست که از سیستم‌های اطلاعاتی برای ارائه انواع خدمات به بیماران و پزشکان بهره می‌برد (سایت رادیور<sup>۱۱</sup>، ۲۰۱۴).

بنابراین، داده کاوی می‌تواند با ارائه تکنیک خود، یکی از راهکارها برای استخراج اطلاعات از مقادیر زیادی از داده‌ها باشد و قادر است کیفیت تصمیم‌گیری مدیریت داده‌ها را بهبود بخشد (اسلامی و همکاران، ۲۰۱۱). هدف اصلی داده کاوی پیش‌بینی است. درخت تصمیم‌گیری، به عنوان یکی از تکنیک‌های داده کاوی ثابت کرده که در میان دیگر تکنیک‌ها، از جمله شبکه‌های عصبی مصنوعی و مدل رگرسیون، دقیق‌ترین پیش‌بینی را ارائه می‌دهد (دلن<sup>۱۲</sup> و همکاران، ۲۰۰۵). داده کاوی پیشگویانه شایع‌ترین نوع داده کاوی و دارای بیشترین آمار برنامه‌های کاربردی تجاری است (کیودسی و همکاران، ۲۰۱۷).

### ۳. روش‌شناسی پژوهش

روش تحقیق این پژوهش داده کاوی می‌باشد. داده کاوی یک روش و تکنولوژی است که از سال ۱۹۹۴ معرفی شده است (ترایبولا<sup>۱۳</sup>، ۱۹۹۷) روش و تکنولوژی داده کاوی حجم زیادی از داده‌ها را به اطلاعات معنی‌دار برای پشتیبانی کارآمد از تصمیم‌گیری تبدیل می‌کند (کوه<sup>۱۴</sup> و همکاران، ۲۰۱۱). در این تحقیق، از الگوریتم طبقه‌بندی داده‌ها (الگوریتم C4.5) برای برای پیش‌بینی عوامل موثر بر بیماری نادر بر اساس چندین شاخص از قبیل سن، جنسیت، مدت زمان اقامت و بیماری استفاده شده است.

و کا (WEKA) فرآیند داده کاوی را با استفاده از درخت تصمیم‌گیری که برای کمک به تصمیم‌گیری سازمان براساس عوامل طبقه بندی شده است، مورد استفاده قرار می‌دهد (ریچا<sup>۱۵</sup>، ۲۰۱۳). داده‌های جمع‌آوری شده از طریق برخی جلسات مصاحبه با آنالیزورهای اطلاعات و مدیران فناوری اطلاعات بنیاد بیماری‌های نادر ایران از طریق تعاملات چهره به چهره، تلفن یا ایمیل به دست آمده است. هدف اصلی این جلسات بحث و بررسی، دستیابی به بینش در مورد نوع داده‌ها برای مرور برخی از احتمالات در خصوص کارآیی داده کاوی است که می‌تواند برای رادو-آی-آر سودمند باشد. شکل ۱ طرح تحقیقاتی را نشان می‌دهد.

<sup>10</sup> Julie

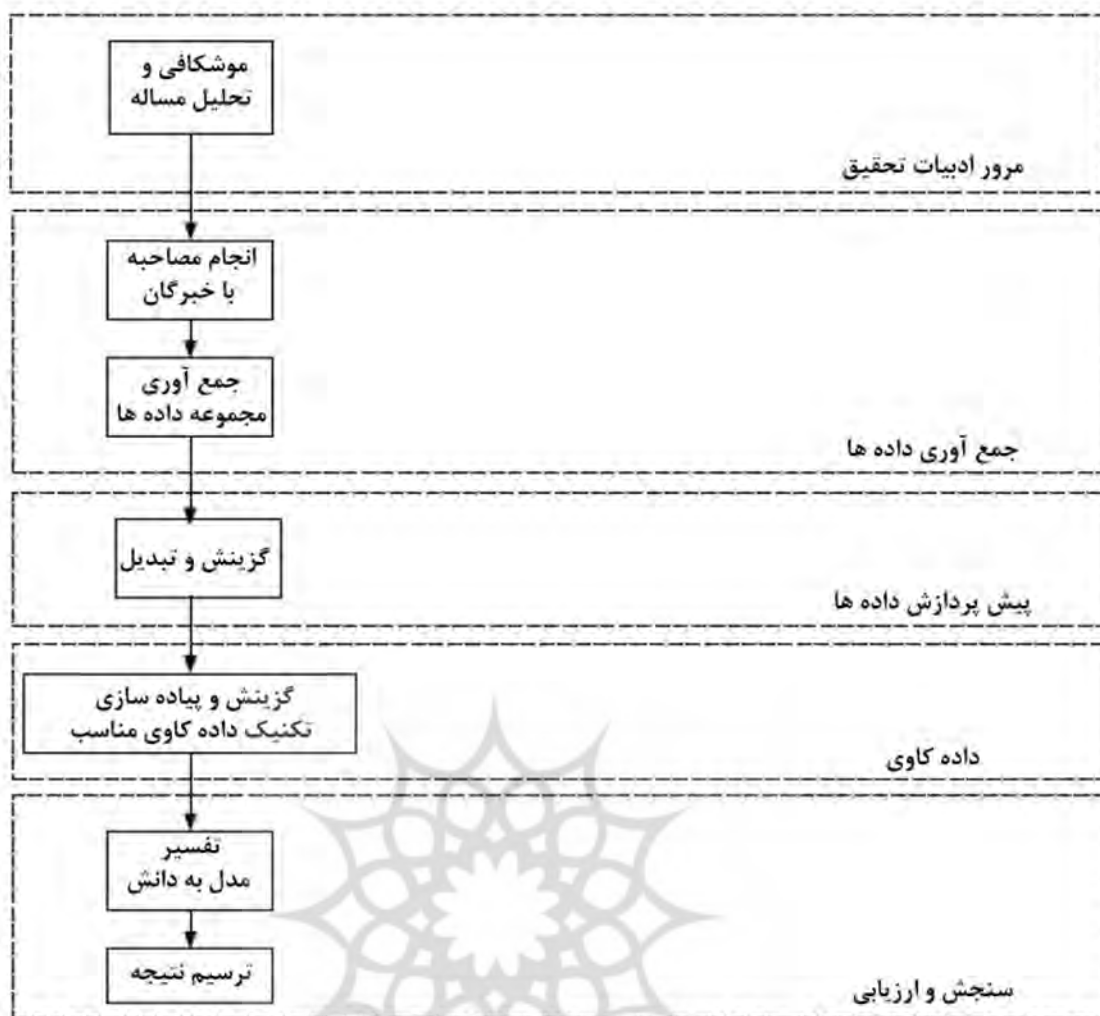
<sup>11</sup> Radoir.org

<sup>12</sup> Delen

<sup>13</sup> Trybula

<sup>14</sup> Koh

<sup>15</sup> Richa



شکل ۱: طرح تحقیق

#### ۴. یافته‌های پژوهش

(۱) جمع آوری و درک اطلاعات داده‌های جمع آوری شده به صورت وضعیت کلینیک درمانی بنیاد بیماری‌های نادر ایران از دوره ۱ فروردین ۱۳۹۸ تا ۱ خرداد ۱۳۹۸، به عنوان داده‌های خام برای این تحقیق مورد استفاده قرار گرفته است. مجموع ۲۳۵۲ گزارش اطلاعات مربوط به بیماران، شامل ۶۰۵ بیمار بستری و ۱۷۴۷ اطلاعات سرپایی است که می‌تواند مورد بررسی و تحلیل قرار گیرند. از ۲۳۵۲ گزارش اطلاعات، ۹۷۶ بیماری نادر و ۱۳۷۶ عدد آن بیماری‌های غیر نادر است. این داده‌ها به عنوان ویژگی‌های ورودی برای پردازش و اجرای الگوریتم درخت تصمیم‌گیری C4.5 در طبقه‌بندی و پیش‌بینی مدل استفاده می‌شود. خروجی مدل پیش‌بینی "بیمار نادر" یا "غیر نادر" است. در این تحقیق، داده‌های خام با چهارچوبی مجزا از علامت (CSV) ذخیره می‌شوند. سه فایل پایگاه داده بیماران بستری از مراکز وابسته به بنیاد، مرکز تصویربرداری پزشکی پارسین، آزمایشگاه بیماران نادر، کلینیک فوق تخصصی بیماران نادر و یک فایل پایگاه داده بیماران سرپایی در یک جدول در مایکروسافت اکسل ادغام شده و به فرمت CSV تبدیل شده تا با ابزار داده‌کاوی و کا (WEKA) سازگار و به راحتی پردازش گردد.

پس از اتمام ادغام داده‌ها، پاک‌سازی داده‌ها انجام شد، داده‌های نامتناسب، داده‌های فاقد ارزش مقدار (فاقد پرونده) و داده‌های بیش از حد (داده‌های حاوی بیش از یک پرونده با همان مقادیر) حذف شده‌اند زیرا وجود آنها می‌تواند کیفیت یا دقت نتایج داده‌کاوی را کاهش دهد.

پس از ادغام و پاک‌سازی داده‌ها، یک جدول چهار ستونی که شامل (سن، جنسیت، مدت زمان اقامت و بیماری) است تعریف می‌شود. جزئیات این ویژگی‌ها در جدول ۱ دیده می‌شود.

جدول ۱ - توصیف ویژگی‌ها در مجموعه داده‌ها			
ویژگی‌ها	توصیف	نوع	مقادیر ممکن
Age	سن بیمار	عددی	(۱ و ۲ و ۳ و ...)
Gender	جنسیت بیمار	اسمی	مرد، زن
Length Of Stay	مدت زمان اقامت بیمار در مراکز درمانی	عددی	۰ روز، ۱ روز، ۲ روز و ...
Disease	آیا بیماری نادر است یا خیر؟	اسمی	نادر، غیر نادر

(۲) تجزیه و تحلیل تجربی

(۱) پیش پردازش داده‌ها

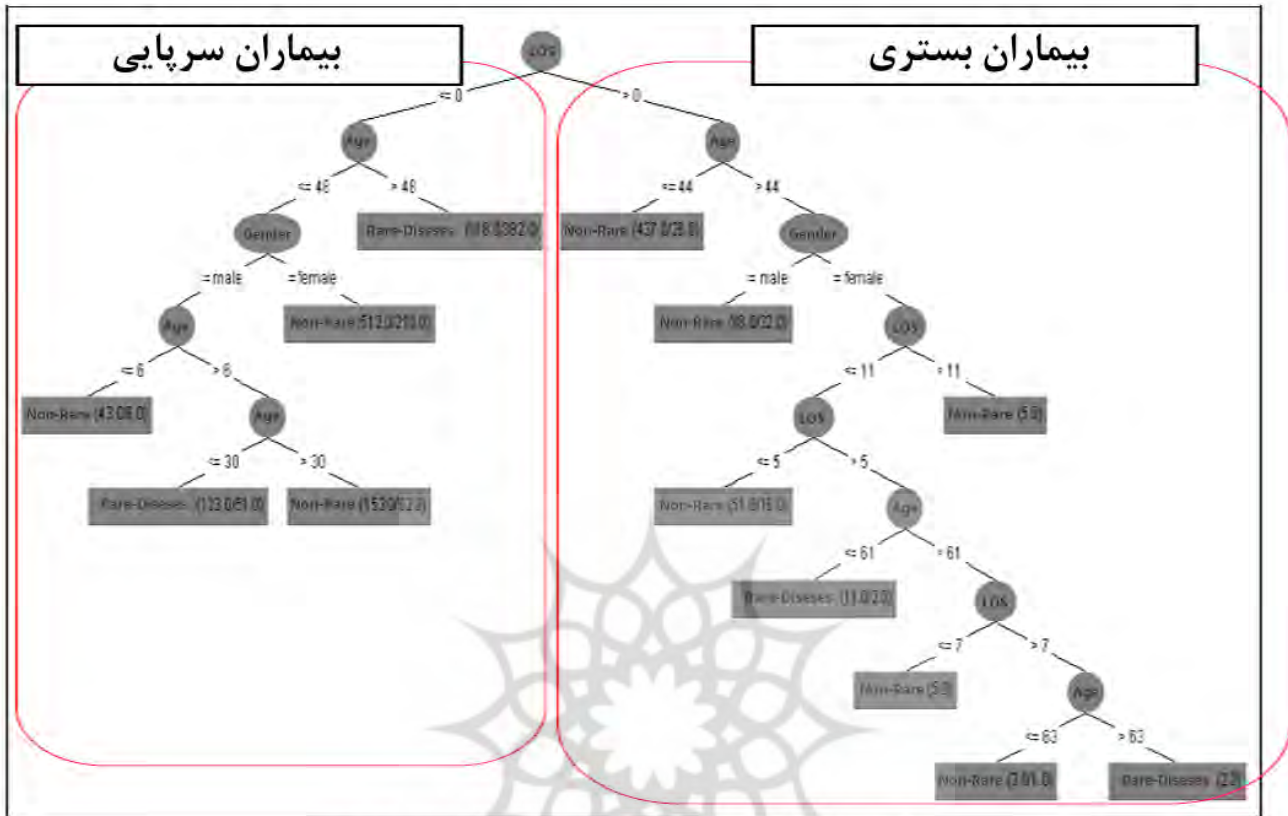
تجزیه و تحلیل تجربی بر روی داده‌های بیماران با استفاده از نرم‌افزار وکا برای به دست آوردن تصویری از درخت تصمیم-گیری و سنجش دقت طبقه‌بندی، ماتریس درهم ریختگی و نرخ TP (شخص بیمار، به درستی بیمار تشخیص داده شود) و نرخ TN (شخص سالم، به درستی سالم تشخیص داده شود) انجام شد. پس از بارگذاری فایل به پنجره اکسپلورر وکا، هیستوگرام تولید می‌شود.

این گراف نشان می‌دهد که تعداد بیماران زن بیشتر از مردان با تعداد ۱۳۴۸ نفر است. بیماران مردی که از بیماری‌های نادر رنج می‌برند تقریباً همانند کسانی می‌باشند که از بیماری‌های مزمن رنج نمی‌برند، درحالی که تعداد بیماران زن مبتلا به بیماری نادر کمتر از افرادی است که از آنها رنج نمی‌برند. همچنین یافته‌ها نشان می‌دهد که حداقل میزان اقامت بیمار در مراکز خدماتی ۱ روز و حداکثر ۳۷ روز است و تعداد بیماران مبتلا به بیماری نادر کمتر از کسانی است که دارای بیماری غیر نادر با تعداد ۹۷۶ بیمار نادر در مقایسه با سال ۱۳۷۶ بیمار غیر نادر می‌باشند.

پس از اتمام فرایند آماده سازی داده‌ها، مدل‌های طبقه‌بندی ساخته شدند. روش درخت تصمیم‌گیری ارائه شده توسط وکا که طبقه بندی الگوریتم C4.5 (J4.8) می‌باشد در مجموعه داده‌ها برای تولید درخت طبقه‌بندی انجام پذیرفته است. داده‌های طبقه‌بندی شده توسط مجموعه داده‌های آموزشی، ارزیابی می‌شود.

درخت طبقه‌بندی تولید شده در شکل ۲ نشان داده شده است. برخی از گره‌ها در شکل ۲ با عنوان مدت زمان اقامت (LOS)، به عنوان مهم‌ترین عامل پیش‌بینی بیماری می‌باشد. گره مدت زمان اقامت (LOS) ریشه درخت تصمیم‌گیری و طبقه‌بندی می‌شود، زیرا از میان دیگر ویژگی‌ها بیشترین اطلاعات را به خود اختصاص داده است. از این رو، قوانین طبقه-

بندی شده از درخت تصمیم گیری در دو قسمت توضیح داده می‌شود، یکی بیماران سرپایی (برای بیمارانی که در یک روز اقامت دارند) و دیگری تحت عنوان بیماران بستری (برای بیماران که بیش از یک روز اقامت دارند) توضیح داده شده است.



شکل ۲. درخت تصمیم گیری که در آن مدت زمان اقامت (LOS) به عنوان مهمترین عامل است

## ۲) خلاصه ارزیابی

عملکرد الگوریتم با استفاده از مصورسازی درخت تصمیم گیری، تجزیه و تحلیل بیماران سرپایی و بستری، دقت طبقه بندی، ماتریس درهم ریختگی و نرخ TP و TN مورد بررسی قرار می‌گیرد. می‌توان این جمع بندی را نمود که طول مدت اقامت (LOS) مهم ترین عاملی است که بر وضعیت بیماری ها تأثیر می‌گذارد و بنیاد را قادر می‌سازد تا طول مدت اقامت بیمار را با توجه به عوامل مختلف مانند سن، جنس و نوع بیماری شناسایی نماید. جدول ۲ عوامل موثر بر بیماری نادر را نشان می‌دهد که می‌توان از آن نتیجه گیری کرد که اکثر بیماران مبتلا به بیماری های نادر سرپایی، دارای سن بیش از ۴۸ سال می‌باشند.

جدول ۲ - عوامل مؤثر بر بیماری های نادر				
نوع	گروه سنی (سال)	جنسیت	مدت زمان اقامت (روز)	تعداد
بیمار سرپایی	۳۰-۷	مرد	۰	۷۲
بیمار سرپایی	بیش از ۴۸	هر دو	۰	۵۳۴
بیمار بستری	۶۱-۴۵	زن	۱۱-۶	۱۲
بیمار بستری	بیش از ۶۳	زن	۱۱-۸	۲

میزان دقت طبقه بندی برای پیش بینی عملکرد، نسبتاً بالاتر از میانگین و در حدود ۶۶,۳۷٪ است. ۱۵۶۱ مورد از ۲۳۵۲ به درستی طبقه بندی شدند که در آن درصد بیماران غیرنادر (۶۸/۴ درصد) بیشتر از بیماران نادر (۶۳/۵ درصد) است. نتیجه تحقیق در مقایسه با نتایج تحقیقات قبلی تفاوت چندانی را در پی نداشت در حالی که نتایج دقت حاصل از انجام آن تحقیقات در حدود ۶۷٪ بوده است؛ به طور مثال هوانگک، بیماری نادر را از طریق آنومالی و بر اساس علائم حیاتی بیماران مانند ECG پیش بینی نموده است (لی و همکاران، ۲۰۰۸).

مقدار آماری کاپا (Kappa) برابر 0.3154 است که ارزیابی میزان سازش میان مقادیر پیش بینی شده توسط مدل می باشد. مقادیر خطای زیر (حاشیه خطای مطلوب، خطای جذر میانگین مربعات، خطای مطلق نسبی و خطای جذر میانگین نسبی مربعات) خطای پیش بینی را تخمین می زند. علاوه بر این، محدوده ROC برای هر دو دسته درخت تصمیم گیری برابر ۰,۷۱۳ است که نشان می دهد اعتبار طبقه بندی بالا است. از آنجا که ROC نزدیکتر به ۱ است، توان تفکیک کننده طبقه بندی بیشتر است.

و کا خطاهای حاصل از طبقه بندی را مشخص می نماید. خطای طبقه بندی، تعداد عوامل مورد نیاز برای تجزیه و تحلیل را که باید به طور گسترده مورد توجه قرار گیرد نشان می دهد. به عنوان مثال، بیماران جوانتر مبتلا به بیماری های غیر نادر نیز وجود داشتند که به احتمال زیاد و به طور اشتباه به عنوان بیمار نادر طبقه بندی شده اند. شاید این مساله به دلیل فقدان اطلاعات بیماران نادر به ویژه برای بیماران جوانتر با توجه به نسبت بین بیماران جوان (کمتر از ۲۵ سال) مبتلا به بیماری های نادر و غیر نادر ۳۸۷:۱۴۰ رخ داده باشد. در نتیجه، الگوریتم برای پیش بینی بیماری های نادر برای بیماران جوانتر اطلاعات کافی را در اختیار ندارد.

داده کاوی به طور معمول یک فرآیند تکراری است که در آن چندین مرحله باید چندین بار تکرار شود (دینی، ۲۰۱۶). در این مورد مطالعاتی، با توجه به دقت که زیاد بالا نیست، پیش پردازش داده ها به منظور افزایش دقت مجدداً انجام می شود، مانند حذف خطا به صورت یک به یک (بر اساس خطای طبقه بندی)، حذف داده های متراکم، تعیین داده های خارج از محدوده و به حداقل رساندن داده های عددی با طبقه بندی سن به کودکان نوپا، خردسالی، اوایل نوجوانی، اواخر نوجوانی، اوایل جوانی، اواخر جوانی، اوایل سالمندی، اواخر سالمندی و کهنسالان (بر اساس طبقه بندی گروه سنی اداره بهداشت جمهوری اسلامی ایران)؛ اما حتی پس از انجام پردازش داده ها برای دومین بار، هنوز هم پارامتر دقت آنچنان که باید، افزایش نیافت.

جدول ۳ ماتریس درهم ریختگی را برای هر طبقه بندی کننده جهت تفسیر و درک نتایج نمایش می دهد.

جدول ۳ - ماتریس درهم ریختگی			
کلاس های واقعی			
ب	الف		
$\frac{FP}{435}$	$\frac{TP}{941}$	الف: بیمار غیر نادر	کلاس های پیش بینی شده
$\frac{TN}{620}$	$\frac{FN}{356}$	ب: بیمار نادر	
۱۰۵۵	۱۲۹۷	جمع کل	J4.8
٪۵۸,۷۶	٪۷۲,۵۵	دقت	

یک ماتریس درهم‌ریختگی، اطلاعاتی را به صورت عددی ارائه می‌دهد که برای توصیف عملکرد مدل طبقه‌بندی بسیار مناسب است و باعث می‌شود مقایسه عملکرد مدل‌های مختلف آسان‌تر شود. علاوه بر این، ماتریس درهم‌ریختگی نه تنها نشان می‌دهد که مدل پیش‌بینی خوبی دارد، بلکه جزئیاتی را در طول فرآیند استخراج اطلاعات که ممکن است اشتباه باشد نمایان می‌سازد (دیوید بو، ۲۰۱۲).

ماتریس درهم‌ریختگی در جدول ۳ نشان داده شده است که در آن ستون‌ها طبقه‌های واقعی و سطرها طبقه‌های پیش‌بینی شده را نشان می‌دهد. در جدول ۳ اعداد ۹۴۱ و ۶۲۰ شاخص تعداد مواردی است که مقادیر واقعی و پیش‌بینی شده مشابه هستند. ۹۴۱ شاخص بیماریانی است که با بیماری غیرنادر پیش‌بینی می‌شوند و واقعا بیماران غیرنادر هستند، در حالی که ۶۲۰ بیماریانی هستند که با بیماری‌های نادر پیش‌بینی می‌شوند و واقعا بیماران نادر هستند؛ به عبارت دیگر، قطر ماتریس همه پیش‌بینی‌های درست را نشان می‌دهد.

از طرفی تعداد ۴۳۵ شاخص مواردی است که نتیجه واقعی بیماری، غیر نادر بوده اما به عنوان یک بیماری نادر پیش‌بینی شده و تعداد ۳۵۶ شاخص مواردی که در آن نتیجه یک بیماری نادر بوده اما به عنوان بیماری غیرنادر پیش‌بینی شده است. لذا می‌توان چنین تجزیه و تحلیل کرد که درصد دقت در رده‌بندی بیماری غیر نادر (۷۲٫۵۵٪) نسبت به بیماری نادر (۵۸٫۷۶٪) بالاتر است.

ضریب اطمینان آزمون پیش‌بینی را می‌توان با محاسبه حساسیت و ویژگی اندازه‌گیری نمود (پاتیل، ۲۰۰۹). در جدول ۳، مقدار TP (مثبت درست) ۹۴۱ و مقدار FP (مثبت نادرست)، ۴۳۵ تعیین شده است. علاوه بر این، ۳۵۶ مورد از FN (منفی نادرست) و ۶۲۰ مورد از TN (منفی نادرست) یافت شده است. سلول‌هایی که با TP (مثبت درست) برچسب‌گذاری شده‌اند، تعداد موارد واقعی است که توسط الگوریتم طبقه‌بندی به‌طور دقیق پیش‌بینی شده است، در حالی که سلول‌های دیگر، FN، FP و TN به روش‌های مشابه تفسیر می‌شوند. عملکرد الگوریتم را نه تنها با محاسبه دقت اجرا، بلکه با حساسیت و ویژگی نیز می‌توان مورد ارزیابی قرار داد.

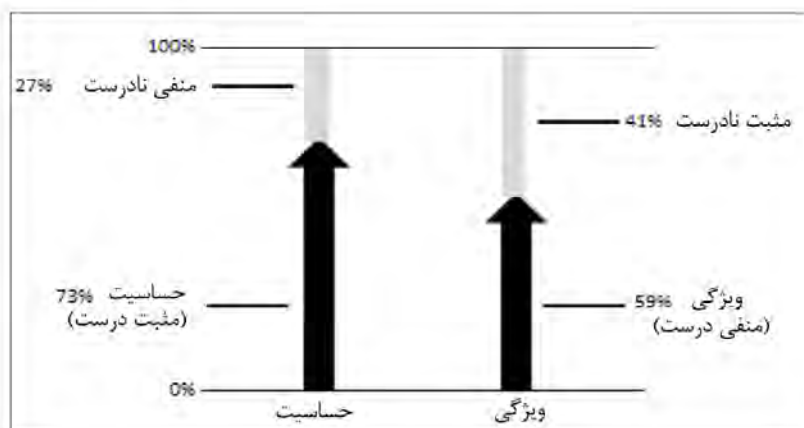
- حساسیت به توانایی آزمون پیش‌بینی در شناسایی درست بیماران مبتلا به بیماری غیر نادر را مشخص می‌کند. حساسیت برابر با نرخ مثبت درست است:

$$\text{حساسیت} = \frac{TP}{TP + FN} \Rightarrow \frac{941}{941+356} = 0.73$$

- ویژگی به توانایی آزمون پیش‌بینی در شناسایی بیماران مبتلا به بیماری نادر اشاره می‌کند. ویژگی برابر با نرخ منفی درست است:

$$\text{ویژگی} = \frac{TN}{TN + FP} \Rightarrow \frac{620}{620+435} = 0.59$$





شکل ۳: نسبت حساسیت به ویژگی

شکل ۳ حساسیت بالا اما ویژگی کم (مثبت نادرست بیشتر) را نشان می دهد که بیان کننده تعداد مواردی است که نتیجه واقعی آن غیر نادر بوده اما به عنوان بیماری نادر پیش بینی شده است. این امر ممکن است ناشی از داده های نامتقارن باشد که در آن تعدادی از داده های غیر نادر بسیار بالاتر از داده های نادر است. مسئله ناپایداری داده ها گسترش یافته و به طور جهانی رخ داده است که بر نتیجه فرایند داده کاوی تأثیر می گذارد (چاولا، ۲۰۰۵). به عنوان مثال، در این تحقیق تعداد بیماران سرپایی بسیار بزرگتر از داده های بیماران بستری با نسبت ۱۷۴۷:۶۰۵ می باشد و تعداد بیماران مبتلا به بیماری های نادر بسیار کمتر از بیماری های غیر نادر با نسبت ۹۷۶:۱۳۷۶ است. با توجه به نتیجه بالا، نرخ TP به حساسیت و از آنجا که نرخ TP به حساسیت برابر ۰,۷۳ است لذا پیش بینی تمامی موارد توسط مدل درست است.

### ۵. بحث و نتیجه گیری

هدف اولیه این تحقیق، پیش بینی عوامل جمعیت شناختی بر بیماری نادر است که بر اساس نادر یا غیر نادر بودن طبقه بندی می شود؛ اما پس از انجام فرایند داده کاوی، طول مدت اقامت (LOS)، به عنوان دانش حاصل از مدل پیش بینی، به عنوان مهمترین عامل که یک بیماری نادر را پیش بینی می کند، ظاهر می شود؛ به عبارت دیگر، طول مدت اقامت (LOS) قوی ترین شاخص پیش بینی بیماری های نادر است که این بنیاد را قادر می سازد تا طول مدت اقامت (LOS) بیماران را با توجه به وضعیت بیماری، جنسیت و سن را پیش بینی نماید؛ بنابراین، ثابت می شود که داده کاوی بیش از یک موضوع آماری می باشد و به عنوان یک ابزار تجزیه و تحلیل پیش بینی مطرح می باشد.

علاوه بر این، الگوریتم C4.5 یک مدل طبقه بندی را در قالب مصورسازی (درخت تصمیم گیری) با موفقیت ایجاد می کند که می تواند به راحتی برای پیش بینی بیماری نادر مورد مطالعه و تفسیر قرار گیرد. دقت پیش بینی آن ۶۶,۳۷ درصد بوده است. با این حال، اندازه حساسیت تا ۰,۷۳ به دست آمد که این بدان معناست که ضریب اطمینان آزمون پیش بینی، خوب است و معتبر بودن داده ها ۱۰۰ درصد ثابت شده اند؛ بنابراین، علیرغم دقت پایین مدل مشتق شده، نتیجه برخی از عواملی را که بنیاد بیماری های نادر می تواند در هنگام پیش بینی هزینه های درمان بیمار مورد توجه قرار دهد، نشان می دهد. به عنوان مثال، مدل نشان می دهد که شاخص جنس یا سن، به شدت تحت تأثیر مدت زمان بستری بیمار قرار می گیرد. بیمار سالخورده با بیماری نادر باید در مراکز درمانی تحت مراقبت و مدت زمان بستری طولانی تری نسبت به افراد جوان تر قرار گیرد. در نتیجه، درخت تصمیم گیری یک روش مفید و آموزنده برای انجام داده کاوی پیش بینی شده است.

مزیت داده کاوی برای مراکز درمانی بیماران نادر را می توان اینگونه بیان کرد؛ دانشی که می تواند از الگوهای ایجاد شده پس از داده کاوی شناسایی شود، احتمالاً می تواند به عنوان یک سیاست جدید برای بنیاد بیماری های نادر ایران (رادو-آی-آر) مورد توجه قرار گیرد. مصورسازی از طریق درخت تصمیم گیری که خواندن و تفسیر را آسان می نماید نیز به افراد شاغل در بنیاد در ترسیم اطلاعات، کمک شایانی را می نمایند. برخی پیشنهادهایی که می تواند در مراکز درمانی بیماران نادر به کار گرفته شود به شرح زیر می باشد؛ اینکه با این روش ممکن است بنیاد بیماری های نادر ایران قادر به شناسایی مدت زمان بستری شدن بیماران مبتلا به بیماری نادر شود که این موضوع می تواند مخاطب را در اختصاص میزان بودجه بیشتر برای خدمات بهداشتی و داروها به بیماران نادر بستری شده در مقابل بیماران سرپایی مبتلا به بیماری های نادر، جهت برآورد نمودن میزان پوشش بیمه ای که باید برای یک بیمار بستری شده در نظر گرفته شود و نیز برآورد نمودن میزان پرداختی هزینه حق بیمه پذیرش شدگان در هر ماه، کمک و راهنمایی نماید. همچنین مشاوره های بهداشتی می بایست صورت گیرند، مخصوصاً برای زنانی که مستعد ابتلا به بیماری های نادر در مقایسه با همتایان مرد خود می باشند، دانش، توانایی، آگاهی از بیماری های نادر و درک پذیرش شدگان می بایست افزایش یابد و حداقل یک بار در سال برای پذیرش شدگان بالای ۴۰ سال، غربالگری پزشکی انجام شود. هدف این مطالعه شناسایی عوامل مخاطره آمیز برای بیماری های نادر است.

قبل از انجام داده کاوی، اطمینان حاصل نمودن از اینکه نسبت نمونه محاسبه شده متعادل باشد حائز اهمیت است. نیازی نیست این نسبت یکی باشد اما برای سیستم محاسبه گر مناسب است تا چگونگی طبقه بندی را برای هر دسته ایجاد شده فرابگیرد. داده های مورد استفاده در این مطالعه عمدتاً توسط یک طبقه غالب از نمونه ها انتخاب شده اند. در نتیجه، داده های نامتعادل منجر به کمبود دقت و طبقه بندی نادرست شده است و در نتیجه تعداد بیماری های غیرنادر بیش از بیماری نادر شده است.

بنیاد می تواند از روش های مختلفی برای افزایش دقت استفاده کند. به عنوان مثال، یکی از روش ها جنگل تصادفی یا ماشین بردار پشتیبانی (SVM) است که یکی از گزینه های موجود برای افزایش دقت می باشد (ادکی، ۲۰۱۲). با این حال، روش های دیگر ممکن است دقت بیشتری داشته باشند، اما تصویری مانند درخت تصمیم گیری را ارائه نمی دهد. به ویژه برای مراکزی که افراد با داده کاوی آشنایی ندارند، استخراج اطلاعات از درخت تصمیم مصور شده، ساده تر است. بر اساس این مطالعه، مشخص شد که داده های دریافتی از مراکز درمانی بنیاد، به لحاظ کیفی مورد انتظار نیست. در نتیجه، یک سری تکنیک های پاکسازی داده اتخاذ گردید. این موضوع نشان داد که عملکرد روش های داده کاوی بستگی به کیفیت داده ها دارد که بنیاد باید این مساله را در نظر بگیرد. داده های مورد استفاده در این تحقیق حاوی بسیاری از مقادیر گم شده، داده های نامناسب یا اطلاعات ناقص است که بر کیفیت داده هایی که باید قبل از انجام فرایند داده کاوی پاک شوند (پاک سازی داده ها) تاثیر می گذارد؛ بنابراین، برای اطمینان از اینکه داده های نگهداری شده پاک شده اند، لازم است برخی از مکانیزم ها مورد استفاده قرار گیرند. لذا داده کاوی می تواند به راحتی داده ها را بدون نیاز به مداخله نوع خاصی از توجیه یا اعمال نظر شخصی که می تواند مانع ادامه فعالیت شود، در صورتی که بنیاد پیش بینی نماید که داده های پاکسازی شده را از ابتدا در اختیار داشته باشد، پردازش نماید.

داشتن یک مدیریت داده خوب، یکی از روش های حفظ کیفیت داده است. کیفیت داده ها می تواند زمانی که رادو-آی-آر آر میل به تامین اعتبار برای روند تجزیه و تحلیل جهت درک شکاف میان سطح بلوغ مدیریت داده ها و بهبودی الگوها و

روش های مدیریت داده ها دارد، بکارگرفته شود. سطح بلوغ مدیریت داده ها می تواند برای کمک به ارزیابی وضعیت فعلی مدیریت داده ها بنیاد استفاده شود. بر اساس یافته های این تحقیق، تاکنون بنیاد بیماری های نادر ایران هیچ فرایندی را در مراکز درمانی خود برای دانستن اینکه آیا داده ها پاک سازی شده اند یا خیر، در اختیار ندارد؛ بنابراین، بر اساس پنج سطح بلوغ کیفیت داده شده توسط روش لری انگلیسی (ادلمن، ۲۰۰۵)، کیفیت داده ها، بلوغ را در سطح دوم نشان نمی دهد که در آن بنیاد نه تنها نیاز به انجام نمایه سازی داده و پاکسازی داده ها دارد بلکه باید پشتیبانی گسترده ای از بنیاد برای بهبود کیفی داده ها صورت پذیرد

## ۶. منابع و مآخذ

1. Adell, A., Ahmadi, P., & Sebt, M. (2010). Desining Model for Choosing human resources with data mining approach. *Journal of Iranian Technology*, 2(4), 5.
2. Bali, R. K. (Ed.). (2005). *Clinical knowledge management: opportunities and challenges*. IGI Global.
3. Bowes, D., Hall, T., & Gray, D. (2012, September). Comparing the performance of fault prediction models which report multiple performance measures: recomputing the confusion matrix. In *Proceedings of the 8th international conference on predictive models in software engineering* (pp. 109-118).
4. Chawla, N. V. (2009). Data mining for imbalanced datasets: An overview. In *Data mining and knowledge discovery handbook* (pp. 875-886). Springer, Boston, MA.
5. Cleary, R. J. (2006). *Applied data mining: statistical methods for business and industry*.
6. Delen, D., Walker, G., & Kadam, A. (2005). Predicting breast cancer survivability: a comparison of three data mining methods. *Artificial intelligence in medicine*, 34(2), 113-127.
7. Edeki, C., & Pandya, S. (2012). Comparative study of data mining and statistical learning techniques for prediction of cancer survivability. *Mediterranean journal of Social Sciences*, 3(14), 49-49.
8. Giudici, P. (2005). *Applied data mining: statistical methods for business and industry*. John Wiley & Sons.
9. Helmerhorst, H. J., Schultz, M. J., van der Voort, P. H., Bosman, R. J., Juffermans, N. P., de Wilde, R. B., ... & de Keizer, N. F. (2016). Effectiveness and clinical outcomes of a two-step implementation of conservative oxygenation targets in critically ill patients: a before and after trial. *Critical care medicine*, 44(3), 554-563.
10. Julie, D. M., & Kannan, B. (2012). *Statistical Machine Learning Techniques for the Prediction of Learning Disabilities in School-Age Children* (Doctoral dissertation, Cochin University of Science and Technology).
11. Koh, H. C., & Tan, G. (2011). Data mining applications in healthcare. *Journal of healthcare information management*, 19(2), 65.
12. Li, J., Huang, K. Y., Jin, J., & Shi, J. (2008). A survey on statistical methods for health care fraud detection. *Health care management science*, 11(3), 275-287.
13. Moss, L. T., Abai, M., & Adelman, S. (2005). *How to improve data quality*. Data Strategy; Addison-Wesley Professional: Boston, MA, USA.
14. Patil, B. M., Toshniwal, D., & Joshi, R. C. (2009, March). Predicting burn patient survivability using decision tree in weka environment. In *2009 IEEE International Advance Computing Conference* (pp. 1353-1356). IEEE.
15. Piatetsky-Shapiro, G. (2000). Knowledge discovery in databases: 10 years after. *Acm Sigkdd Explorations Newsletter*, 1(2), 59-61.
16. Price, S., Dobbs, D., Oliveira, J., Beidas, S., Burkhart, T., & Sharp, J. (2013). *Clinical & Business Intelligence: Data Management—A Foundation for Analytics*.
17. Qudsi, D. H. (2016). Predictive Analytics Data Mining in Imbalanced Medical Dataset. *Jurnal Komputer Terapan*, 2(2), 195-204.
18. Qudsi, D. H., Kartiwi, M., & Saleh, N. B. (2017). Predictive data mining of chronic diseases using decision tree: a case study of health insurance company in Indonesia. *International Journal of Applied Engineering Research*, 12(7), 1334-1339.

19. Rogers, G., & Joyner, E. (1997, March). Mining your data for health care quality improvement. In SAS User Group International Conference (pp. 641-647).
20. Sharma, R., Ghosh, A., & Joshi, P. K. (2013). Decision tree approach for classification of remotely sensed satellite data using open source support. *Journal of Earth System Science*, 122(5), 1237-1247.
21. Trybula, W. J. (1997). Data mining and knowledge discovery. *Annual review of information science and technology (ARIST)*, 32, 197-229.
22. Yousofi, M. H., Esmaili, M., & Sharifian, M. S. (2016). A study on image mining; its importance and challenges. *American Journal of Software Engineering and Applications*, 5(3-1), 5-9.



# Decision Tree-Based Use of Data Mining Techniques to Determine & Predict the Rare Patients Status

## The Case Study: Rare Patients Registry System (SABNA)

Mohsen Sojoudi<sup>1</sup>  
Fariba Ebrahimbabaie<sup>\*2</sup>

Date of Receipt: 2020/07/15 Date of Issue: 2020/07/18

### Abstract

The health database contains a broad range of clinical data through which the discovered communication patterns and algorithms are led to new medical data achievements. Nowadays due to the existing integrated data as well as the informatics technology growth, it has turned into be a significant emergence. This study aims to identify the potential advantages with which data mining techniques can be applied for health and treatment trends using RADOIR's patients' data as a case study. The most usual method is WEKA based data pre-processing, data mining and classification with decision tree to create prediction modalities via picturing a tree in order to analyze a rare disease. Result & Also we have used WEKA to evaluate the prediction processing via measuring and analyzing the sensitivity rate. As a result, there are some factors which a patient supporting center can take them into consideration while predicting the patients' treatment costs and expenses.

### Keyword

Data mining, Prediction Model, Rare Disease, Decision Tree, Sensitivity Analysis, WEKA

1. Operation Management: Research (PhD) Ferdowsi University of Mashad – Iran ([mohsened@gmail.com](mailto:mohsened@gmail.com))
2. Phd student of psychology, Lorestan University, Iran ([faribaebrahimbabaie@yahoo.com](mailto:faribaebrahimbabaie@yahoo.com))

پژوهشگاه علوم انسانی و مطالعات فرهنگی  
پرتال جامع علوم انسانی