

## Discovering Hidden Cluster Structures in Patients with Cirrhosis Based on Laboratory Data

Mina Pazouki<sup>1</sup>, Mohammad Mehdi Sepehri<sup>2</sup>, Mehdi Saberifiroozi<sup>3</sup>

<sup>1</sup> MSc, Department of Industrial Engineering, South branch Islamic Azad University, Tehran, Iran.

<sup>2</sup> Associate Professor, Department of Industrial Engineering, Faculty of Engineering, Tarbiat Modares University, Tehran, Iran

<sup>3</sup> Professor, Digestive Disease Research Center, Shariati Hospital, Tehran University of Medical Sciences, Tehran, Iran.

### ABSTRACT

#### **Background:**

Liver cirrhosis was one of the most important liver diseases. Other chronic liver diseases could be lead to liver cirrhosis. Liver cirrhosis could be lead one kind of liver cancers named hepatocellular carcinoma. Cirrhosis in the early stages just by laboratory and imaging testes could be diagnosed. In this study cirrhotic patients were classified based on laboratory symptoms. For this purpose data mining approach has been used in this research. Data mining was an interdisciplinary science that discovers the hidden knowledge in the data.

#### **Materials and Methods:**

We used K-Means algorithm to categorize the statues of cirrhotic patients. In order to determine the quality of clustering results and to find the best number of clusters, we have used silhouette indices.

#### **Results:**

Our data consists of 410 records which have been collected from Dr. Shariati hospital The number of features in this study are 11 items and sampling were divided into two main groups.

#### **Conclusion:**

At first, we have done clustering based on 21 attributes and the average silhouette was 41 percent. We improved the model, in order to reach a reasonable structure. Finally, based on 7 attributes, a reasonable clustering model was derived. The new model provides 64 percent average silhouette, and based on patients' status, patients are divided into 2 main categories. The risk of HCC in the first cluster is 23 percent and in the second cluster is 14 percent.

**Keywords:** Liver cirrhosis; Data mining; clustering; K-means; Silhouette

*please cite this paper as:*

Pazoki M, Sepehri MM, Saberifiroozi M. Discovering hidden cluster structures in patients with cirrhosis based. *Govaresh* 2014;18:191-7.

#### **Corresponding author:**

Mohammad Mehdi Sepehri, MD

Department of Industrial Engineering, Tarbiat

Modares University (TMU), Jalal-e Al-e Ahmad

Highway, Tehran 1411713114, Iran

Telefax: + 98 21 82883379

E-mail: mehdi.sepehri@gmail.com

Received: 05 May 2014

Edited: 11 Aug. 2014

Accepted: 12 Aug. 2014

# کشف ساختارهای خوشه‌ای پنهان در بیماران مبتلا به سیروز کبدی بر پایه نشانه‌های آزمایشگاهی

مینا پازوکی<sup>۱</sup>، محمد مهدی سپهری<sup>۲</sup>، مهدی صابری فیروزی<sup>۲</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد، دانشکده مهندسی صنایع، دانشگاه آزاد اسلامی واحد تهران جنوب، تهران، ایران  
<sup>۲</sup> دانشیار، بخش مهندسی صنایع، دانشکده فنی و مهندسی، دانشگاه تربیت مدرس، تهران، ایران  
 استاد، مرکز تحقیقات گوارش و کبد، بیمارستان دکتر شریعتی، دانشگاه علوم پزشکی تهران، تهران، ایران

## چکیده

### زمینه و هدف:

بیماری سیروز کبدی مهم ترین بیماری کبدی است که سایر بیماری های کبدی به صورت مزمن می‌توانند زمینه ساز بروز آن باشند و خود می‌تواند منشا ابتلا به سرطان کبد شود. این بیماری درمان قطعی ندارد و در مراحل اولیه هیچ گونه نشانه‌ای ندارد و تنها بر اساس آزمایش ها و تصویر برداری ها قابل تشخیص است. در این پژوهش جهت طبقه بندی وضعیت این بیماران بر اساس نشانه‌های آزمایشگاهی از رویکرد داده کاوی استفاده شده است. داده کاوی یک علم میان رشته ای است که به کشف دانش پنهان در داده‌ها می‌پردازد.

### روش بررسی:

در این مطالعه با استفاده از رویکرد خوشه بندی و الگوریتم k-means وضعیت بیماران را دسته بندی کردیم و برای تعیین کیفیت خوشه بندی و پیدا کردن بهترین تعداد خوشه از شاخص سیلوئت استفاده شد.

### یافته‌ها:

داده های ما شامل ۴۱۰ مورد بیماران مبتلا به سیروز کبدی می‌باشد که از بیمارستان دکتر شریعتی جمع آوری شده اند. تعداد مشخصه های مورد بررسی شامل ۲۱ مورد می‌باشد. در این مطالعه وضعیت بیماران مبتلا به سیروز کبدی به دو دسته اصلی تقسیم شد.

### نتیجه گیری:

خوشه بندی بر اساس ۲۱ مشخصه انجام شد و میانگین معیار سیلوئت ۴۱ درصد شد. اما با کاهش تعداد مشخصه ها در نهایت با در نظر گرفتن ۷ مشخصه به کیفیت خوشه بندی قابل قبول با متوسط معیار سیلوئت ۶۴ درصد رسیدیم و وضعیت بیماران بر اساس نشانه‌های آزمایشگاهی به دو دسته تقسیم شد که احتمال ابتلا به سرطان در خوشه اول ۲۳ درصد و در خوشه دوم ۱۴ درصد می‌باشد.

**کلید واژه:** سیروز کبدی، داده کاوی، خوشه بندی، سیلوئت، K-means

گوارش/ دوره ۱۹، شماره ۳/ پاییز ۱۳۹۳/ ۱۹۷-۱۹۱

### زمینه و هدف:

داده کاوی<sup>۱</sup> به عنوان رویکرد کشف دانش در داده ها می‌کاود تا گنج پنهان دانش را پیدا کند. داده کاوی در اواخر سال ۲۰۰۹ به عنوان موضوع

#### 1. Data mining

### نویسنده مسئول: محمد مهدی سپهری

تهران، بزرگراه جلال آل احمد، پل نصر، دانشگاه تربیت مدرس،

کد پستی ۱۴۱۱۷۱۳۱۱۴

تلفن و نامبر: ۰۲۱-۸۲۸۸۳۳۷۹

پست الکترونیک: mehdi.sephri@modares.ac.ir

تاریخ دریافت: ۹۳/۲/۱۵

تاریخ اصلاح نهایی: ۹۳/۵/۲۰

تاریخ پذیرش: ۹۳/۵/۲۱

تحقیقات پزشکی در سال ۲۰۱۰ ثبت شد. امروزه یکی از مهم ترین کاربرد های داده کاوی در زمینه پزشکی است و به پژوهشگران کمک می‌کند که به بینش عمیقی از مجموعه داده های پزشکی دست پیدا کنند. (۱) در سال های اخیر داده کاوی علاوه بر این که توانسته است در زمینه های مختلف پزشکی در فرایند شناخت بیماری، تشخیص بیماری، پیش بینی روند بیماری، انتخاب روش درمان، پیش بینی طول مدت اقامت در بیمارستان، ارتباط بین بیماری ها و ارتباط بین داروها کمک شایانی به ارائه دهندگان خدمات بهداشتی و درمان نماید. (۷-۲) به دلیل افزایش کیفیت خدمات بهداشتی و درمانی باعث کاهش ریسک های ناشی از تشخیص غلط بیماری و یا روش درمان و در نهایت کاهش هزینه و بهبود وضعیت سلامت بیماران شده است.

بر حسب این که در فرایند داده کاوی استنتاج چه نوع دانشی از مجموعه آموزشی مورد نظر است از روش های مختلف داده کاوی می‌توان

بهره جست. این روش ها از نظر شیوه یادگیری به دو دسته اصلی تقسیم می شوند. (۸)

الف- الگوریتم های یادگیری با نظارت

ب- الگوریتم های یادگیری بدون نظارت

در روش های داده کاوی مبتنی بر الگوریتم های یادگیری با نظارت هدف از داده کاوی مشخص است و در حقیقت تحلیل گر می داند که به دنبال چه دانشی می گردد. بر خلاف آن در روش های مبتنی بر یادگیری بدون نظارت، هدف از کشف دانش کاملا تعریف شده نیست. (۹)

خوشه بندی یک تکنیک داده کاوی است که از رویکرد یادگیری غیر نظارتی برای تحلیل داده ها استفاده می کند. (۱۰) داده ها بر اساس اصل حداکثر کردن شباهت داخل گروه ها و حداقل کردن شباهت بین گروه ها، خوشه بندی می شوند. خوشه بندی با پیدا کردن روابط پنهان در داده ها که برای پیدا کردن الگوها و پیش بینی وقایع آینده یا تبدیل رکورد داده ها به خوشه های معنی دار مفید هستند، به ارزش پایگاه داده های موجود می افزاید. (۱۱)

موضوع مورد بحث در این پژوهش از حوزه بهداشت و سلامت، بیماری سیروز کبدی<sup>۱</sup> است. سیروز کبدی یکی از جدی ترین بیماری های کبدی است که تقریباً تمام بیماری های کبدی اگر منجر به از بین رفتن سلول های کبدی شوند می توانند به سیروز تبدیل شوند. (۱۲) این بیماری در اثر التهاب مزمن بافت کبد ایجاد می شود که به دلایل گوناگونی اتفاق می افتد. هر گونه التهابی به صورت کم یا زیاد باعث تخریب بافت کبد می شود. در اثر ادامه این فرایند به طور مزمن و مکرر، ساختار منظم و یک شکل کبد در هم می ریزد و به تدریج بافت کبد سفت می شود. فعالیت های احیا کننده تحت سیروز کبدی ادامه می یابند ولی سرعت پیشرفت از دست دادن سلول ها بیش از جایگزینی آنها است. (۱۳) علل عمده سیروز کبدی شامل هیپاتیت B، هیپاتیت C و هیپاتیت اتوایمیون<sup>۲</sup> است. ولی داروها و مواد سمی، کبد چرب، بیماری ویلسون، تنگی مجاری صفراوی، مشکلات قلبی و انسداد وریدهای کبدی از علل دیگر این بیماری می باشند. (۱۲) در مراحل اولیه، بیماری هیچ گونه نشانه ای ندارد و تنها با آزمایش های دقیق و بررسی های رادیولوژیکی و بافت شناسی قابل تشخیص است. (۱۴) به تدریج که بیماری پیشرفت می کند نشانه ها و عوارض بیماری نیز افزایش می یابد. از مهم ترین نشانه های سیروز کبدی آسیت<sup>۳</sup> و تجمع مایع در شکم، آنسفالوپاتی کبدی<sup>۴</sup>، واریس مری<sup>۵</sup> و خونریزی واریس<sup>۶</sup> و یا عوارض دیگری که بر روی سیستم قلب و عروق، ریه ها و کلیه ها ایجاد می کند، می باشد. (۱۵) درمان بیماری سیروز کبدی بر اساس درمان بیماری زمینه ای آن می باشد و درمان قطعی ندارد (۱۶).

سیروز کبدی می تواند منجر به ابتلا به سرطان کبد شود که هپاتوسلولار کارسینوما (HCC) رایج ترین نوع آن می باشد. بر اساس گزارش آژانس بین المللی تحقیقات سرطان، HCC در مردان پنجمین سرطان شایع در دنیا است و در زنان هفتمین سرطان شایع است. (۱۷) این سرطان سومین علت شایع مرگ و میر در اثر سرطان در دنیا می باشد. (۱۸) میزان ابتلای بیماران سیروویک به HCC در غرب سالانه بین ۲ تا ۳ درصد و در خاور دور بین ۶ الی ۱۱ درصد می باشد. (۱۹)

در این پژوهش به کمک رویکرد خوشه بندی به دنبال دسته بندی وضعیت بیماران سیروز کبدی بر اساس نشانه های آزمایشگاهی آن ها هستیم. که می تواند به درک پزشکان از وضعیت بیماری بیماران بر اساس چند مورد از نشانه های آزمایشگاهی و تصمیم گیری آن ها در فرایند درمان تاثیر به سزایی داشته باشد.

### روش بررسی:

#### معرفی داده ها:

کل داده های جمع آوری شده شامل ۴۱۰ رکورد است که از بیماران مبتلا به سیروز مراجعه کننده به بیمارستان دکتر شریعتی جمع آوری شده است. داده های جمع آوری شده مربوط به سال های ۸۹ تا ۹۲ می باشند و ۵۰ مورد از این بیماران در طی این ۴ سال به HCC مبتلا شدند. تعداد مشخصه ها ۲۱ مورد است که در جدول ۱ معرفی شده اند.

#### آماده سازی داده و پیش پردازش:

اغلب به دلیل خطاهای عملیاتی و پیاده سازی سیستم ها، داده های به دست آمده پر غلط، ناقص و ناسازگارند. بنابراین نیاز است که در ابتدا این داده ها تمیز شوند که این کار با اجرای مراحل پیش پردازش داده ها و در طی مراحل زیر انجام شد:

#### الف) پاکسازی داده ها:

پاکسازی داده در واقع مرحله کنترل کیفی قبل از تحلیل داده است. در این پژوهش قبل از انجام عملیات پیش پردازش کلیه داده ها از نظر صادق بودن با مقادیر امکان پذیر کنترل شده اند و مقادیر اشتباه اصلاح شده اند. و هم چنین رکوردهایی از داده ها که بسیاری از مشخصه های آن ها خالی بود و یا اعداد اشتباه وارد شده بودند حذف شدند و در این مرحله تعداد کل رکوردهای ما به ۳۶۴ مورد رسید.

#### پُر کردن مقادیر مفقوده:

در این مرحله برای پر کردن مقادیر از دست رفته از الگوریتم دسته بندی k نزدیک ترین همسایگی<sup>۷</sup> استفاده شد و بر اساس تعداد ۳ نزدیکترین همسایگی و مد آن ها، مقادیر مفقود شده پر شدند. علت انتخاب عدد ۳ این بود که انتخاب تعداد کمتر ممکن است روش را متاثر

1. Liver Cirrhosis
2. Autoimmune hepatitis
3. Ascites
4. Encephalopathy
5. Varices
6. Varical bleeding

جدول ۱: مشخصه‌های مورد بررسی

شماره	۱	۲	۳	۴	۵	۶	۷	۸	۹
مشخصه	Age	Height	Weight	BMI	Platelet	WBC	AST	ALT	Serum Albumin
شماره	۱۰	۱۱	۱۲	۱۳	۱۴	۱۵	۱۶	۱۷	۱۸
مشخصه	Alkph	BiLirubinT	BilirubinD	Creatinine	BUN	INR	FBS	TG	Cholesterol
شماره	۱۹	۲۰	۲۱						
مشخصه	AFP	Na	k						

جدول ۲: تفسیر مقادیر میانگین معیار سیلوئت (۱۶)

میانگین معیار سیلوئت	تفسیر
۰/۷۱-۱	ساختار قوی
۰/۵۱-۰/۷	ساختار منطقی (مناسب)
۰/۲۵-۰/۵	ساختار ضعیف
<۰/۲۵	هیچ ساختار قابل توجهی وجود ندارد

### الگوریتم K-means:

پارامتر  $k$  (تعداد خوشه) را به عنوان ورودی گرفته و مجموعه  $n$  شیء را به  $k$  خوشه افراز می‌کند. به طوری که سطح شباهت داخلی خوشه‌ها بالا بوده و سطح شباهت اشیا بین خوشه‌ها پایین باشد. شباهت هر خوشه نسبت به متوسط اشیا آن خوشه سنجیده شده که این متوسط مرکز خوشه نیز نامیده می‌شود. (۱۶)

### اعتبار سنجی خوشه‌ها:

شاخص‌های اعتبارسنجی برای سنجش میزان صحت نتایج خوشه‌بندی به منظور مقایسه بین روش‌های خوشه‌بندی مختلف یا مقایسه‌ی نتایج حاصل از یک روش با پارامترهای مختلف مورد استفاده قرار می‌گیرند. دو معیار پایه‌ی اندازه‌گیری پیشنهاد شده برای ارزیابی و انتخاب خوشه‌های بهینه عبارتند از:  
**تراکم:** داده‌های متعلق به یک خوشه بایستی تا حد ممکن به یکدیگر نزدیک باشند. معیار رایج برای تعیین میزان تراکم داده‌ها واریانس داده‌ها است.  
**جدایی:** خوشه‌ها خود بایستی به اندازه کافی از یکدیگر جدا باشند. در این پژوهش از شاخص اعتبارسنجی سیلوئت<sup>۳</sup> استفاده شده است که در ادامه مختصراً تشریح می‌شود:

**شاخص سیلوئت:** یکی از معیارهای متداول اعتبارسنجی خوشه‌بندی است. و دو معیار فواصل درون خوشه‌ای و برون خوشه‌ای را همزمان در نظر می‌گیرد. (۱۷) جدول ۲ تفسیر مقادیر مختلف معیار سیلوئت را نشان می‌دهد:

### یافته‌ها:

در ابتدا مدل‌سازی بر اساس ۲۱ مشخصه جدول ۱-۳ و بر اساس الگوریتم K-means انجام شد. شکل ۱ نتایج با توجه به نمودار سیلوئت در شکل ۱، در این مرحله متوسط کیفیت خوشه‌مدل خوشه‌بندی را

1. Compactness
2. Separation
3. Silhouette

از داده‌های مغشوش کرده و مقدار بزرگتر آن هم ممکن است رفتارهای نزدیک‌تر را در نظر نگیرد. (۸)

### شناسایی نقاط پرت:

جهت شناسایی نقاط پرت از الگوریتم خوشه‌بندی K-means استفاده شد. و داده‌ها به ۴ خوشه تقسیم شدند و ۲۰ داده‌ی پرت شناسایی و حذف شدند. در این مرحله تعداد داده‌های ما به ۳۴۴ مورد رسید.

### (ب) یکپارچه‌سازی داده:

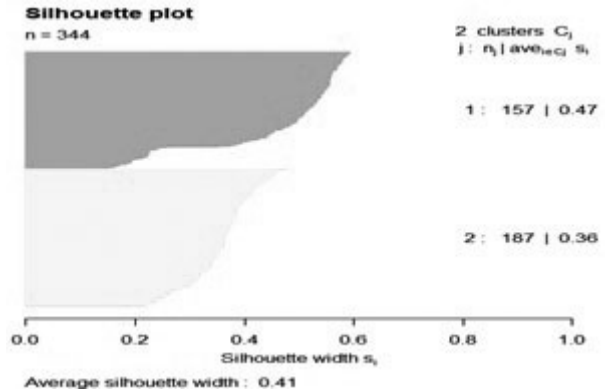
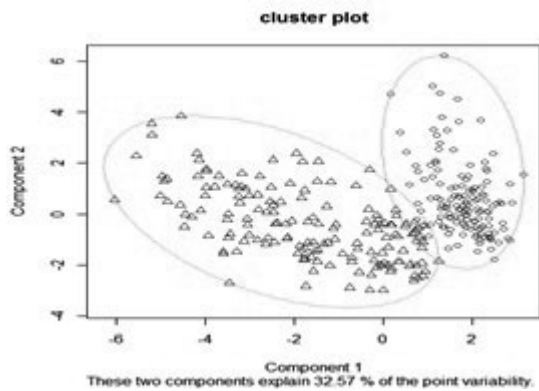
این عملیات شامل یکپارچه‌سازی چندین پایگاه داده است. که در این جا داده‌های مربوط به بیمارستان شریعتی و بیمارستان امام خمینی در یک فایل اکسل مرتب شدند.

### (ج) تبدیل داده:

در این مرحله فرایند نرمال‌سازی روی داده‌ها انجام شد. نرمال‌سازی تغییر مقیاس داده‌ها به گونه‌ای است که آن‌ها را به یک فاصله‌ی کوچک و معین نگاشت می‌کند و باعث می‌شود که داده‌ها با مقیاس بزرگ نتایج را به سمت خود منحرف نکنند. در این پژوهش با استفاده از روش Min-Max (۸) نرمال‌سازی داده‌های کمی که شامل مشخصه‌های نتایج آزمایشگاهی بود انجام شد و داده‌ها در بازه عددی بین صفر و یک قرار گرفتند.

### مدل‌سازی:

روش مدل‌سازی در این پژوهش بر اساس الگوریتم افرازی K-means می‌باشد.

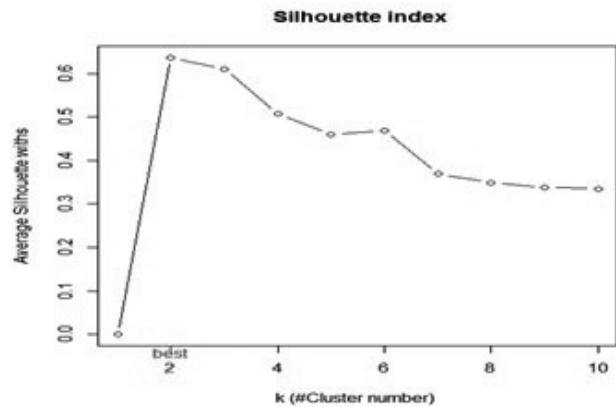


شکل ۱: الف) نمودار معیار سیلوئت برای دو خوشه ها، ب) شکل خوشه ها.

جدول ۳: مشخصه های مربوط به مدل با ۷ مشخصه

شماره	۱	۲	۳	۴	۵	۶	۷
مشخصه	AST	ALT	Serum Albumin	Creatinine	AFP	Na	K

ها را به دو دسته طبقه بندی کنیم. جدول ۲-۴ متوسط مقدار هر مشخصه را در هر خوشه نشان می دهد. براساس نتایج فوق وضعیت بیماران به دو دسته تقسیم شد که با توجه به جدول ۴ به طور متوسط خوشه ی ۱ وضعیت بدتری را نسبت به خوشه ی ۲ نشان می دهد و ۷ مشخصه ی، K، Na، AFP، ALT، AST، Serum Albumin و Creatinine در بین نشانه های آزمایشگاهی در نظر گرفته شده برای تعیین وضعیت بیماران از اهمیت بالاتری برخوردارند.



شکل ۲: نمودار روند میزان تغییرات معیار سیلوئت برای تعداد خوشه های مختلف.

### بحث:

بیماری سیروز کبدی مهم ترین بیماری کبدی است که سایر بیماری های کبدی به صورت مزمن می توانند زمینه ساز بروز آن باشند و خود می تواند منشا ابتلا به سرطان کبد شود. در این پژوهش جهت بررسی مهم ترین نشانه های آزمایشگاهی در تعیین وضعیت بیماران مبتلا به سیروز کبدی و دسته بندی وضعیت آن ها به کمک رویکرد های غیر نظارتی و بر اساس رفتار داده ها، از رویکرد خوشه بندی استفاده کردیم. برای این منظور نتایج آزمایش مربوط به ۴۱۰ بیمار مبتلا به سیروز کبدی در نظر گرفته شد و وضعیت بیماران مبتلا به سیروز کبدی را بر اساس نشانه های آزمایشگاهی آنها و رویکرد خوشه بندی که یک رویکرد کاملاً غیر نظارتی است تقسیم بندی نمودیم. کلیه نشانه های آزمایشگاهی در نظر گرفته شده در این مطالعه بر اساس نتایج آزمایش اولین مراجعه بیماران می باشد ولی آنزیم های کبدی TSA و TLA به علت این که در اثر مصرف دارو ها و بیماری های عفونی و بیماری های ایسکمیک و... مقادیرشان تغییر می کند،

نشان می دهد. ها ۴۱ درصد است. که بر اساس جدول ۲-۳ این خوشه بندی ضعیف است.

نمودار شکل ۲ روند تغییرات معیار سیلوئت با در نظر گرفتن تعداد خوشه های مختلف را مشخص می کند که بهترین تعداد خوشه در این مرحله ۲ می باشد. با کاهش تعداد مشخصه ها و حذف مشخصه های ناسازگار مدل خوشه بندی بهبود پیدا کرد و در نهایت به مدلی با ۷ مشخصه ی معرفی شده در جدول ۳ رسیدیم. نتایج مدل خوشه بندی در این مرحله در شکل ۳ آورده شده است:

بر اساس نمودار سیلوئت در شکل ۳ برای تعداد دو خوشه، متوسط کیفیت خوشه بندی ۶۴ درصد است و در این مدل، خوشه بندی ما در محدوده ی قابل قبول قرار می گیرد. شکل ۴ نمودار روند معیار سیلوئت با در نظر گرفتن تعداد خوشه های مختلف را نشان می دهد. بر اساس شکل ۴ بهترین تعداد خوشه در این مرحله ۲ می باشد. بنابراین با در نظر گرفتن نشانه های آزمایشگاهی بیماران مبتلا به سیروز کبدی توانستیم وضعیت آن



HCC و هم چنین این که مقادیر آنزیم های AST و ALT در اثر عوامل مختلف بسیار تغییر پذیر می باشند و این که نتایج آزمایش بیماران مربوط به آزمایشگاه های مختلف است و می تواند خطاهایی را به همراه داشته باشد، ضعف هایی دارد. از نقاط قوت این پژوهش می توان به این که اولین پژوهشی است که به کاربرد داده کاوی در بیماری سیروز کبدی پرداخته است و این که در ایران اولین پژوهشی است که به طور متوالی داده های مربوط به چهار سال بیماران سیروتیک را مورد بررسی قرار داده است، اشاره کرد.

. گرچه انجام این آزمایش ها تاثیر فوری برای مراقبت از بیماران دارد ولی تاثیر این آزمایش ها در عاقبت بیماران و مراقبت طولانی مدت آنها کاملا مشخص نمی باشد. بنابراین در این پژوهش هدف تنها بررسی وضعیت بیماران بر اساس نشانه های آزمایشگاهی بوده است و هم چنین کشف غیر نظارتی نشانه های مهم در تعیین وضعیت بیماران و باید صحت نتایج آن در مورد بیماران جدید که به درمانگاه ها مراجعه می کنند بررسی شود. این مطالعه، یک مطالعه مقدماتی است که می تواند در تحقیقات آتی با در نظر گرفتن نشانه های بالینی بیماران توسعه داده شود. این پژوهش به علت کم بودن تعداد داده ها و تعداد بیماران مبتلا به

## REFERENCES

1. Yoo I, Alafairet P, Marinov M, Pena-Hernandez K, Gopidi R, Chang JF, et al. Data mining in healthcare and biomedicine: a survey of the literature. *J Med Syst* 2012;36:2431-48.
2. Anunciação O, Gomes BC, Vinga S, Gaspar J, Oliveira AL, Rueff J. A Data Mining Approach for the Detection of High-Risk Breast Cancer Groups. *Adv Bioinformatics* 2010;74:43-51.
3. Gauthier E, Brisson L, Lenca P, Ragusa S. Breast cancer risk score: a data mining approach to improve readability. *International Conference on Data Mining* 2011:15-21.
4. Bhatla N, Jyoti, K. An Analysis of Heart Disease Prediction using Different Data Mining Techniques. *International J Engineering* 2012;1:8.
5. Agrawal A, Misra S, Narayanan R, Polepeddi L, Choudhary A. A lung cancer outcome calculator using ensemble data mining on SEER data. In Proceedings of the Tenth International Workshop on Data Mining in Bioinformatics 2011:5.
6. Kaur AR. A Study of Detection of Lung Cancer Using Data Mining Classification Techniques. *International J* 2013; 3(3).
7. Ahmed K, Jesmin T, Zamilur Rahman M. Early Prevention and Detection of Skin Cancer Risk using Data Mining. *International J Computer Applications* 2013;62:1-6.
8. Han J, Kamber M, Pei J. Data mining: concepts and techniques. *Morgan kaufmann* 2006.
9. Marinakis Y, Marinaki M, Matsatsinis N. A stochastic nature inspired metaheuristic for clustering analysis. *International J Business Intelligence Data Mining* 2008;3:30-44.
10. Wagstaff K, Cardie C, Rogers S, Schrödl S. Constrained k-means clustering with background knowledge. *ICML* 2001 ;1:577-84.
11. Suh SC, Saffer S, Adla NK. Extraction Of Meaningful Rules In A Medical Database. Seventh International Conference 2008:450-6.
12. Rajeswari P, Reena G. Analysis of Liver Disorder Using Data mining Algorithm. *Global J Computer Science Technology* 2010;10:4.
13. D'Amico G, Garcia-Tsao G, Pagliaro L. Natural history and prognostic indicators of survival in cirrhosis: a systematic review of 118 studies. *J Hepatol* 2006;44:217-31.
14. Singh V, Nagpal S. A Guided clustering Technique for Knowledge Discovery—A Case Study of Liver Disorder Dataset. *International J of Computing Business Research* (1).
15. Asrani SK, Kim WR. Organ allocation for chronic liver disease: model for end-stage liver disease and beyond. *Curr Opin Gastroenterol* 2010;26:209-13.
16. Rousseeuw PJ. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *J computational applied Mathematics* 1987;20:53-65.
17. El-Serag HB. Epidemiology of viral hepatitis and hepatocellular carcinoma. *Gastroenterology* 2012;142:1264-73.
18. Pietrangelo A, Dohil R, Newbury R, Fox L, Bastian J, Aceves S. Reviews Basic Clin Gastroenterology.
19. Chiamonte M, Stroffolini T, Vian A, Stazi MA, Floreani A, Lorenzoni U, et al. Rate of incidence of hepatocellular carcinoma in patients with compensated viral cirrhosis. *Cancer* 1999;85:2132-7.