

سبحان



Data mining: Case Presentation

Dr O. Pournik MD, MPH, MSc, PhD

pournik@gmail.com

Case Presentation

Screening Models

استفاده از روش‌های داده‌کاوی
برای ساخت مدل غربال‌گری
رتینوپاتی نارسایی

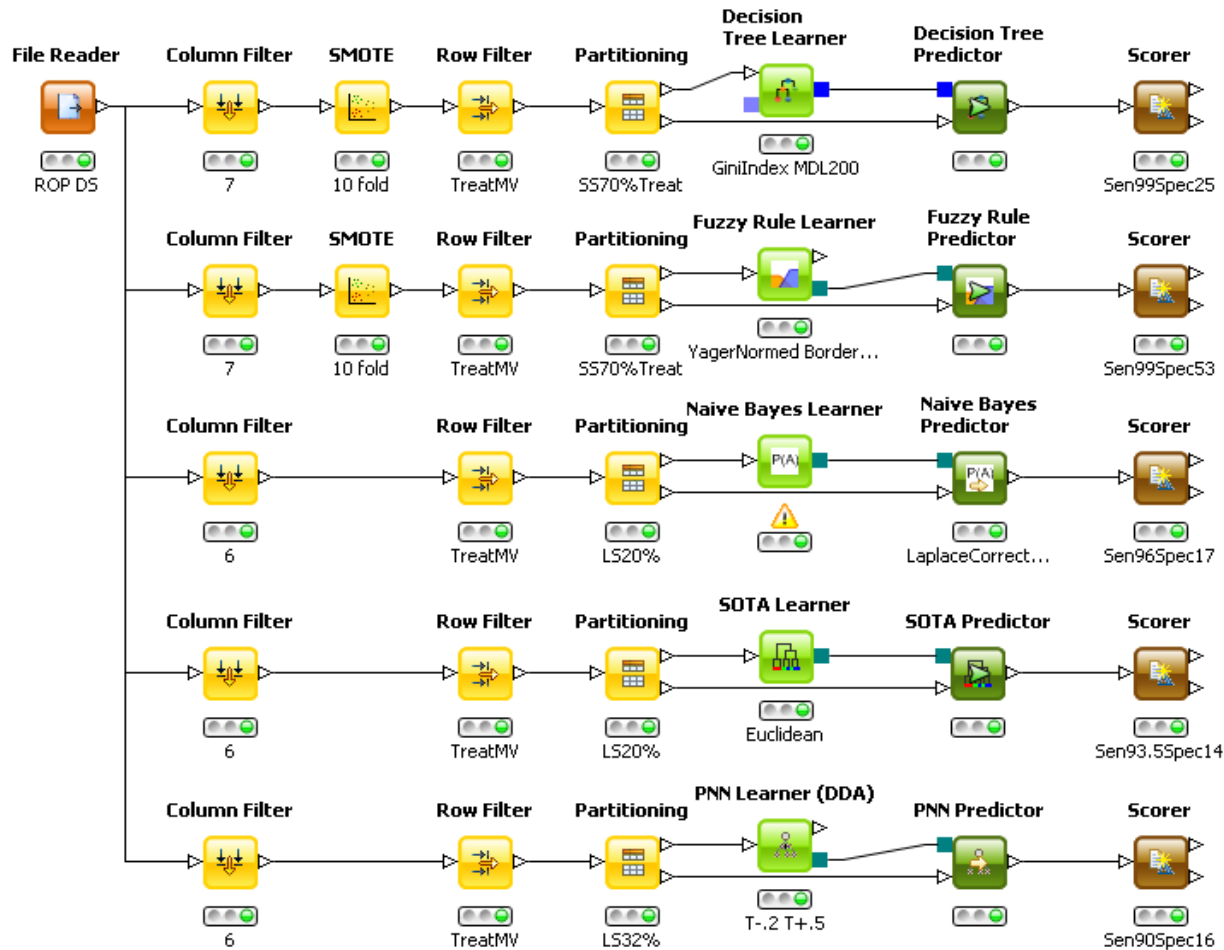


مدل سازی

- پس از تایید کیفیت داده‌ها، مدل سازی با توجه به ساختار categorical متغیر هدف انجام شد.
- مجموعه داده‌ها به صورت خوشه‌ای (stratified) به دو گروه داده های آموزشی (۷۰٪) و ارزیابی (۳۰٪) بر اساس متغیر هدف تقسیم شدند.
- برای مدل سازی از نرم افزار KNIME (نسخه ۲.۳.۱) و Clementine ۱۲ استفاده شد.
- مدل سازی با ۵ متد Decision tree، Fuzzy rule، Naïve Bayes، SOTA و PNN انجام شد.

مدل سازی

فرآیند و جریان ساخت مدل‌ها



مدل سازی

متغیرهای مورد استفاده در هر یک از مدل ها

متغیرهای استفاده شده در مدل	مدل
سن بارداری، وزن تولد	درخت تصمیم (Decision tree)
سن بارداری، سابقه سقط، آپگار بدو تولد، باروری با روشهای کمکی	قاعده فازی (Fuzzy Rule)
سن بارداری، وزن تولد، جنس، آپگار بدو تولد، آپگار دقیقه ۵	Naïve Bayes
سن بارداری، وزن تولد، جنس، آپگار بدو تولد، آپگار دقیقه ۵	الگوریتم درختی خود سازمان یافته (SOTA)
سن بارداری، وزن تولد، جنس، آپگار بدو تولد، آپگار دقیقه ۵	شبکه عصبی احتمالی (PNN)

ارزیابی

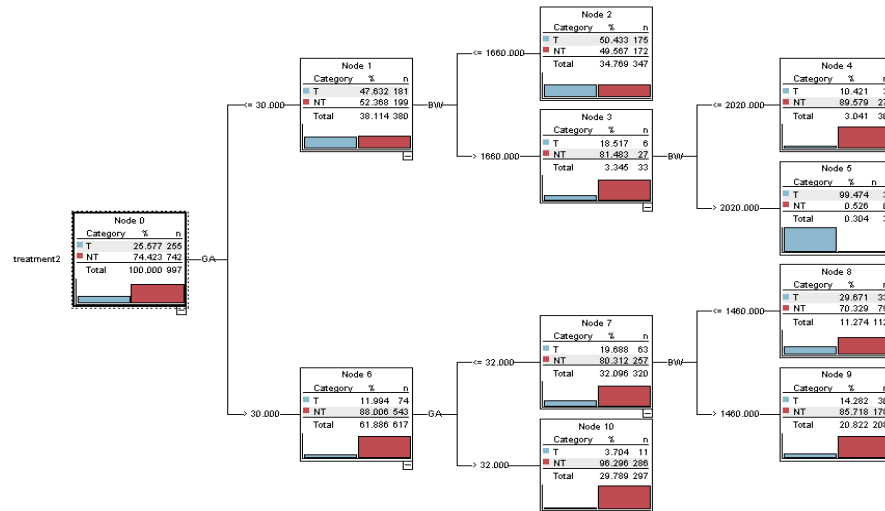
- مجموعه داده‌ها به صورت stratified به دو گروه داده‌های آموزشی (۷۰٪) و ارزیابی (۳۰٪) بر اساس متغیر هدف تقسیم شدند.

مدل	حساسیت	ویژگی	تعداد فیچرها در مدل	سطح زیر نمودار AUC
درخت تصمیم (Decision tree)	۰,۹۹۱	۰,۲۴۴	۲	۰,۸۷۵
قاعده فازی (Fuzzy Rule)	۰,۹۸۷	۰,۵۳۲	۴	۰,۹۰۶
Naïve Bayes	۰,۹۵۶	۰,۱۷۲	۵	۰,۸۱۵
الگوریتم درختی خود سازمان یافته (SOTA)	۰,۹۳۵	۰,۱۳۲	۵	۰,۸۱۴
شبکه عصبی احتمالی (PNN)	۰,۸۳۲	۰,۱۶۱	۵	۰,۷۶۵

مدل سازی-۲

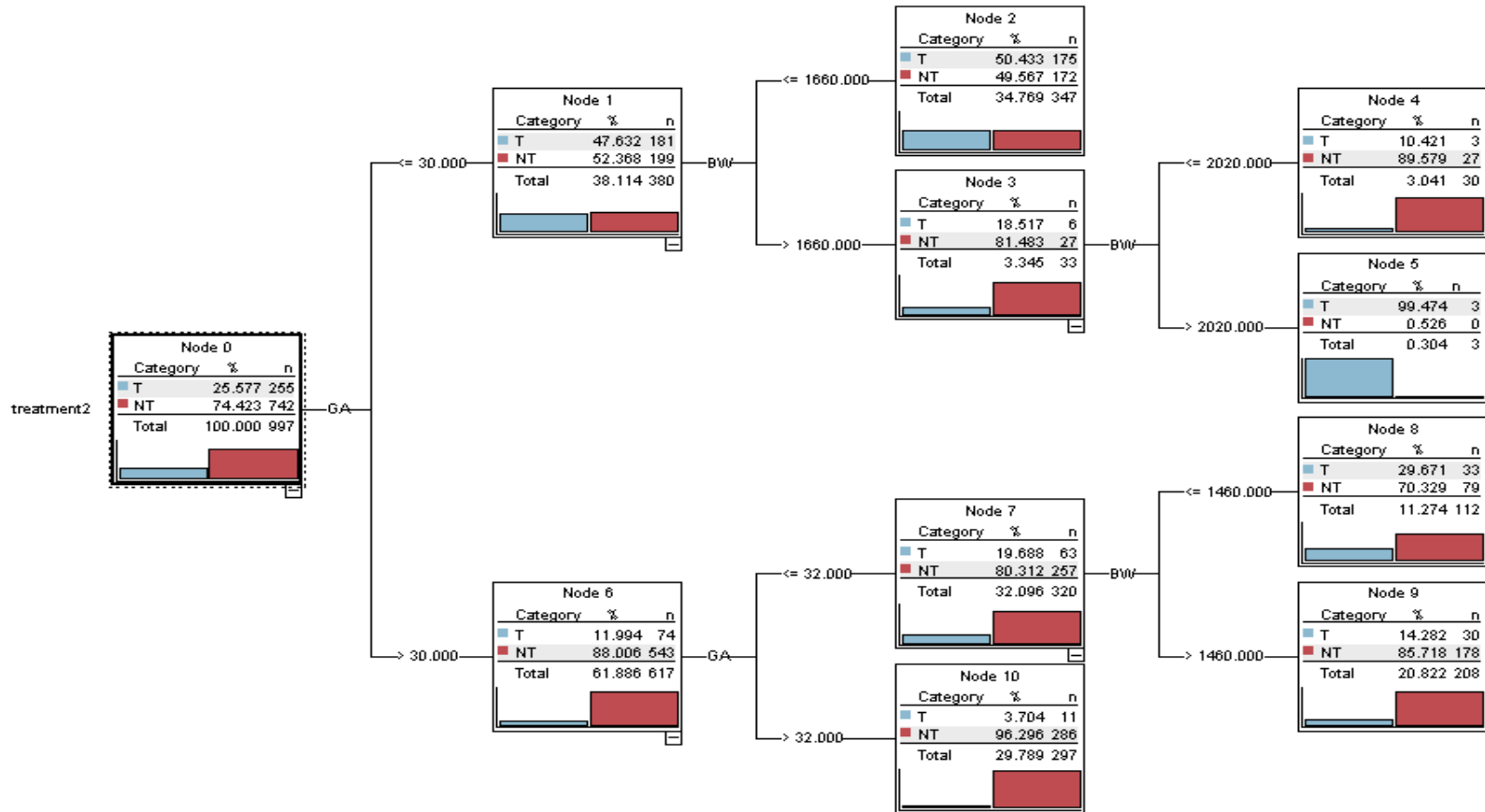
ارتقای مدل کاربردی

- با توجه به امکان شناخت ساده نحوه عمل کرد مدل درخت تصمیم و امکان به کارگیری آن به وسیله افراد بدون نیاز به رایانه، این مدل مورد بررسی بیشتری قرار گرفت.
- مدل اولیه درخت تصمیم گیری، مدل پیچیده‌ای بود که با بهره‌گیری از شاخص کیفیت GINI و بعد از حرس کردن (Pruning) با متد MDL (Minimum description Length) به شکل زیر تبدیل شد. در این مدل حساسیت به ۹۴٪ کاهش پیدا کرد و ویژگی به ۴۸.۷٪ افزایش یافت.



مدل سازی-۲

نمونه مدل کاربردی



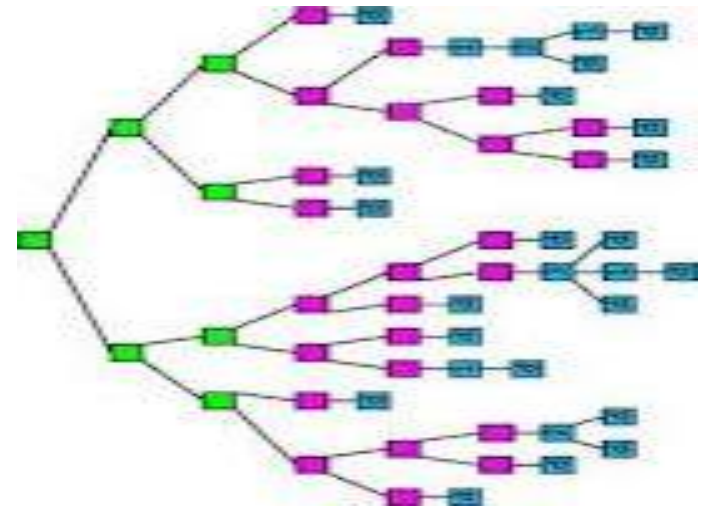
استقرار - نتیجه گیری

- تکنیک‌های داده کاوی می توانند در طراحی مدل‌های بهینه برای غربالگری نوزادان و کشف موارد ROP موثر باشد.
- در بین مدل‌های طراحی شده در این مطالعه، عملکرد مدل Fuzzy Rule بهتر از سایر مدل‌ها است و بهره‌گیری از آن می‌تواند علاوه بر کمک به متخصصین در کشف موارد نیازمند مداخله، به هزینه‌های اثربخش‌تر شدن غربالگری ROP منجر شود.
- هم‌چنین مدل درختی کاربردی شده نیز می‌تواند با توجه به حساسیت قابل قبول و عدم نیاز به رایانه در صورت نیاز به وسیله سیستم بهداشتی-درمانی مورد استفاده قرار گیرد.

Case Presentation

Diagnostic Models

طراحی مدل تشخیصی HTLV1
با استفاده از شمارش کامل سلول
های خونی



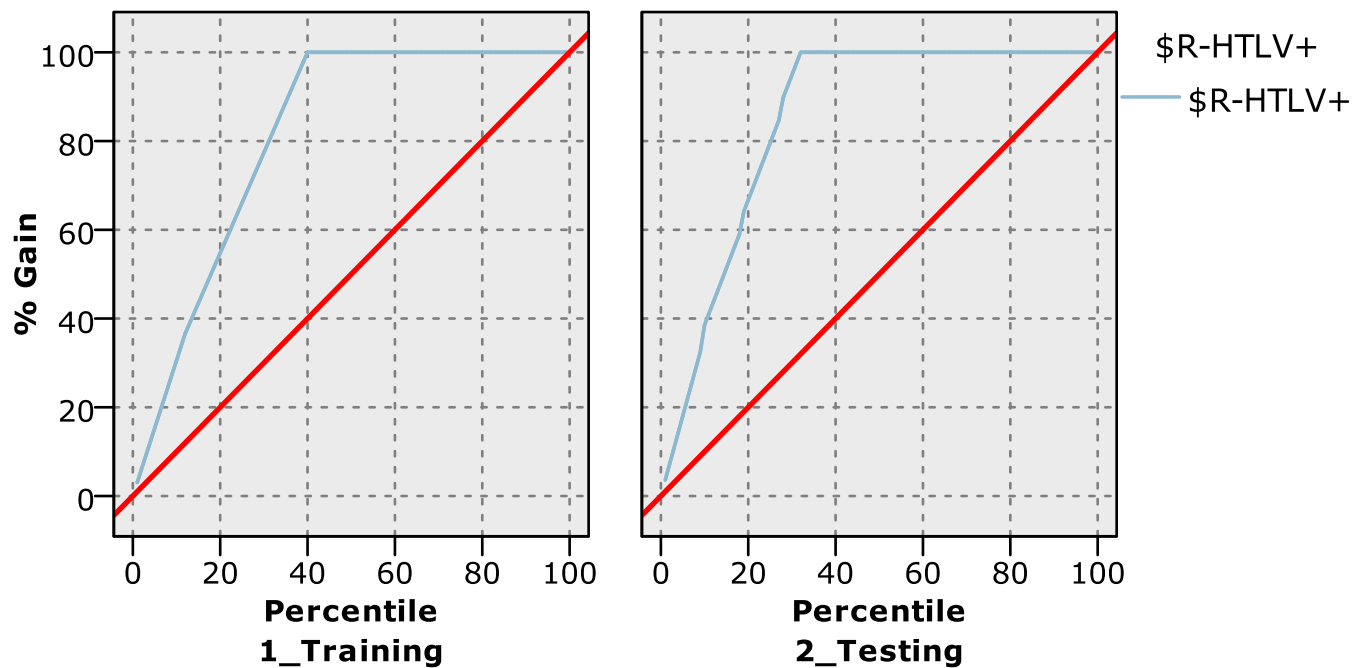
- برای انجام این مطالعه اطلاعات مربوط به نتایج ۵۹۹ مورد نمونه CBC شامل ۲۴۶ نمونه از بیماران مبتلا به HTLV1، ۱۴۲ نمونه از بیماران مبتلا به ALL، ۱۱۰ نمونه از بیماران مبتلا به ATL و ۱۰۱ نمونه از افراد سالم جمع آوری گردید.
- نتایج آزمایش‌های افراد از سیستم اطلاعات بیمارستانی بیمارستان‌های قائم و امام رضا (ع) مشهد استخراج شد.

- مدل سازی بر روی این داده‌ها با استفاده از متد CHAID و به کمک نرم‌افزار Clementine نسخه ۱۲ انجام شد. روش CHAID که مختصر شده عبارت Chi-squared Automatic Interaction Detector می‌باشد یکی از انواع درخت تصمیم است.

- بانک داده‌های مورد استفاده در این مطالعه محتوی نتایج آزمایش CBC ۳۰۲ نفر شامل ۹۴ بیمار مبتلا به لوکمی، ۱۰۱ فرد نرمال و ۱۰۷ بیمار مبتلا به HTLV1 بود. متغیرهای مورد بررسی شامل تعداد گلبول‌های سفید، تعداد گلبول‌های قرمز، هموگلوبین، هماتوکریت، MCH، MCV، MCHC، پلاکت، درصد پلی مورفونوکلئرها، درصد لنفوسیت‌ها، درصد مونوسیت‌ها و درصد ائوزینوفیل‌ها بود.

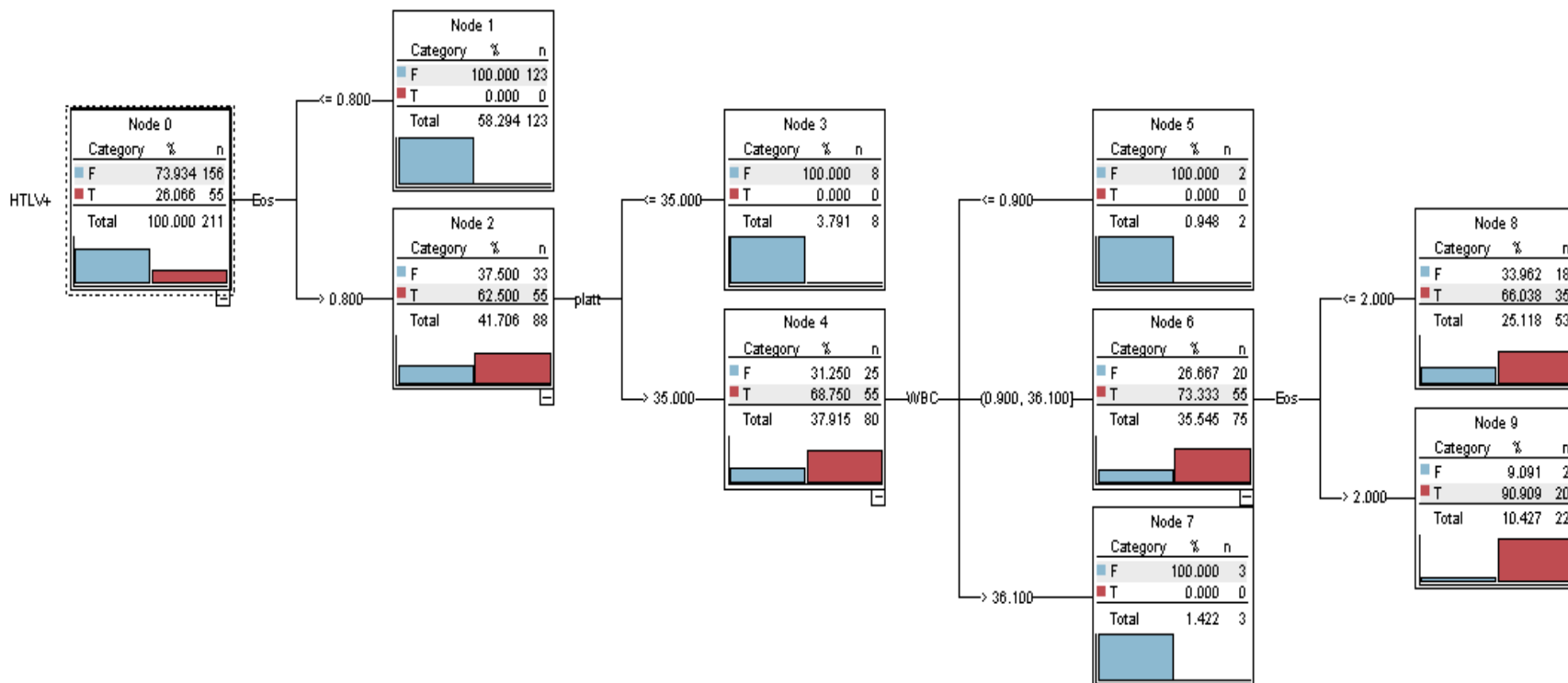
	HTLV1	Leukemia	Healthy	P value
WBC(10000)	13.7	66.2	8.7	0.698
RBC	4.0	3.6	4.6	<0.001
Hgb	11.6	9.8	13.2	<0.001
MCV	88.7	89.7	85.1	0.002
MCH	29.2	42.3	28.9	0.366
MCHC	34.0	32.9	33.9	0.797
Platelet	178	1009	292	0.286
Lymphocyte	29.3	27.8	32.4	0.142
Monocyte	4.7	5.4	6.2	0.504
Eosinophil	3.6	4.8	2.7	0.074

- میانگین متغیرها در گروه های مختلف و سطح معنی داری تفاوت بین گروه های مختلف بر اساس آنالیز واریانس



'HTLV+' = "T"

ارزیابی عملکرد مدل (Gain Chart)



• مدل نهایی به دست آمده در قالب درخت تصمیم

- این مدل بر اساس روش CHAID ساخته شده که یکی از انواع روش های مبتنی بر درخت تصمیم می باشد. مزیت کلی این روش در سادگی نمایش کلیه مراحل طبقه بندی است. همین امر بهره گیری از مدل به دست آمده را ممکن می سازد.

- با اجرای این مدل بر روی نتایج CBC افراد بدون علامت یا بیماران متفرقه در مناطق اندمیک کشور می توان افراد ناقل احتمالی را کشف و برای آزمایش‌های تکمیلی معرفی کرد. با توجه به احتمال بروز انواع بدخیمی‌ها و اختلالات سیستم اعصاب در افراد مبتلا، کشف این موارد می تواند به کاهش بار بیماری‌های مذکور و ارتقای سطح سلامت افراد به ویژه در مناطق اندمیک منجر شود. همچنین پیشنهاد می شود با استفاده از داده‌های بیشتر امکان توسعه روش و کیفیت مدل فراهم شود.

نتیجه گیری

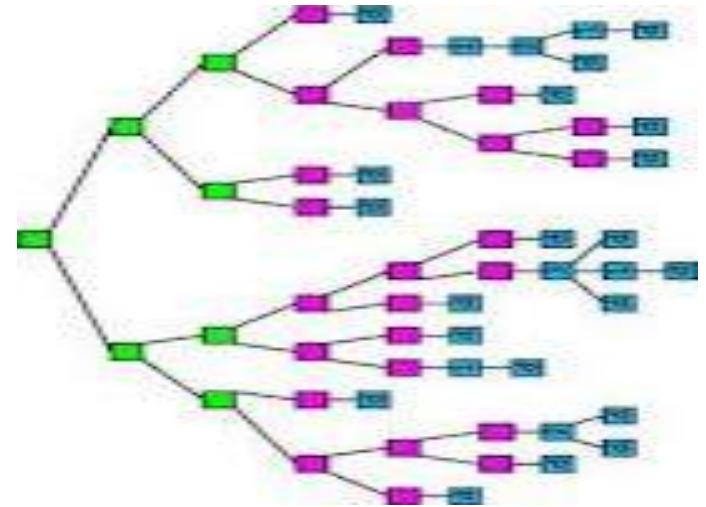
- مدل های تصمیم گیری مبتنی بر داده کاوی و به طور خاص تکنیک به کار گرفته شده در این مقاله و درخت تصمیم به دست آمده می توانند در شناسایی بیماران مبتلا به بیماری HTLV1 کمک کننده بوده و به پزشکان و متخصصین در شناسایی سریع تر و آسان تر این افراد کمک نماید.

Case Presentation

Diagnostic Models

کاربرد داده کاوی پیش بینی
کننده در پزشکی:

مقایسه الگوهای منتخب در شناسایی
زودهنگام بیماری منثريت باکتریایی
در ایران

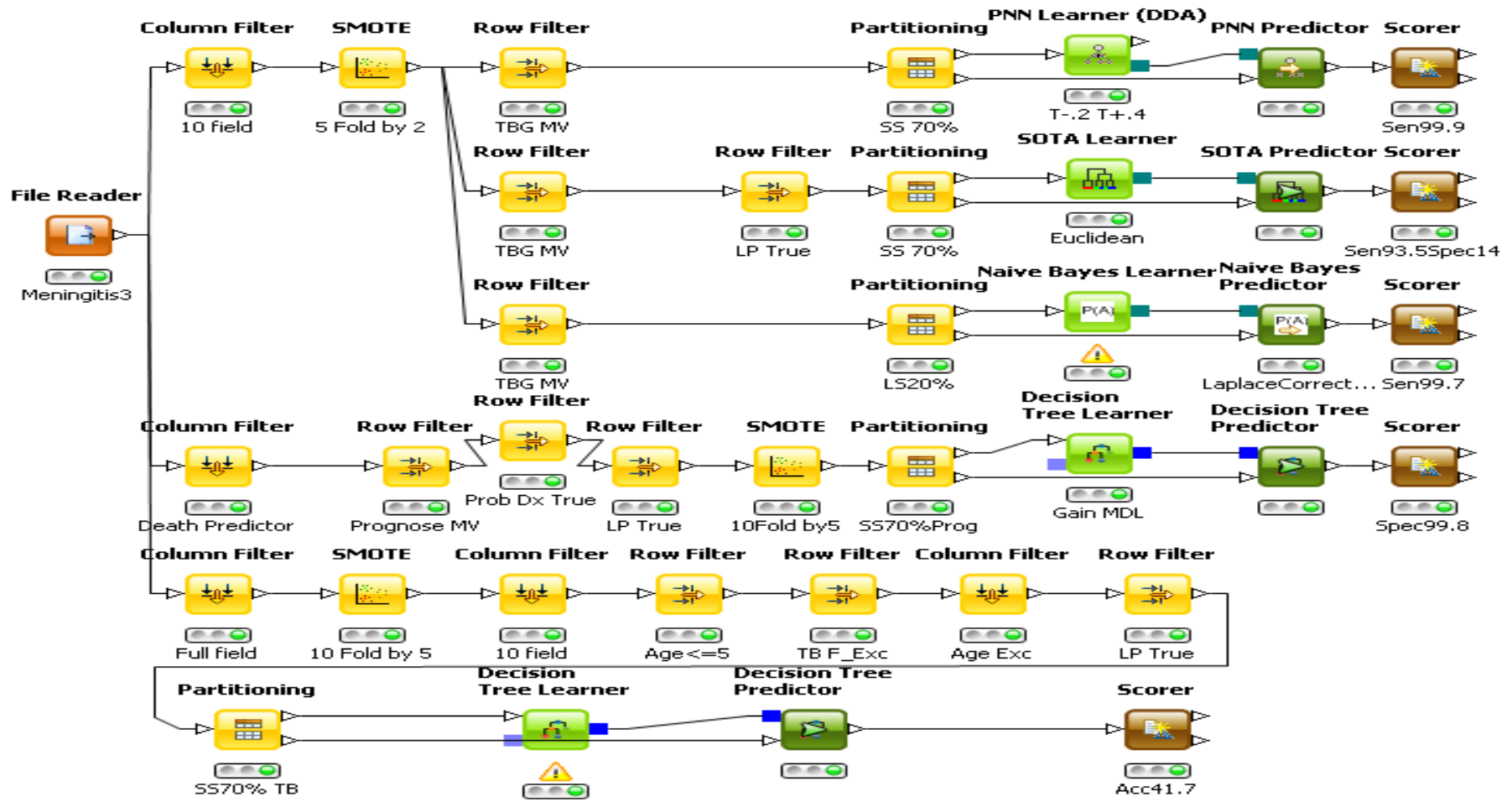


ادراک داده ها

نام متغیر	نام متغیر
مقدار قند	سن (سال)
کشت CSF	روز
کشت خون	جنس
لاتکس	شغل
رنگ آمیزی گرم	استان
تشخیص محتمل	دانشگاه
تشخیص قطعی	شهر
تشخیص نهایی	منطقه
آنتی بیوتیک قبل	انجام LP
آنتی بیوتیک بعد	ظاهر نمونه
نتیجه درمان	تعداد سلول
دریافت واکسن Hib	درصد PMN
دریافت واکسن مننگو کوکی	درصد لنفوسیت
	مقدار پروتئین

- بانک اطلاعاتی بیماران مننژیتی مرکز مدیریت بیماری های واگیر وزارت بهداشت در سال های ۱۳۸۸ و ۱۳۸۹ از کل کشور و مشتمل بر ۷۵۷۴ مورد

مدل سازی پیش گوئی کننده بیماری مننژیت باکتریایی



بحث و نتیجه گیری

Sensitivity	Accuracy	نام تکنیک
٪.۹۴/۷	٪.۸۸/۸	SOTA
٪.۹۹/۷	٪.۹۳/۱	Naïve Bayes

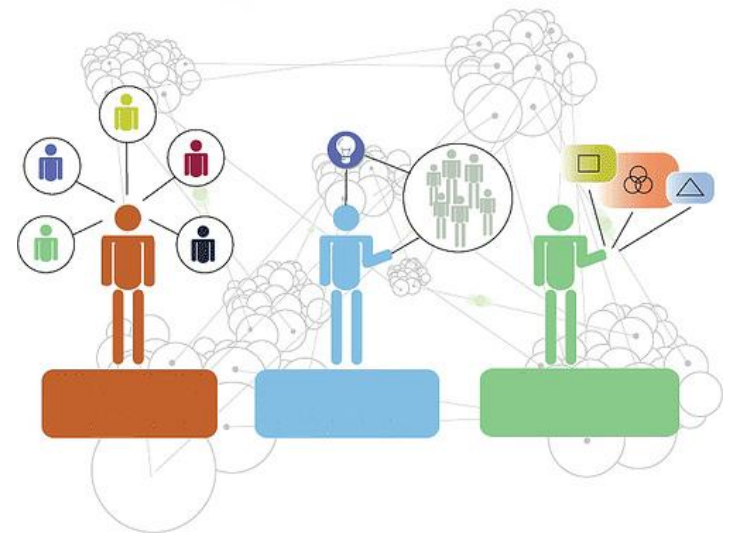
پیشنهادات در خصوص استفاده از نتایج این پژوهش

- ❖ نتایج حاصل از این پژوهش کاربرد موفق مدل‌های **Predictive** موثر را در تشخیص و کنترل بیماری‌ها نشان داد.
- ❖ این مدل‌ها با داده‌های بیشتر می‌توانند بهینه شده و در سیستم‌های بهداشت و درمان مانند سیستم‌های پایش و نظارتی و **HIS** و **DSS**‌ها مورد استفاده قرار گیرند.

Case Presentation

Fraud or Anomaly Detection

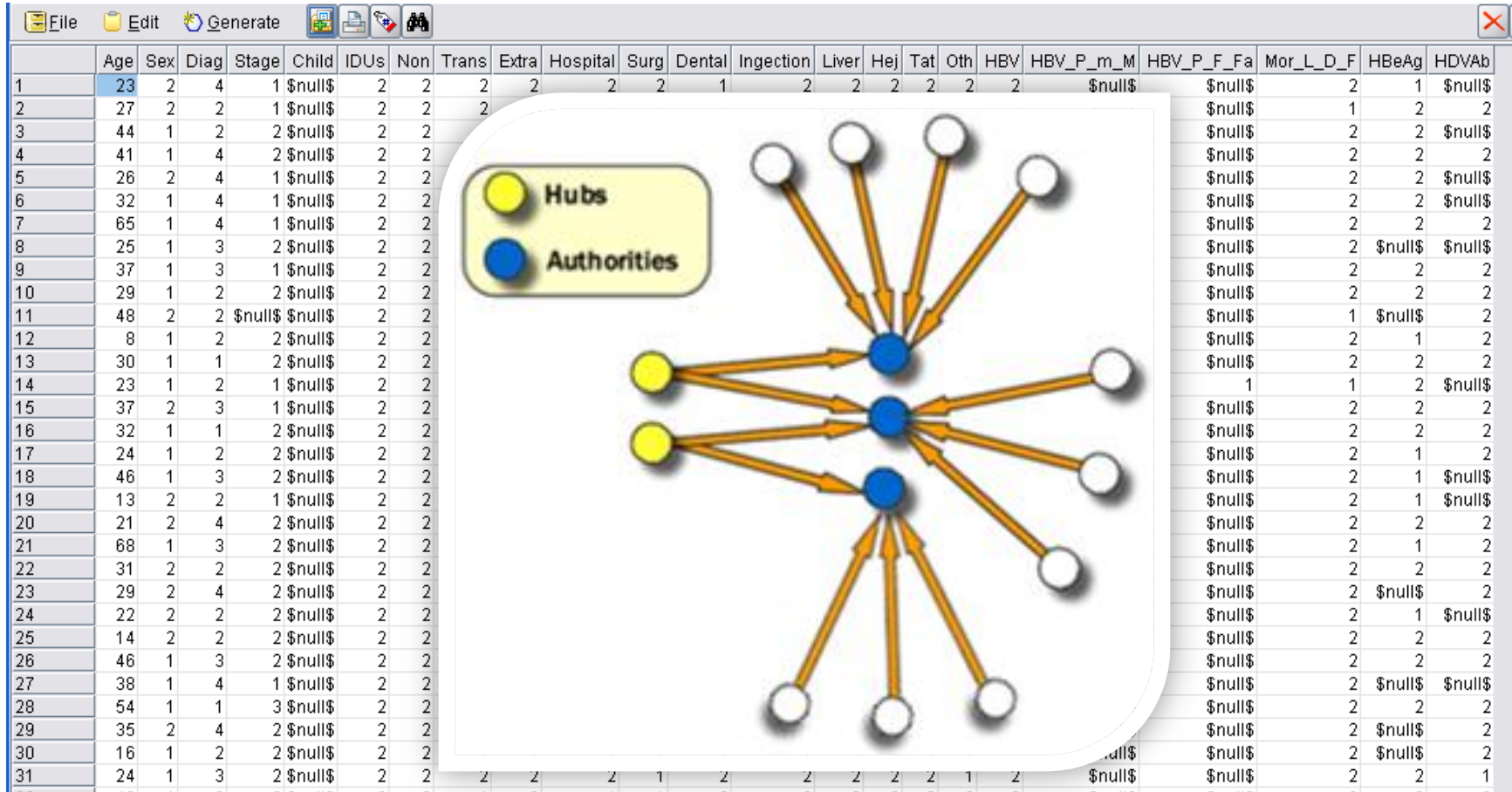
*Hepatitis B familial
transmission analysis
A Data mining approach*



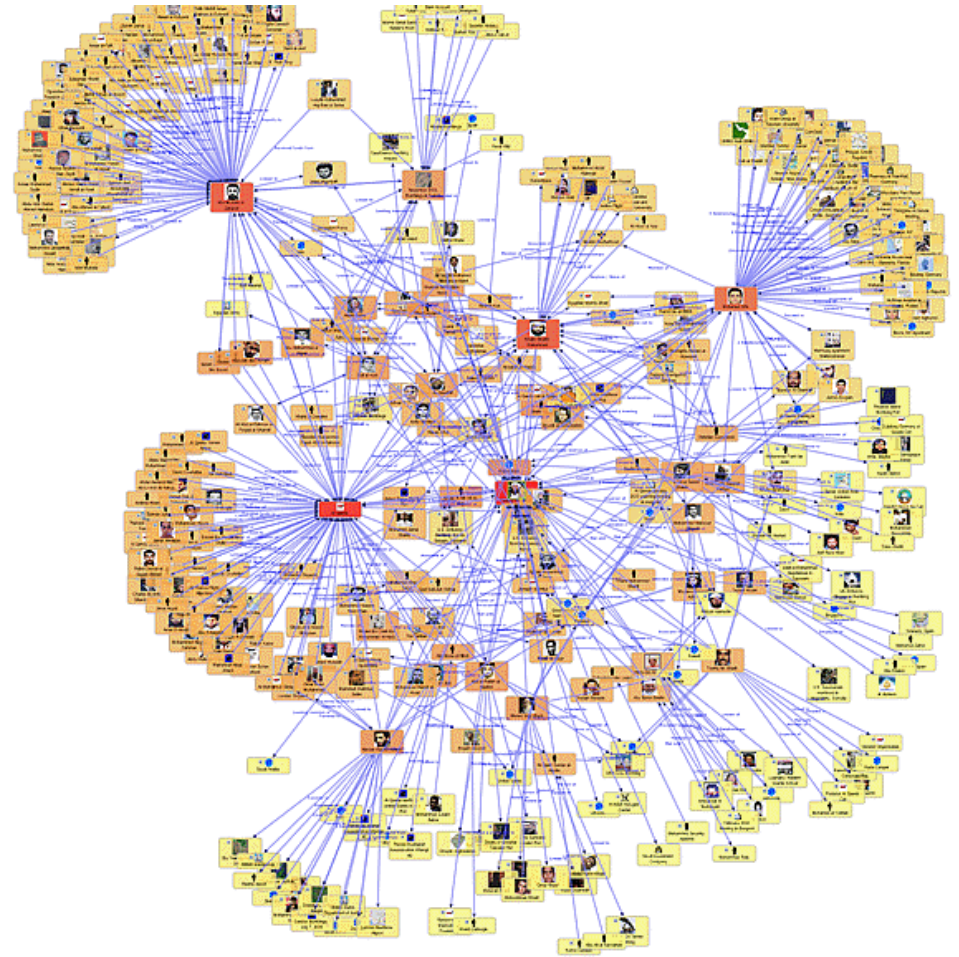
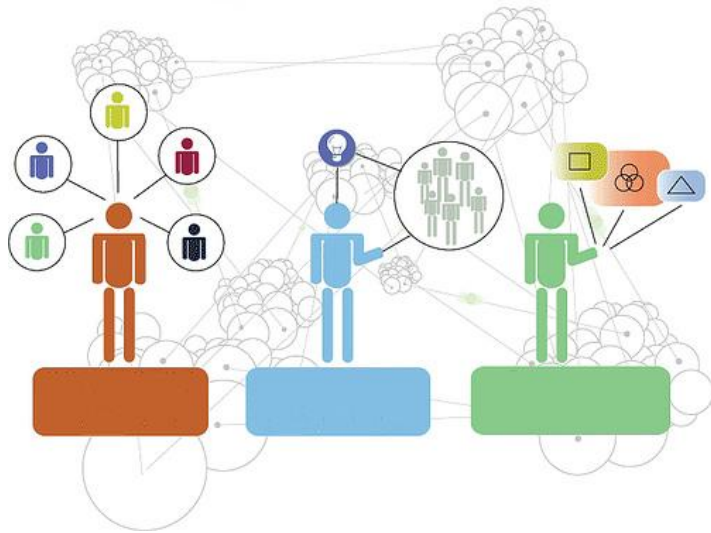
Data Set Specification

- **330 hepatitis B patients**
 - **Demographic data**
 - **Viral markers**
 - **Stage of disease**
 - **Family history of liver disease mortality**
 - **Morbidity**
- **Demographic data and viral markers of their offspring were also collected.**

Traditional Data Analysis View



Network Analysis



Feature Selection Analysis

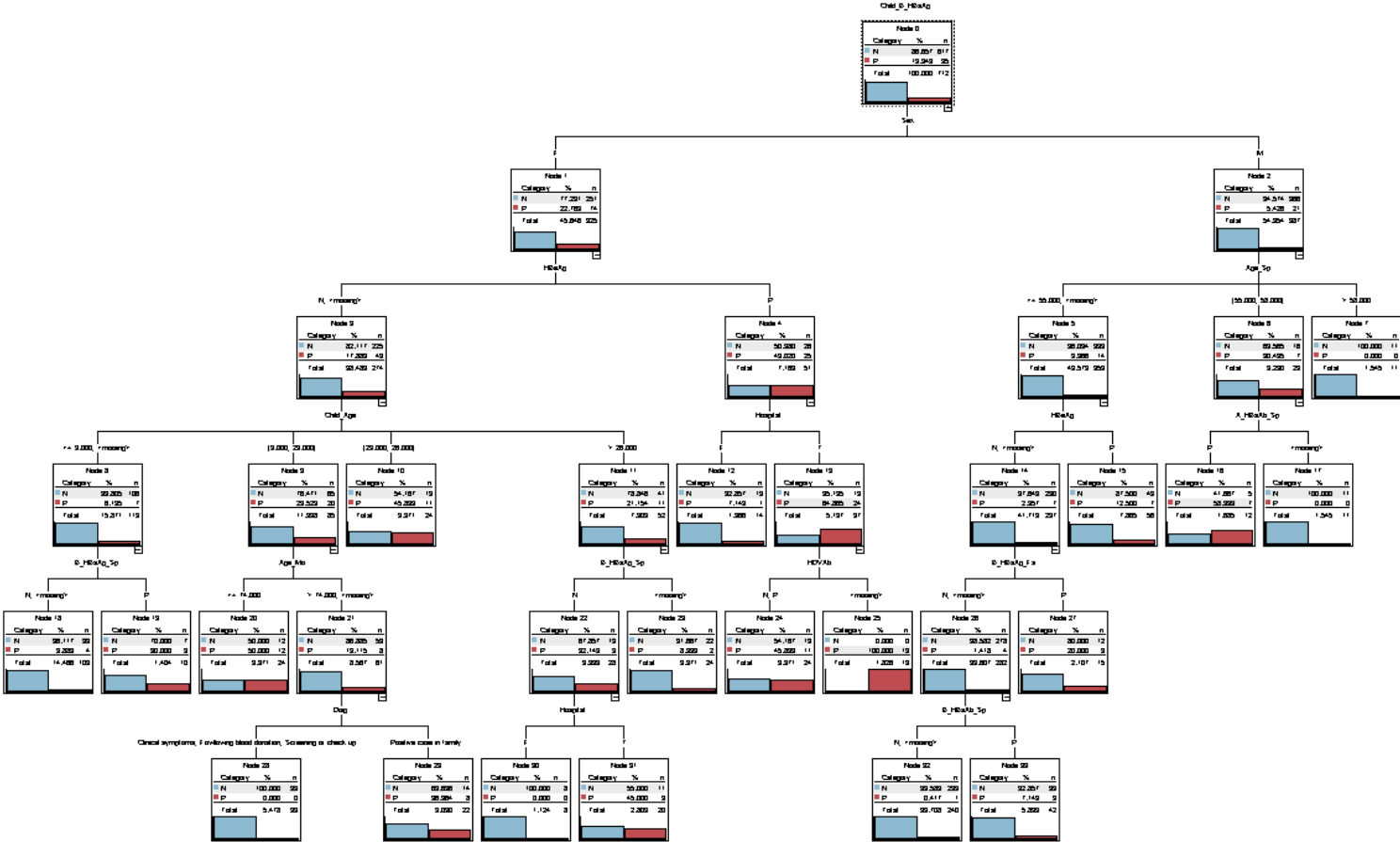
The screenshot shows a software interface for feature selection analysis. The main window displays a table of 20 features ranked by importance. A callout box highlights a subset of features, showing their relative ranking and importance within that subset.

Rank	Field	Type	Importance	Value
1	Sex	Flag	Important	1.0
2	HBeAg	Flag	Important	1.0
3	Diag	Set	Important	1.0
4	B_HBcA...	Flag	Important	1.0
5	Hospital	Flag	Important	0.998
6	Child_Age	Range	Important	0.99
7	Hej	Flag		
8	NonIDU	Flag		
9	Age_Sp	Range		
10	Dental	Flag		
11	Oth	Flag		
12	Liver	Flag		
13	Age	Range		
14	Du_Marr	Range		
15	Ict_His_...	Flag		
16	Re_Spo...	Flag	Unimport...	0.571
17	Trans	Flag	Unimport...	0.504
18	HDVAb	Flag	Unimport...	0.495
19	Stage	Set	Unimport...	0.319
20	Mor_L_D...	Flag	Unimport...	0.311

Rank	Field	Type	Importance	Value
1	Sex	Flag	Important	1.0
2	HBeAg	Flag	Important	1.0
3	Child_Age	Range	Important	0.99
4	Age	Range	Marginal	0.917
5	HDVAb	Flag	Unimport...	0.495

Selected fields:11 Total fields available:51

Findings (Tree View)

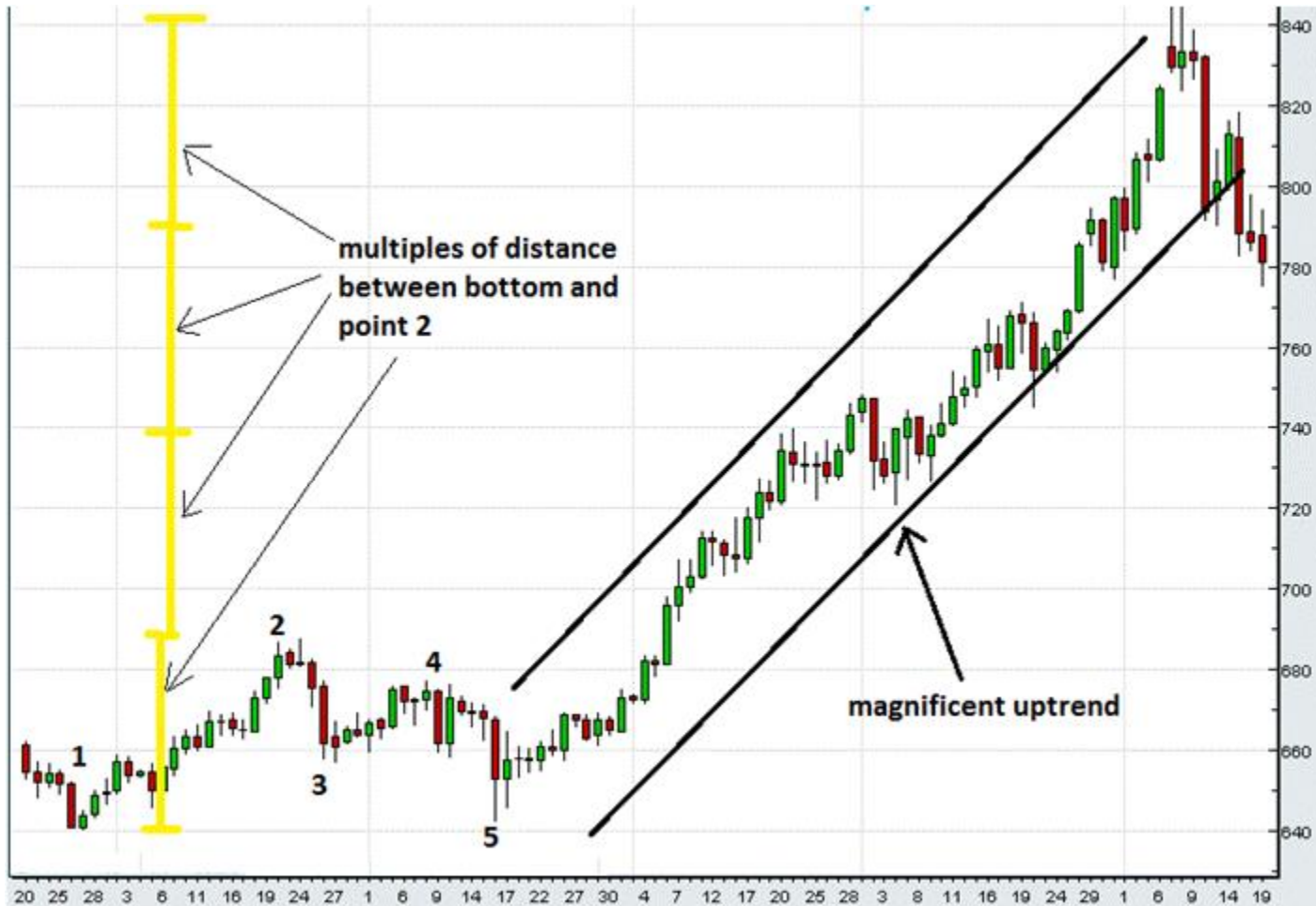


Findings

The most dominant finding in the extracted patterns was the difference in transmission pattern of **male and female index cases**.

- In male index cases,
 - **Family history of liver disease mortality and morbidity** was the main predictive factor of disease transmission to offspring.
 - In the next step, proxies of **exposure** duration such as offspring age were important.
- In female index cases,
 - Rate of transmission was significantly higher in female index cases and the main determining factor was **HBeAg**.
 - In HBeAg negative patients, the pattern of transmission was different in patients based on their **HDVAb**. Transmission rate was 25% in HDVAb negative and 5% in HDVAb positive cases.

Anomaly Evaluation Approach



Case Presentation

Health Policy Management

پیش بینی موفقیت درمان بیماری سل به کمک شبکه عصبی

اعلام اشرفی _ دکتر امید یوزنگ _ دکتر امیر زما زبون
دانشجوی دکتری انفورماتیک پزشکی، دکتری انفورماتیک پزشکی

چکیده:
سل، یکی از پرخطرترین عفونت‌های جهانی است. هدف از این مطالعه پیش‌بینی میزان موفقیت درمان سل در استان تهران است. شبکه عصبی عمیق (PNN) برای تشخیص موفقیت درمان سل با توجه به متغیرهای مختلف درمانی (158 آزمون سل) و وضعیت تشخیصی (100 آزمون سل) استفاده شد. در این مطالعه، میزان موفقیت درمان سل در استان تهران با استفاده از شبکه عصبی عمیق (PNN) به میزان 92 درصد برآورد شد. نتایج این مطالعه نشان می‌دهد که استفاده از شبکه عصبی عمیق (PNN) برای پیش‌بینی موفقیت درمان سل در استان تهران می‌تواند به پزشکان در تصمیم‌گیری‌ها کمک کند.

مقدمه:
در سده‌های اخیر، فرآیند کشف دانش (Knowledge discovery process) از اهمیت ویژه‌ای برخوردار شده است. داده‌های حجیم و متنوع در دسترس قرار گرفته‌اند. هدف از این مطالعه، پیش‌بینی موفقیت درمان سل در استان تهران است. در این مطالعه، از شبکه عصبی عمیق (PNN) برای تشخیص موفقیت درمان سل استفاده شد. نتایج این مطالعه نشان می‌دهد که استفاده از شبکه عصبی عمیق (PNN) برای پیش‌بینی موفقیت درمان سل در استان تهران می‌تواند به پزشکان در تصمیم‌گیری‌ها کمک کند.

مکان قرارگیری استان تهران در کشور ایران

روش و نتایج:
در این مطالعه، از شبکه عصبی عمیق (PNN) برای تشخیص موفقیت درمان سل استفاده شد. نتایج این مطالعه نشان می‌دهد که استفاده از شبکه عصبی عمیق (PNN) برای پیش‌بینی موفقیت درمان سل در استان تهران می‌تواند به پزشکان در تصمیم‌گیری‌ها کمک کند.

نتیجه گیری:
نتایج این مطالعه نشان می‌دهد که استفاده از شبکه عصبی عمیق (PNN) برای پیش‌بینی موفقیت درمان سل در استان تهران می‌تواند به پزشکان در تصمیم‌گیری‌ها کمک کند.

CRISP-DM (Cross Industry Standard Process for Data Mining)

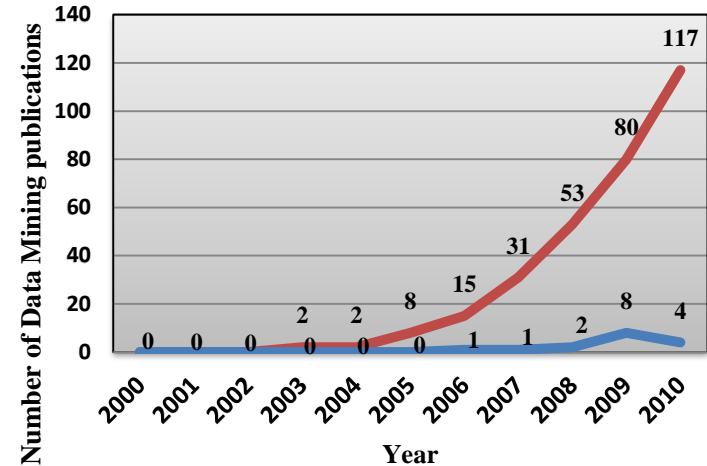
منابع:
1. ...
2. ...
3. ...

تشکر و قدردانی:
از همه کسانی که در فرآیند تهیه این مقاله کمک کردند تشکر می‌کنیم.

Case Presentation

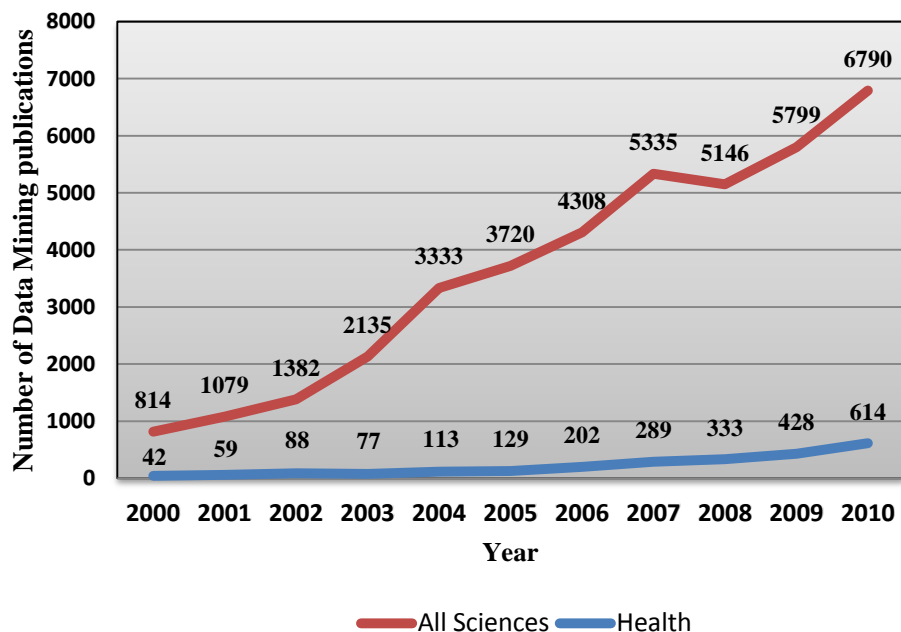
Scientometric Approach

Data Mining in Medicine: A Scientometrics Perspective

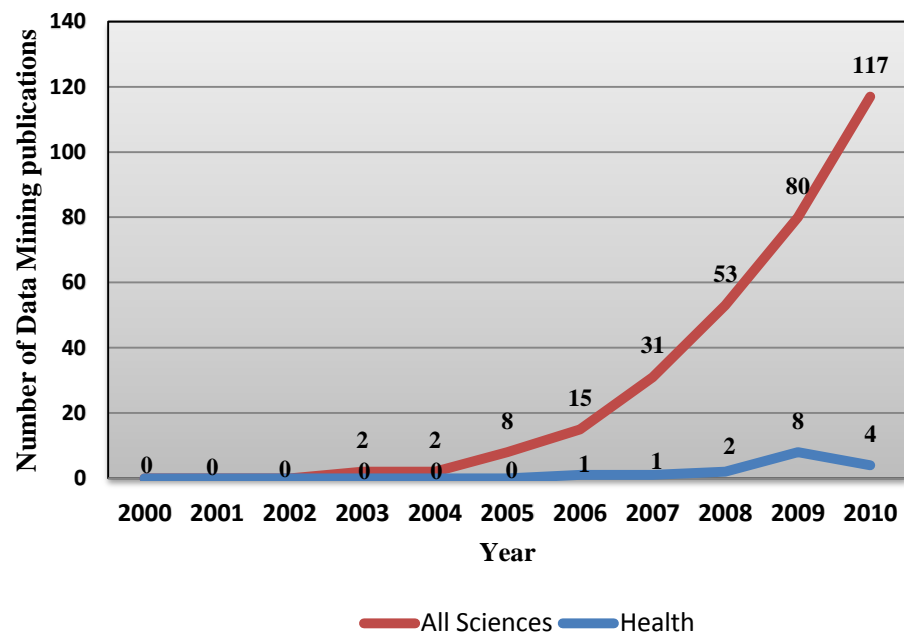


Comparing growth rate of Data Mining publications in health and all sciences in Iran

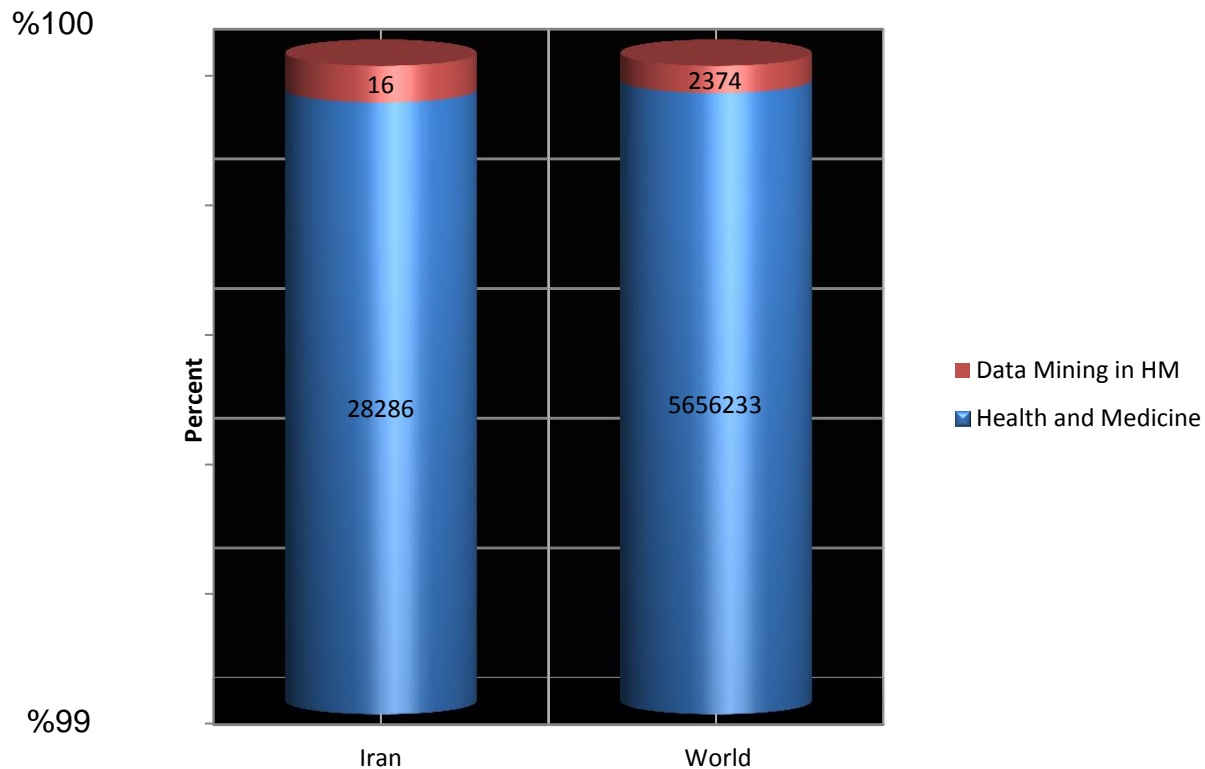
World Trend



Iran Trend



Share of Data Mining publications in total number of articles in Medicine



Findings and Suggestions of Scientometric Study of DataMining in Medicine

- Publication Trends in Iran is somehow in accordance with world growth trend
- But the volume of publication can be much higher Specially in health and medicine domain.



**Any
Questions?**