

سید ابن الحنفیہ



دانشگاه قم

دانشکده فنی و مهندسی

پایان نامه دوره کارشناسی ارشد مهندسی فناوری اطلاعات

عنوان:

**استفاده از داده‌کاوی برای کشف الگوهای پنهان در
داده‌های سرطان سینه
(بیمارستان شهید رجایی آران و بیدگل)**

استاد راهنما:

دکتر بهروز مینایی بیدگلی

استاد مشاور:

دکتر فریبا مینایی بیدگلی

نگارنده:

فاطمه السادات طباطبائی عینکی

تابستان / ۱۳۹۰



«صورت جلسه دفاع از پایان نامه کارشناسی ارشد»

با تأییدات خداوند متعال و با استعانت از حضرت ولی عصر(عجل الله تعالی فرجه الشریف)

جلسه دفاعیه پایان نامه کارشناسی ارشد آقای / خانم: رشته:

تحت عنوان:

با حضور هیأت داوران در محل دانشگاه قم در تاریخ: / / تشکیل گردید.

در این جلسه، پایان نامه با نمره (به عدد، به حروف)

بدرجه: عالی بسیار خوب خوب قابل قبول مورد دفاع قرار گرفت.

نام و نام خانوادگی	سمت	مرتبه علمی	امضاء
	استاد راهنما		
	استاد مشاور		
	استاد ناظر		
	استاد ناظر		
	نماینده کمیته تحصیلات تکمیلی		

مدیر امور آموزش و تحصیلات تکمیلی
نام و امضاء

معاون آموزشی و پژوهشی دانشکده
نام و امضاء

تقدیم به:

ساحت مقدس امام زمان روحی فداه

و همچنین

تقدیم به روح پدر بزرگوارم

که الطافش همیشه شامل حالم بوده و هست،

مادر مهربانم

که همواره بر مشکلاتم صبوری کرد،

و همسر فداکارم

که با گذشت و ایثارش مرا یاری نمود

تشکر و قدردانی:

خداوند منان را شاکرم که این عبد ناچیز و حقیرش را توفیق تحصیل علم داد، و تا این مراحل و مدارج علمی رساند، و نیز از اساتید محترم دکتر بهروز مینایی استاد راهنما و دکتر فریبا مینایی استاد مشاور که بنده را در خلق این اثر کمک و راهنمایی کردند، و دیگر اساتید دوران کارشناسی ارشد که با مشاوره‌های خویش بنده را مساعدت نمودند، کمال تشکر و سپاس را دارم و برای این عزیزان توفیق روزافزون و عمر با برکت از درگاه احدیت مسئلت دارم.

چکیده:

سرطان، یکی از دلایل اصلی مرگ در سراسر جهان است. در این بیماری، سلول‌های بدن به صورت غیر قابل کنترل رشد می‌کنند. گروهی از سلول‌هایی که به سرعت تکثیر شده‌اند، ممکن است یک توده، جرم یا بافت اضافی ایجاد کنند. این توده‌ها، تومور نامیده می‌شوند. توده‌ها می‌توانند سرطانی یا غیرسرطانی باشند. سرطان سینه، دومین عامل مرگ زنان است. این بیماری در میان مردان و زنان دیده شده، گرچه موارد ابتلا در میان مردان نادر است. در ایران نیز سرطان سینه، شایع‌ترین سرطان در میان زنان است. البته میزان ابتلا به این بیماری در ایران، بسیار کمتر از آمریکا و سایر کشورهای اروپایی است. همچنین سن ابتلا به این بیماری در ایران، نسبت به کشورهای توسعه‌یافته، یک دهه پائین‌تر است. فاکتور اصلی در کاهش مرگ و میر ناشی از این بیماری، تشخیص سریع و صحیح است. در این تحقیق، یک سیستم رده‌بندی ارائه شده که می‌تواند در تشخیص تومور به پزشکان کمک کند. به علت چند رده‌ای و نامتوازن بودن داده‌های مورد استفاده در این تحقیق، رهیافت جدیدی برای یادگیری در داده‌های چند رده‌ای و نامتوازن ارائه شده است. حساسیت سیستم رده‌بندی ارائه شده، ۸۶.۵۳ درصد و دقت آن، ۸۹.۲۵ درصد می‌باشد. همچنین در این تحقیق با استفاده از الگوریتم *Apriori*، به شناسایی عوامل موثر در بروز توده‌های خوش‌خیم و بدخیم پرداخته شده و عوامل افزایش سن، منارک زودرس، یائسگی دیررس، داشتن سن بالا هنگام اولین بارداری، سابقه‌ی طولانی مدت مصرف قرص‌های ضدبارداری، چاقی، نداشتن سابقه‌ی بارداری، زایمان و شیردهی و تعداد دفعات شیردهی بالا، به عنوان عوامل موثر در بروز توده، شناخته شده‌اند. اگر چه بیشتر این عوامل قبلاً به عنوان فاکتورهای خطرزا در بروز سرطان سینه معرفی شده بودند، فاکتور تعداد دفعات شیردهی بالا، در این مطالعه به عنوان یک فاکتور خطرزا تعیین گردیده که با یافته‌های محققان قبلی در تعارض می‌باشد. این مسئله می‌تواند به دلیل تعداد بسیار اندک نمونه‌های مبتلا به تومور بدخیم باشد.

کلمات کلیدی: داده‌کاوی، رده‌بندی، قوانین انجمنی، تشخیص سرطان سینه، داده‌های

نامتوازن و چند رده‌ای

فهرست مطالب

صفحه	عنوان
۱	فصل ۱: مقدمه.....
۲	۱-۱. موضوع تحقیق.....
۳	۲-۱. اهمیت و ضرورت تحقیق.....
۵	۳-۱. قلمرو تحقیق.....
۵	۴-۱. فرضیات تحقیق.....
۶	۵-۱. سوالات تحقیق.....
۷	۶-۱. اهداف و کاربردهای تحقیق.....
۸	۷-۱. نوآوری در تحقیق.....
۹	۸-۱. چالش‌ها و مشکلات انجام تحقیق.....
۱۱	۹-۱. محدودیت‌های تحقیق.....
۱۲	۱۰-۱. ساختار پایان‌نامه.....
۱۴	فصل ۲: ادبیات تحقیق.....
۱۵	۱-۲. مقدمه.....
۱۵	۲-۲. داده‌کاوی.....
۱۷	۱-۲-۲. مفهوم داده‌کاوی.....
۱۸	۲-۲-۲. مراحل داده‌کاوی.....
۱۸	۳-۲-۲. پیش‌پردازش.....
۲۱	۴-۲-۲. داده‌کاوی.....
۲۲	۵-۲-۲. پس‌پردازش.....
۲۲	۶-۲-۲. کاربردهای داده‌کاوی.....
۲۳	۳-۲. داده‌کاوی در پزشکی.....
۲۵	۴-۲. سرطان.....

۲۶	سرطان سینه	۱-۴-۲
۲۷	عوامل بروز سرطان سینه	۲-۴-۲
۲۸	آمار ابتلا و مرگ و میر ناشی از سرطان سینه در ایران و جهان	۳-۴-۲
۳۴	الگوریتم‌های رده‌بندی	۵-۲
۳۵	درخت تصمیم	۱-۵-۲
۳۸	PSO	۲-۵-۲
۴۰	DPSO	۳-۵-۲
۴۱	ماشین بردار پشتیبان	۴-۵-۲
۴۴	شبکه‌ی عصبی	۵-۵-۲
۴۵	BAGGING	۶-۵-۲
۴۷	ADABOOST	۷-۵-۲
۵۱	قوانین انجمنی	۶-۲
۵۲	APRIORI	۱-۶-۲
۵۷	فصل ۳: مرور کارهای انجام شده	
۵۸	۱-۳ مقدمه	
۵۸	۲-۳ کارهای مرتبط	
۶۵	فصل ۴: داده‌های نامتوازن	
۶۶	۱-۴ مقدمه	
۶۶	۲-۴ مفهوم داده‌ی نامتوازن	
۶۸	۳-۴ روش‌های یادگیری در داده‌های نامتوازن	
۶۹	۱-۳-۴ نمونه‌برداری	
۷۳	۲-۳-۴ روش‌های حساس به هزینه	
۷۵	۴-۴ معیارهای ارزیابی رده‌بند در داده‌های نامتوازن	
۷۹	۵-۴ معیارهای ارزیابی رده‌بند در داده‌های نامتوازن و چند رده‌ای	

۸۱ میانگین گیری میکرو	۴-۵-۱
۸۲ میانگین گیری ماکرو	۴-۵-۲
۸۳ فصل ۵: پیش پردازش داده ها	
۸۴ ۱-۵ مقدمه	
۸۴ ۲-۵ روش جمع آوری داده	
۸۵ ۳-۵ ویژگی ها	
۹۱ ۴-۵ نحوه توزیع داده ها بر اساس ویژگی ها	
۹۱ ۱-۴-۵ نوع تومور	
۹۲ ۲-۴-۵ سن	
۹۳ ۵-۵ پیش پردازش های انجام شده	
۹۳ ۱-۵-۵ اصلاح مقادیر گم شده	
۹۳ ۲-۵-۵ اصلاح داده های پرت	
۹۳ ۳-۵-۵ ایجاد ویژگی	
۹۴ ۴-۵-۵ نرمال سازی	
۹۴ ۵-۵-۵ حذف ویژگی های اضافی	
۹۴ ۶-۵ SAMPLING	
۹۶ فصل ۶: الگوریتم پیشنهادی	
۹۷ ۱-۶ مقدمه	
۹۸ ۲-۶ EASYENSEMBLE	
۱۰۰ ۳-۶ MODIFIEDBAGGING	
۱۰۱ ۴-۶ الگوریتم پیشنهادی بر مبنای EASYENSEMBLE	
۱۰۶ فصل ۷: نتایج و یافته های تحقیق	
۱۰۷ ۱-۷ مقدمه	

۱۰۷ رده‌بندی	۲-۷
۱۰۸ نتیجه‌ی حاصل از اجرای الگوریتم پیشنهادی	۱-۲-۷
۱۱۶ مقایسه با الگوریتم‌های پایه	۲-۲-۷
۱۱۹ مقایسه با سایر روش‌های یادگیری در داده‌های نامتوازن	۳-۲-۷
۱۲۵ مقایسه کارایی الگوریتم پیشنهادی با الگوریتم DPSO	۴-۲-۷
۱۲۸ قوانین انجمنی	۳-۷
۱۲۸ گسسته‌سازی ویژگی‌های پیوسته	۱-۳-۷
۱۲۹ حذف ویژگی‌های غیرمفید	۲-۳-۷
۱۲۹ تجمیع مقادیر ویژگی‌ها	۳-۳-۷
۱۳۰ تبدیل ویژگی‌ها	۴-۳-۷
۱۳۰ معرفی قوانین	۵-۳-۷
۱۳۵ فصل ۸: نتیجه‌گیری و پیشنهادها	
۱۳۶ مقدمه	۱-۸
۱۳۶ نتیجه‌گیری	۲-۸
۱۳۷ رده‌بندی	۱-۲-۸
۱۳۹ قوانین انجمنی	۲-۲-۸
۱۴۰ پیشنهادها	۳-۸
۱۴۰ مجموعه‌ی داده	۱-۳-۸
۱۴۰ داده‌کاوی	۲-۳-۸
۱۴۲ مراجع	
۱۵۵ اصطلاح‌نامه	
۱۶۲ پیوست الف: نتایج تفصیلی	
۱۶۵ پیوست ب: قوانین انجمنی	

فهرست جداول

عنوان	صفحه
جدول ۱-۲: تعداد زنان ایرانی مبتلا به ۱۰ نوع از رایج‌ترین سرطان‌ها در سال‌های ۲۰۰۳ تا ۲۰۰۶	۳۳
جدول ۱-۴: ماتریس هزینه‌ی چند رده‌ای.....	۷۴
جدول ۲-۴: ماتریس اغتشاش برای مسائل دو رده‌ای.....	۷۶
جدول ۳-۴: ماتریس اغتشاش برای مسائل چند رده‌ای.....	۸۰
جدول ۱-۵: ویژگی‌های موجود در مجموعه‌ی داده.....	۸۵
جدول ۱-۷: قوانین به دست آمده برای تومور بدخیم با حداقل پشتیبان ۰.۰۶ و حداقل اطمینان ۷۰	۱۳۱
جدول ۲-۷: قوانین به دست آمده برای تومور بدخیم با حداقل پشتیبان ۰.۰۳ و حداقل اطمینان ۷۰ (بدون مشخصات توده)	۱۳۲
جدول ۳-۷: قوانین به دست آمده برای تومور خوش‌خیم با حداقل پشتیبان ۰.۰۶ و حداقل اطمینان ۸۰ (بدون مشخصات توده)	۱۳۲

فهرست تصاویر

عنوان	صفحه
تصویر ۱-۲: ارتباط داده‌کاوی با سایر علوم.....	۱۸
تصویر ۲-۲: فرآیند کشف دانش از پایگاه داده.....	۱۸
تصویر ۳-۲: آغاز سرطان.....	۲۵
تصویر ۴-۲: مقایسه‌ی میزان ابتلا به سرطان سینه در سراسر جهان در سال ۲۰۰۸.....	۳۲
تصویر ۵-۲: نمونه‌ای از یک درخت تصمیم.....	۳۵
تصویر ۶-۲: نحوه‌ی کدگذاری ذرات در الگوریتم DPSO.....	۴۱
تصویر ۷-۲: یک مجموعه‌ی آموزش دوبعدی که داده‌های آن به صورت خطی قابل جداسازی است.....	۴۲
تصویر ۸-۲: دو خط جداساز با حاشیه‌های مختلف.....	۴۳
تصویر ۹-۲: افزایش صحت مدل با استفاده از BAGGING.....	۴۵
تصویر ۱۰-۲: شبه‌کد الگوریتم BAGGING.....	۴۶
تصویر ۱۱-۲: شبه‌کد الگوریتم ADABOOST.....	۵۰
تصویر ۱۲-۲: شبه‌کد الگوریتم APRIORI.....	۵۳
تصویر ۱۳-۲: شبه‌کد الگوریتم تولید کاندید.....	۵۴
تصویر ۱-۴: (a) یک مجموعه‌ی داده با عدم توازن بین رده‌ای (b) یک مجموعه داده‌ی پیچیده دارای عدم توازن بین رده‌ای و داخل رده‌ای.....	۶۸
تصویر ۲-۴: (a) k تا از نزدیک‌ترین همسایه‌های x_i با فرض $k = 6$ (b) تولید داده بر اساس فاصله‌ی اقلیدسی.....	۷۲
تصویر ۱-۵: نمونه‌برداری طبقه‌بندی شده.....	۹۵
تصویر ۱-۶: شبه‌کد الگوریتم EASYENSEMBLE.....	۹۹
تصویر ۲-۶: شبه‌کد الگوریتم MODIFIEDBAGGING.....	۱۰۰
تصویر ۳-۶: فلوچارت الگوریتم پیشنهادی.....	۱۰۳

تصویر ۴-۶: شبه‌کد الگوریتم پیشنهادی..... ۱۰۴

فهرست نمودارها

صفحه	عنوان
۱۶	نمودار ۱-۲: تعداد صفحات وب فهرست شده توسط موتور جستجوی گوگل از سال ۱۹۹۸ تا ۲۰۰۳.....
۱۶	نمودار ۲-۲: تعداد صفحات وب فهرست شده توسط موتور جستجوی گوگل در طول سه سال اخیر.....
۲۹	نمودار ۳-۲: مقایسه‌ی انواع سرطان از لحاظ میزان ابتلا و مرگ و میر (WHO).....
۳۰	نمودار ۴-۲: مقایسه‌ی میزان ابتلا و مرگ و میر ناشی از سرطان سینه بین نژادهای مختلف.....
۳۱	نمودار ۵-۲: میزان ابتلا و مرگ و میر ناشی از سرطان سینه در نقاط مختلف جهان (WHO).....
۳۴	نمودار ۶-۲: شیوع سرطان سینه در ایران، به تفکیک جنسیت و در گروه‌های سنی مختلف (۱۳۸۴).....
۹۲	نمودار ۱-۵: توزیع داده‌ها بر اساس نوع تومور.....
۹۲	نمودار ۲-۵: توزیع داده‌ها بر اساس سن.....
۱۰۹	نمودار ۱-۷: انتخاب پارامترهای a_1 ، a_2 و a_3
۱۱۰	نمودار ۲-۷: نتیجه‌ی افزایش T روی مجموعه‌ی آموزش.....
۱۱۰	نمودار ۳-۷: نتیجه‌ی افزایش T روی مجموعه‌ی آزمون.....
۱۱۲	نمودار ۴-۷: مقایسه‌ی کارایی الگوریتم پیشنهادی در تشخیص تمام رده‌ها با انتخاب رده‌بندهای پایه‌ی مختلف.....
۱۱۲	نمودار ۵-۷: مقایسه‌ی حساسیت الگوریتم پیشنهادی با انتخاب رده‌بندهای پایه‌ی مختلف (به تفکیک رده).....
۱۱۳	نمودار ۶-۷: مقایسه‌ی دقت الگوریتم پیشنهادی با انتخاب رده‌بندهای پایه‌ی مختلف (به تفکیک رده).....
۱۱۵	نمودار ۷-۷: تأثیر انتخاب نسبت‌های مختلف بین مجموعه‌ی آموزش و آزمون بر میزان حساسیت.....

- نمودار ۷-۸: تأثیر انتخاب نسبت‌های مختلف بین مجموعه‌ی آموزش و آزمون بر میزان دقت. ۱۱۵
- نمودار ۷-۹: مقایسه‌ی کارایی الگوریتم پیشنهادی با سایر الگوریتم‌های پایه (حاصل اعمال مدل روی مجموعه‌ی آزمون)..... ۱۱۶
- نمودار ۷-۱۰: مقایسه‌ی کارایی الگوریتم‌ها در تشخیص رده‌های مختلف..... ۱۱۷
- نمودار ۷-۱۱: مقایسه‌ی کارایی الگوریتم پیشنهادی با سایر الگوریتم‌های پایه (حاصل اعمال مدل روی مجموعه‌ی آموزش)..... ۱۱۹
- نمودار ۷-۱۲: مقایسه‌ی نتایج حاصل از اجرای الگوریتم پیشنهادی و سایر روش‌های یادگیری در داده‌های نامتوازن روی مجموعه‌ی آزمون..... ۱۲۱
- نمودار ۷-۱۳: مقایسه‌ی حساسیت الگوریتم پیشنهادی و سایر روش‌های یادگیری در داده‌های نامتوازن روی مجموعه‌ی آزمون به تفکیک رده‌ها..... ۱۲۲
- نمودار ۷-۱۴: مقایسه‌ی دقت الگوریتم پیشنهادی و سایر روش‌های یادگیری در داده‌های نامتوازن روی مجموعه‌ی آزمون به تفکیک رده‌ها..... ۱۲۲
- نمودار ۷-۱۵: مقایسه‌ی معیار F حاصل از اجرای الگوریتم پیشنهادی و سایر روش‌های یادگیری در داده‌های نامتوازن روی مجموعه‌ی آزمون به تفکیک رده‌ها..... ۱۲۴
- نمودار ۷-۱۶: استفاده از معیارهای صحت، حساسیت، دقت و معیار F ، بعنوان تابع تناسب (مقایسه‌ی نتایج بر اساس معیارهای G و میانگین ماکروی حساسیت، دقت و معیار F)..... ۱۲۶
- نمودار ۷-۱۷: مقایسه‌ی نتایج حاصل از DPSO در صورت استفاده از توابع تناسب مختلف بر اساس حساسیت در تشخیص رده‌ها..... ۱۲۶
- نمودار ۷-۱۸: مقایسه‌ی نتایج حاصل از DPSO در صورت استفاده از توابع تناسب مختلف بر اساس دقت در تشخیص رده‌ها..... ۱۲۷
- نمودار ۷-۱۹: مقایسه‌ی کارایی الگوریتم DPSO با الگوریتم پیشنهادی..... ۱۲۸

فهرست نشانه‌ها و اختصارها (Abbreviations)

PSO	Particle Swarm Optimization
DPSO	Discrete Particle Swarm Optimization
SVM	Support Vector Machine
LS-SVM	Least Square Support Vector Machine
k-NN	k Nearest Neighbor
C&RT	Classification and Regression Trees
CHAID	Chi-squared Automatic Interaction Detector
QUEST	Quick, Unbiased, Efficient Statistical Tree
G-mean	Geometric mean
SMOTE	Synthetic Minority Over-Sampling
RUS	Random Under-Sampling
ROS	Random Over-Sampling
OSS	One-Sided Selection
EE	EasyEnsemble
MB	ModifiedBagging
MEE	Multi-class EasyEnsemble
NN	Neural Network
KDD	Knowledge discovery in databases
WHO	World Health Organization
SEER	Surveillance, Epidemiology and End Results
NCI	National Cancer Institute
BMI	Body Mass Index
FNA	Fine Needle Aspiration
MHT	Menopausal Hormone Therapy
CRM	Customer Relationship Management

فصل ۱: مقدمه

۱-۱. موضوع تحقیق

سرطان^۱ نوعی بیماری است که در آن سلول‌های بدن، به صورت غیر قابل کنترل رشد می‌کنند، تغییر می‌کنند و تکثیر می‌شوند. سرطان، معمولاً با بخشی از بدن که درگیر این بیماری است، نامگذاری می‌شود. در نتیجه، سرطان سینه^۲ به معنای رشد غیرعادی سلول‌ها در بافت سینه است. گروهی از سلول‌هایی که به سرعت تکثیر شده‌اند، ممکن است یک توده^۳، جرم^۴ یا بافت اضافی^۵ ایجاد کنند. این توده‌ها، تومور^۶ نامیده می‌شوند. توده‌ها می‌توانند سرطانی (بدخیم^۷) یا غیرسرطانی (خوش‌خیم^۸) باشند. تومورهای سرطانی، سرایت می‌کنند و بافت‌های سالم بدن را نابود می‌کنند. در واقع سرطان سینه به معنای تومورهای بدخیمی است که از سلول‌های سینه گسترش یافته‌اند. این سرطان، در میان مردان و زنان دیده شده است، گرچه موارد ابتلا در میان مردان نادر است. دلیل واقعی بروز سرطان سینه، مشخص نیست اما برخی فاکتورها، ریسک ابتلا به سرطان سینه را افزایش می‌دهند^۹.

سرطان سینه، دلیل اصلی مرگ زنان ۴۰ تا ۵۵ ساله و دومین عامل مرگ زنان (پس از سرطان ریه^{۱۰}) است^{۱۱}. طبق آمار سازمان بهداشت جهانی^۱، سالانه بیش از ۱.۲ میلیون زن در

¹ Cancer

² Breast Cancer

³ Lump

⁴ Mass

⁵ Extra tissue

⁶ Tumor

⁷ Malignant

⁸ Benign

⁹ Cancer Council Victoria. **Breast cancer: for people with cancer, their families and friends.** (Melbourne: Cancer Council Victoria, 2009), Page 1.

American Cancer Society. **Cancer Facts & Figures 2009.** (Atlanta: American Cancer Society, 2009), Page 1.

American Cancer Society. **Breast Cancer Facts & Figures 2009-2010.** (Atlanta: American Cancer Society, 2010), Page 1.

¹⁰ Lung Cancer

¹¹ http://www.imaginis.com/breasthealth/breast_cancer.asp

سراسر دنیا، مبتلا به سرطان سینه، تشخیص داده می‌شوند.^۲ در ایران نیز سرطان سینه، شایع-ترین سرطان در میان زنان است.^۳ خوشبختانه در سال‌های اخیر، نرخ مرگ و میر ناشی از سرطان سینه، به دلیل تاکید بیشتر روی تکنیک‌های موثرتر تشخیص و معالجه، کاهش یافته است. فاکتور اصلی در این روند، تشخیص سریع و صحیح است.^۴

تکنیک‌های داده‌کاوی^۵ علاوه بر این که می‌توانند در تشخیص سریع این بیماری موثر باشند، می‌توانند از طریق شناسایی عوامل موثر در بروز توده‌های خوش‌خیم و بدخیم، به پیش‌گیری از این بیماری کمک کنند و به این ترتیب باعث ارتقای سلامت جامعه شده و از تحمیل هزینه‌های سنگین تشخیص و درمان بیماری، بر بیماران جلوگیری کنند.

۲-۱. اهمیت و ضرورت تحقیق

پس از ۲۱ قرن، هنوز سرطان دلیل اصلی مرگ در جهان است و سرطان سینه به دومین عامل اصلی مرگ سرطانی در میان زنان تبدیل شده است.^۶ گرچه، گسترش تکنولوژی‌های پزشکی در دهه‌ی گذشته، موجب کاهش میزان مرگ و میر ناشی از این بیماری شده است. به دلیل تشخیص زود هنگام و درمان بهینه، میزان بهبودی افزایش یافته است. حدود ۹۷ درصد زنان می‌توانند به مدت ۵ سال یا بیشتر بهبود یابند.^۷ بعلاوه میزان بهبودی ۵ ساله برای بیماران با هر دو نوع تومور به میزان ۱۰ تا ۱۵ درصد افزایش یافته است و این به دلیل تشخیص زودهنگام بیماری است. بنابراین روش موثر برای کاهش مرگ در اثر سرطان سینه، تشخیص زودهنگام این بیماری است.^۸ تشخیص زودهنگام نیز به یک روش تشخیص صحیح و قابل

¹ World Health Organization (WHO)

² Mehmet Fatih Akay. "Support vector machines combined with feature selection for breast cancer diagnosis". **Expert Systems with Applications**, 36, 2(2009):3240-3247.

³ Shadi Kolahdoozan et al. "Five Common Cancers in Iran". **Archives of Iranian Medicine**, 13, 2(2010):143 - 146.

⁴ Mehmet Fatih Akay, 2009.

⁵ Data Mining

⁶ Wei-Chang Yeh, Wei-Wen Chang & Yuk Ying Chung. "A new hybrid approach for mining breast cancer pattern using discrete particle swarm optimization and statistical method". **Expert Systems with Applications**, 36, 4(2009):8204-8211.

⁷ Dursun Delen, Glenn Walker & Amit Kadam. "Predicting breast cancer survivability: a comparison of three data mining methods". **Artificial Intelligence in Medicine**, 34, 2(2005):113-127.

⁸ Wei-Chang Yeh, Wei-Wen Chang & Yuk Ying Chung, 2009.

اعتماد نیاز دارد که به پزشکان اجازه می‌دهد تومورهای خوش‌خیم را از نوع بدخیم آن شناسایی کنند. بنابراین یافتن یک متد تشخیص صحیح و موثر و همچنین تشخیص عوامل خطر در بروز این بیماری، بسیار با اهمیت است.

مشکلی که در تشخیص این بیماری وجود دارد این است که در ابتدا ممکن است توده‌های سرطانی، بسیار کوچک بوده و قابل لمس نباشند. بدیهی است تا زمانی که یک توده‌ی کوچک، بزرگ شده و قابل لمس شود، بیماری نیز پیشرفت کرده و این تأخیر در تشخیص بیماری، روند درمان را با مشکل مواجه خواهد کرد تا جایی که حتی ممکن است به مرگ بیمار منجر شود. بنابراین نیاز است تا از روش‌های پیشرفته‌تری برای تشخیص زودهنگام سرطان سینه استفاده شود. روش‌های رایج برای تشخیص زودهنگام سرطان سینه عبارتند از: ماموگرافی^۱، سونوگرافی^۲ و نمونه‌برداری. این روش‌ها علاوه بر هزینه‌ی بالایی که بر بیمار تحمیل می‌کنند، ممکن است مضراتی را برای بیمار در پی داشته باشند. مخصوصاً در روش ماموگرافی که از اشعه‌ی ایکس برای تولید تصاویر استفاده می‌شود. به طوری که این روش، معمولاً برای افراد زیر ۴۰ سال که احتمال ابتلا به سرطان سینه در آنان کمتر است، توصیه نمی‌شود.^۳

ماموگرافی، یکی از رایج‌ترین روش‌های تشخیص سرطان سینه است. این روش می‌تواند ناهنجاری‌ها را با صحت^۴ ۸۵ تا ۹۰ درصد شناسایی کند.^۵ اما مشکلی که در این روش وجود دارد، این است که نتیجه‌ای که متخصصان مختلف از این تصاویر استنباط می‌کنند، ممکن است متفاوت باشد و این موضوع چالش دیگری را در امر تشخیص سرطان سینه ایجاد می‌کند.^۶ در چنین مواردی، از روش آسپیراسیون سوزنی ظریف^۷ که یکی از روش‌های نمونه‌برداری است،

¹ Mamography

² Ultrasound

³ محمدرضا صفایی کشتگر و راب استین. اطلاعات عمومی برای پیش‌گیری و نقاشی‌های بالینی برای بیماران مبتلا به سرطان پستان. مریم بزرگر و عطیه اکبری، تهران: پیام سفید، چاپ چهارم، (۱۳۸۸)، صفحه ۱۰

⁴ Accuracy

⁵ Shieu-Ming Chou et al. "Mining the breast cancer pattern using artificial neural networks and multivariate adaptive regression splines". **Expert Systems with Applications**, 27, 1(2004):133-142.

⁶ Joann G. Elmore et al. "Variability in radiologists interpretation of mammograms". **New England Journal of Medicine**, 331, 22(1994):1493-1499.

⁷ Fine Needle Aspiration (FNA)