



گزارش سمینار

کاربرد داده کاوی پروتئوم در دسته بندی سرطان و

کشف زیست نشانگرها

دانشجو

رسول نوروزی

استاد راهنما

دکتر امیر البدوی

تابستان ۱۳۹۶

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

چکیده

سرطان بیماری بسیار پیچیده است که در سطح مولکولی در بدن پدیدار می‌شود. سرطان‌ها در شکل سلول‌ها، الگو و کمیت زیست‌نشانگرها به خصوص پروتئین‌های موجود در بافت و سلول تأثیر می‌گذارند. در این میان علم پروتئومیکس که به بررسی الگوها و میزان بیان پروتئین‌ها در بافت و سلول می‌پردازد می‌تواند یکی از کلیدهای حل مسئله سرطان باشد. پروتئومیکس شامل رویکردها و فن‌آوری‌های متنوعی است که در این بین رویکرد پروتئومیکس شناختی و مقایسه‌ای و فن‌آوری طیف‌سنج جرمی در پیوستگی استفاده از روش‌ها تحلیلی سطح بالا همچون داده‌کاوی در دسته‌بندی سرطان و کشف زیست‌نشانگرهای مربوطه می‌تواند بسیار راهگشا باشد. در این گزارش به بررسی رویکردهای علم پروتئومیکس و فن‌آوری‌های بکار رفته در آن خواهیم پرداخت و سپس تمرکز مطالعات را بر روی پروتئومیکس شناختی و مقایسه‌ای با داده‌های حاصل از فن‌آوری طیف‌سنج خواهیم گذاشت. پس از آن به علم داده‌کاوی و نحوه استفاده از داده‌کاوی با داده‌های حاصل از طیف‌سنج جرمی برای دسته‌بندی و کشف زیست‌نشانگرها خواهیم پرداخت و در آخر هم مروری بر تحقیقات گذشته خواهیم داشت.

کلیدواژه‌ها: طیف‌سنج جرمی، سرطان، داده‌کاوی، زیست‌نشانگر

فهرست مطالب

عنوان	شماره صفحه
فصل اول مفاهیم و کلیات	
۱/۱ مقدمه	۲
۱/۲ انسان از منظر زیست‌شناسی	۲
۱/۲/۱ اسیدهای نوکلئوتید	۳
۱/۲/۲ پروتئین‌ها	۴
۱/۳ سرطان	۴
۱/۳/۱ اساس مولکولی	۵
۱/۴ بیان مسئله	۶
۱/۵ ضرورت مسئله	۷
۱/۶ تعاریف و اصطلاحات	۸
آنتیبادی	۸
آنتیژن	۸
پلیپتاید	۸
جرم مولکولی (M)	۹
نقطه ایزو الکتریک	۹
اثرانگشت جرم-پتاید	۹
طیف جرمی	۹
توصیف پروتئین‌های بیان‌شده	۹
۱/۷ آرایش کلی گزارش	۱۰
۱/۸ خلاصه فصل	۱۱

فصل دوم پروتئومیکس و کاربردهای علوم کامپیوتری در داده های حاصل از طیف سنج جرمی

۱۳	۲/۱ مقدمه
۱۳	۲/۲ انواع پروتئومیکس
۱۳	۲/۲/۱ پروتئومیکس شناختی و مقایسه ای
۱۴	۲/۲/۲ تغییرات پس از ترجمه پروتئین ها
۱۴	2.2.3 مکان یابی پروتئین ها
۱۴	2.2.4 برهم کنش پروتئین ها
۱۵	۲/۳ روشهای نمونه گیری برای کشف نشانگرهای سرطانی
۱۶	۳/۲/۱ مایعات زیست پذیر
۱۶	۲/۳/۲ بافت
۱۶	۲/۳/۳ سل لاین
۱۷	۲/۴ فن آوری های آنالیز پروتئومیکس
۱۷	۲/۴/۱ ریزآرایه های پروتئینی
۱۸	۲/۴/۲ ژل الکتروفورز دوبعدی پلی آکریل آمید (2D-PAGE)
۱۸	۲/۴/۳ طیف سنج جرمی
۲۰	۲/۵ فرصت ها و چالش های آنالیز داده ها در طیف سنج جرمی
۲۱	۲/۶ آنالیز داده های طیف سنج جرمی
۲۲	۲/۶/۱ انتخاب موجک
۲۲	۲/۶/۲ موتورهای جستجو
۲۴	۲/۷ خلاصه فصل

فصل سوم داده کاوی و کاربرد آن در دسته بندی و کشف زیست نشانگرهای سرطانی

۲۶	۳/۱ مقدمه
۲۶	۳/۲ داده کاوی و کاربرد آن در پروتئومیکس
۲۷	۳/۲/۱ پیش پردازش داده ها
۲۷	۳/۲/۲ نفرین بعد
۲۸	۳/۲/۳ روش های انتخاب ویژگی
۲۹	۳/۲/۴ مدل های دسته بندی
۳۱	۳/۲/۵ واریسی اعتبار
۳۲	۳/۲/۶ بررسی عملکرد مدل
۳۳	۳/۳ مروری بر ادبیات کاربرد داده کاوی در دسته بندی و تشخیص زیست نشانگرهای سرطان
۳۳	۳/۳/۱ بررسی مقالات و تحقیقات صورت گرفته
۳۷	۳/۴ خلاصه فصل

فصل چهارم جمع بندی و نتیجه گیری

۳۹	۴/۱ مقدمه
----	-----------

۳۹	_____	۴/۲ مروری بر فصل‌های گذشته
۴۱	_____	۴/۳ بررسی چالش‌ها و پیشنهاد فرصت‌ها
۴۲	_____	۴/۴ داده‌های تحقیق
۴۳	_____	۴/۵ خلاصه فصل
		مراجع
۴۵	_____	مراجع فارسی
۴۵	_____	مراجع انگلیسی
		پیوست‌ها
۴۹	_____	پیوست ((الف))
۵۰	_____	پیوست ((ب))

فهرست شکل ها و نمودارها

موضوع	شماره صفحه
شکل ۱-۱ از سلول تا دی ان ای	۳
شکل ۱-۲ فرآیند ترجمه شدن و ساخت پروتئین	۳
شکل ۱-۳ ایش بینی نرخ رشد سرطان تا سال ۲۰۲۰	۸
شکل ۲-۱ جریان کاری مراحل شناسایی زیست نشانگرها	۱۵
شکل ۲-۲ روش ساندویچ ریز آرایه	۱۸
شکل ۲-۳ روش آرایه های-تک آنتی بادی	۱۸
شکل ۲-۴ طیف جرمی حاصل از سروم خون با فناوری SELD-MS	۲۰
شکل ۲-۵ جریان آنالیز داده های حاصل از طیف سنج جرمی	۲۱
شکل ۳-۱ جریان کاری داده کاوی بر روی داده های طیف جرمی	۲۸
شکل ۳-۲ درخت تصمیم	۳۰
شکل ۳-۳ بردارهای ماشین پشتیبان	۳۱
شکل ۳-۴ شبکه های عصبی	۳۱
شکل ۳-۵ نمونه ای از منحنی ROC	۳۲
شکل ۴-۱ نمونه ای از داده های آزمایشگاه پروتئومیکس دانشگاه کالیفرنیا جنوبی	۴۲

فهرست جدول‌ها

موضوع	شماره صفحه
جدول ۱-۱ اختصارات	۱۰
جدول ۲-۱ ماتریس شدت M/Z	۲۳
جدول ۲-۲ ماتریس بیان پروتئین	۲۴
جدول ۳-۱ مقایسه مدل‌های دسته‌بند	۳۱
جدول ۳-۲ مرور تحقیقات گذشته دسته‌بندی سرطان و کشف زیست‌نشانگر	۳۵
جدول نرم‌افزارهای مورد استفاده در داده‌کاوی داده‌های طیف سنج جرمی	۴۹
فهرست واژگان انگلیسی به فارسی و بالعکس	۵۰

فصل اول

مفاهیم و کلیات

۱,۱ مقدمه

سرطان بیماری پیچیده‌ای است که در سطح مولکولی در بدن انسان ظاهر می‌شود بنابراین برای شناخت بهتر سرطان نیاز است در ابتدای این فصل نگاهی به انسان در سطح مولکولی داشته باشیم در ادامه سرطان را از دیدگاه مولکولی و ارتباط زیست‌نشانه‌ها با سرطان در بدن انسان را تشریح خواهیم کرد در قدم بعدی در راستای مطالب عنوان‌شده بیان مسئله خود و دلیل استفاده از پروتئین‌ها به‌عنوان زیست‌نشانه‌ها منتخب خواهیم پرداخت. در قسمت بعدی از ضرورت پرداختن به بیماری سرطان خواهیم گفت سپس برخی از اصطلاحات که برای وارد شدن به دنیای علم پروتئومیکس نیاز است را تعریف خواهیم کرد، همچنین علائم و اختصارات استفاده شده را فهرست خواهیم کرد پس‌از آن ساختار کلی تحقیق را شرح خواهیم داد و در آخر خلاصه‌ای از مطالب بیان‌شده در فصل را مرور خواهیم کرد.

۱,۲ انسان از منظر زیست‌شناسی

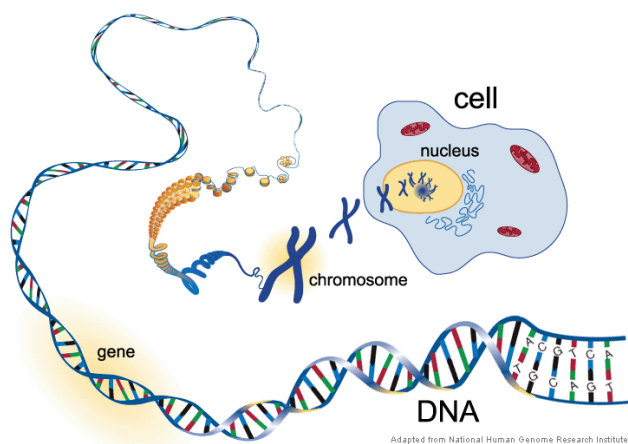
از منظر زیست‌شناسی می‌توان انسان را از منظر چهار ابرمولکول^۱: کربوهیدرات‌ها، پروتئین‌ها، اسیدهای نوکلئوتید و لیپیدها نگریست. گرچه مولکول‌های بسیار دیگری در ساختار بدن وجود دارند اما موارد ذکرشده مهم‌ترین و تأثیرگذارترین آن‌ها می‌باشند. این ابرمولکول‌ها به‌صورت پلیمرهای بلندی می‌باشند که از ترکیب مونومورها (همچنان که مرواریدهای یک گردنبند به کمک یک‌رشته به هم متصل شده‌اند) با کمک سنتز پس‌آش^۲ ایجادشده‌اند. هرکدام از این ابرمولکول‌های بیولوژیکی دارای زیرمجموعه‌ها و به‌تبع آن دارای خواص و وظایف منحصربه‌فردی در بدن می‌باشند. به‌طور مثال کربوهیدرات‌ها شامل سه دسته منوساکاریدها، دی‌ساکاریدها و پلی‌ساکاریدها می‌شوند که وظایفی همچون منبع ذخیره انرژی و تشکیل ساختار دیواره‌های برخی از سلول‌ها را دارا می‌باشند. همچنین لیپیدها که شامل دسته‌های: چربی‌ها و روغن‌ها، واکس‌ها، فسفولیپیدها و آستروئیدها می‌باشند که وظایفی همچون ذخیره انرژی، ایجاد عایق برای سلول‌ها، ایجاد غشای سلولی و... بر عهده‌دارند؛ اما در ارتباط با اسیدهای نوکلئوتید و پروتئین‌ها با توجه به این‌که اساس مطالعه ما در این تحقیق می‌باشند نگاه دقیق‌تر و موشکافانه‌تر خواهیم داشت.

¹ macromolecule

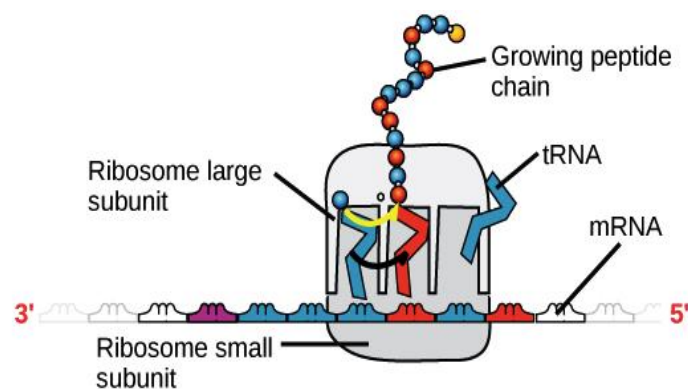
² dehydration synthesis

۱،۲،۱ اسیدهای نوکلئوتید

اسیدهای نوکلئوتید از واحدهایی به نام نوکلئوتید تشکیل شده‌اند که در طبیعت به دو شکل اسید دئوکسی ریبونوکلئیک (DNA) و اسید ریبونوکلئیک (RNA) یافت می‌شوند. هر واحد نوکلئوتید از سه بخش حلقه‌های نیتروژن دار، قند با پنج کربن و حداقل یک گروه فسفات تشکیل شده است. DNA معمولاً به صورت مارپیچ دوگانه یافت می‌شوند. چهار نوع مختلف از مولکول‌های نوکلئوتیدی در زنجیره DNA یافت می‌شوند. این چهار مولکول از نظر ساختمانی شباهتی بسیاری به هم دارند لیکن دارای تفاوت‌های کافی جهت تشخیص‌شان توسط فرآیندهای مربوطه نیز می‌باشند. قسمت‌های ویژه از هر چهار نوکلئوتید که به صورت حلقوی هستند، باز نامیده می‌شود. این حلقه‌ها در اثر برقراری پیوند شیمیایی بین اتم‌های کربن و نیتروژن به وجود آمده‌اند. سایر اتم‌ها (فسفر، هیدروژن و اکسیژن) به طور مستقیم به این حلقه‌ها متصل شده‌اند چهار حلقه مختلف آدنین، تیمین، گوانین و سیتوزین نامیده شده و توسط حروف A, T, C, G نمایش داده می‌شوند. واحدهای نوکلئوتیدی توسط گروه فسفات به هم متصل می‌شوند. ترکیب‌های نوکلئوتیدی قواعدی دارند. تیمین همراه با آدنین و سیتوزین همراه با گوانین جفت می‌شوند. (شکل ۱-۱) این چهار حلقه همچون حروف الفبا برای یادداشت‌برداری و نگهداری اطلاعات ژنتیکی عمل می‌کند. (علی پور، محمد، ۱۳۸۴) RNA برخلاف DNA معمولاً از یک رشته بلند تشکیل شده است. یک نوکلئوتید در RNA می‌تواند شامل ریبوز (قند پنج کربن)، یکی از چهار حلقه نیتروژنی (A, U, G یا C) و یک گروه فسفات شود. RNA انواع مختلفی دارد که هر کدام دارای وظایفی خاصی در سلول می‌باشند مانند tRNA, mRNA و...



شکل ۱-۱ از سلول تا دی‌ان‌ای (وبسایت سازمان ملی ژنتیک انسانی آمریکا)



شکل ۱-۲ فرآیند ترجمه شدن و ساخت پروتئین (وبسایت سازمان ملی ژنتیک انسانی آمریکا)

اسیدهای نوکلئوتید و بخصوص DNA ابرمولکول‌های کلیدی برای ادامه حیات انسان می‌باشند. DNA اطلاعات وراثتی که از والد به فرزند منتقل می‌شود را دارا است همچنان که دربردارنده اطلاعات چگونگی، میزان و زمان ساخت پروتئین‌ها و عملکرد سلول‌ها، بافت‌ها و ارگان‌ها است. در یوکاریت‌ها (موجودات دارای هسته در سلول)، DNA معمولاً به تعدادی از قطعات بسیار بلند و باریک خطی به نام کروموزوم که دارای تعداد مشخصی می‌باشند شکسته می‌شود. برای مثال تعداد کروموزوم‌های سلول انسان ۲۳ جفت (۴۶ تا) است؛ اما شکل کروموزوم در پروکاریت‌ها (موجودات فاقد هسته در سلول مانند ویروس‌ها) به صورت مدور و حلقه‌ای شکل است. کروموزوم‌ها هر کدام می‌توانند حاوی ۱۰ تا ۱۰۰ هزار ژن باشند که هر کدام دستورالعمل ساخت فراورده مورد نیاز سلول را مهیا می‌کنند مانند دستورالعمل ساخت پروتئین‌ها که توالی اسیدهای آمینه آن را مشخص می‌کنند که به این عمل اصطلاحاً

کدگذاری^۳ گفته می‌شود. قبل از این که این اطلاعات استفاده شود برای ساخت پروتئین‌ها، RNA ها یک کپی از ژن‌ها تهیه می‌کنند سپس این کپی از ژن‌ها به ریبوزوم، ماشین سلولی که توالی RNA را می‌خواند (فرآیند ترجمه شدن^۴) و از آن برای ساخت پروتئین‌ها استفاده می‌کند. (شکل ۱-۲)

۱،۲،۲ پروتئین‌ها

پروتئین‌ها از فراوان‌ترین مولکول‌های ارگانیک در بدن موجودات زنده هستند و بیشترین تنوع در ساختار و عملکرد را در بین ابرمولکول‌های زیستی دارند. یک سلول به تنهایی می‌تواند شامل صدها هزار پروتئین با عملکرد مجزا و منحصر باشد. تمام پروتئین‌ها از یک یا تعدادی از زنجیره‌های آمینواسید تشکیل شده است که هر کدام از آن‌ها یک پلی‌پتاید^۵ نامیده می‌شود. آمینواسیدها مونومورهایی هستند که تشکیل‌دهنده پروتئین‌ها می‌باشند. در مجموع تابه‌حال ۲۰ نوع آمینواسید در پروتئین‌ها یافت شده است.

پروتئین‌ها دارای نقش‌های متنوعی در سلول‌ها و ارگان‌های بدن می‌باشند؛ که دودسته مهم آن‌ها به شرح زیر است:

آنزیم‌ها: دسته‌ای از پروتئین‌ها می‌باشند که به‌عنوان کاتالیزور برای سرعت بخشی به واکنش‌های بیوشیمیایی عمل می‌کنند. هر آنزیم یک یا چند مولکول را برای کاتالیز آن در واکنش‌های بیوشیمیایی به‌عنوان زیرمجموعه خود به رسمیت می‌شناسد. به‌طور مثال آنزیم بزاق آمیلاز^۶، برای شکستن مولکول آمیلاز (نوعی نشاسته) به قطعات کوچک قندی است.

هورمون‌ها: دسته دیگری از پروتئین‌های مهم در بدن موجودات زنده می‌باشند که فرآیندهای فیزیولوژی بخصوصی مانند رشد، توسعه، متابولیسم و تولیدمثل را به عهده‌دارند. برخی از هورمون‌ها زیرمجموعه لیپیدها هستند و به پایه-استروئید شناخته می‌شوند و مابقی زیرمجموعه پروتئین‌ها معروف به پایه-پپتید قرار می‌گیرند. (Garrett and Grisham, 2005, Watson et al., 2003)

۱،۳ سرطان

سرطان یک بیماری ژنتیکی است که ۲۷۷ نوع بیماری را شامل می‌گردد. تغییرات ژنتیکی باعث ازهم‌گسیخته شدن نظم طبیعی تقسیم و تمایز سلول‌ها می‌شود. همچنین در محیط زیست ما بیش از یکصد هزار نوع مواد شیمیایی وجود دارد که فقط ۳۵ هزار از آن آنالیز شده و حدود ۳۰۰ عدد از آن‌ها تولید سرطان می‌کند. هنوز ۶۵ هزار مواد شیمیایی باقیمانده در طبیعت آزمایش نشده است. در حال حاضر بیش از ۵۰ درصد بیماری‌های سرطانی را معالجه می‌نماییم مخصوصاً اگر این بیماری در مراحل اولیه تشخیص داده شوند. (پارسا، ۲۰۱۲) موارد زیر برخی از انواع سرطان می‌باشند که در سلول‌های به خصوصی در بدن شروع می‌شوند:

- کارسینوما: در سلول‌های مخاطی به وجود می‌آید که توانایی حمله به ماهیچه‌های اطراف خود را نیز دارا است. سرطان‌های سینه، پروستات، ریه و روده بزرگ از نوع کارسینوما می‌باشند.

³ encode

⁴ translation

⁵ polypeptide

⁶ Salivary amylase

- سارکوما: یک تومور بدخیم در استخوان یا بافت‌های نرم است. سارکوما اساس سرطان‌های استخوانی است.
- لنفاوی: سرطانی است که غدد لنفاوی را مورد حمله قرار می‌دهد و قادر به انتقال به سایر بافت‌های بدن نیز است.
- سرطان خون^۷: نوعی سرطان که بافت سلول‌های خون، گلبول‌های سفید خون و مغز استخوان را شکل می‌دهد.

از عوامل اساسی ایجاد سرطان می‌توانیم به: الکل، دخانیات، اضافه‌وزن و چاقی، مشکلات در سیستم ایمنی بدن، عفونت و ... اشاره کنیم. (Shukla et al., 2016)

سرطان از روش‌های جراحی، شیمی‌درمانی، اشعه درمانی، ایمنو درمانی، ژن درمانی و یا تلفیقی از آن‌ها درمان می‌شود.

۱,۳,۱ اساس مولکولی

اما آنچه در سرطان برای ما بیشترین اهمیت را دارا است اساس مولکولی و سازوکار ایجاد آن است. در یک تعریف ساده می‌شود سرطان را با فعالیت نابجای پروتئین‌های کنترل‌کننده چرخه سلول که موجب ازدیاد در تقسیم سلولی می‌شود نامید. یک سلول سرطانی تفاوت‌های با یک سلول سالم دارا است. سلول‌های سرطانی بدون نیاز به سیگنال‌های لازم برای تکثیر برخلاف سلول‌های سالم، شروع به تقسیم می‌کنند. همچنین سیگنال‌هایی را که برای توقف تکثیر می‌باشند را نادیده می‌گیرند. سلول‌های سالم تنها قادر به تکثیر بین ۴۰-۶۰ بار در طول عمر خود هستند اما برای سلول‌های سرطانی این چرخه به مراتب بیشتر است. سلول‌های سرطانی قابلیت حرکت به سمت سایر بافت‌ها و ارگان‌های بدن را دارا می‌باشند که به آن اصطلاحاً متاستاز^۸ گفته می‌شود همچنان که می‌توانند موجب رشد ماهیچه‌های خونی شوند که اصطلاحاً رگ زایی^۹ گفته می‌شود. سلول‌های سرطانی همچنان از فرمان خودکشی سلولی^{۱۰} که تحت شرایط خاصی مانند معیوب شدن سلول‌ها برای آن برنامه‌ریزی می‌شوند سرپیچی می‌کنند. (Sosa et al., 2014)

سرطان‌های مختلف شامل جهش‌های ژنتیکی^{۱۱} مختلفی نیز می‌باشند و هر کدام تغییرات بخصوصی در ساختار ژن‌ها ایجاد می‌کنند. به‌طور کلی جهش ژنتیکی در دو نوع از تنظیم‌کننده‌های چرخه سلولی موجب توسعه سرطان می‌شود:

۱. تنظیم‌کننده‌های مثبت (آنکوژن‌ها^{۱۲}) که ممکن است بیش از اندازه فعالیت کنند. سلول‌ها تا زمانی که بیش‌فعالی نرسیده باشند پروتو-آنکوژن نامیده می‌شوند ولی ممکن است با تغییر توالی اسیدآمینه‌های پروتئین موجب تغییر در شکل پروتئین موجب همیشه فعال پروتئین شوند. از سوی دیگر ممکن است موجب تقویت شوند، به این صورت که یک سلول مقادیر زیادی کپی از یک ژن تهیه کند و در نتیجه موجب ازدیاد تولید در پروتئین‌ها شود. همچنین بسیاری از پروتئین‌ها سیگنال‌ها و دستورات تکثیر را توسط پروتو-آنکوژن‌ها زمانی که شرایط مساعد باشند منتقل می‌کنند ولی ممکن است یک جهش باعث بیش‌فعالی پروتئین‌ها شود و سیگنال‌های تکثیر را حتی زمانی که شرایط تکثیر مناسب نیست انتقال دهد.

⁷ Leukemia

⁸ metastasis

⁹ angiogenesis

¹⁰ apoptosis

¹¹ mutation

¹² oncogenes

۲. تنظیم‌کننده‌های منفی که سرکوب‌کننده‌ها^{۱۳} تومور نیز نامیده می‌شوند ممکن است غیرفعال شوند. زمانی که DNA آسیب می‌بیند یک پروتئین به نام P53 چرخه سلولی را متوقف می‌کند. این توقف باعث می‌شود تا این پروتئین آنزیم‌های ترمیم‌کننده DNA را فعال کند. اگر آسیب ترمیم پیدا کرد این پروتئین اجازه ازسرگیری فعالیت چرخه سلولی را می‌دهد در غیر این صورت آخرین وظیفه خود یعنی موجب سازی خودکشی سلولی را فراهم می‌کند تا جهش ژنتیکی ادامه پیدا نکند. (Otto and Sicinski, 2017, Joerger and Fersht, 2016)

۱،۴ بیان مسئله

پس از تکمیل پروژه ژنوم انسان، تمرکز محققین به سمت شناخت ساختار، عملکرد و تعامل پروتئین‌ها و نقش آن‌ها در بیماری‌ها که توسط ژن‌ها تولید می‌شود معطوف گردیده است. این تغییر جهت به دلیل: ۱. سطح بیان mRNA معمولاً بیانگر میزان دقیق پروتئین‌های فعال در سلول نیست ۲. توالی ژن‌ها تغییرات پس از ترجمه اصلاحات پروتئین‌ها را که ممکن است برای عملکرد درست پروتئین‌ها ضروری باشد بیان نکند ۳. مطالعات ژنتیک قادر به توصیف فرآیندهای پویای سلولی نیست. (Li et al., 2004) پروتئوم به سری کامل پروتئین‌های بیان‌شده در یک لحظه خاص در یک سلول موردنظر اشاره دارد، بااین حال امروزه سطح این تعریف، از سلول به بافت، اندام و ارگانیسم نیز تعمیم داده شده است. (شیردل et al., 2013) به تبع آن پروتئومیکس یک واژه کلی بوده که به مطالعه پروتئین‌ها در ابعاد گسترده از جنبه‌های، پروتئومیکس شناختی و مقایسه‌ای، هستی‌شناسی^{۱۴} پروتئین‌ها، تعاملات پروتئین-پروتئین، مسیرهای متابولیکی و توصیف پروتئین‌ها است که پژوهش حاضر به پروتئومیکس شناختی و مقایسه‌ای پروتئین‌هایی که با طیف‌سنج جرمی حاصل شده‌اند می‌پردازد. (Lam et al., 2014)

تغییرات بالینی ممکن است در الگوی پروتئومیکس یک ارگان یا بافت بازتاب یابد؛ بنابراین امکان‌پذیر است که پروتئین‌های موجود در یک نمونه خاص برای تشخیص بیماران سرطانی از غیر سرطانی استفاده شود. (Li et al., 2004) با پیشرفت روش‌های نمونه‌گیری و تفکیک پروتئین‌ها میزان داده‌های قابل دسترسی به‌طور فزاینده‌ای افزایش یافته است که به‌سادگی قابل تصویر و تفسیر نیستند و با توجه این‌که در داده‌های پروتئینی معمولاً تعداد نمونه‌ها به دلیل دشواری و هزینه در نمونه‌گیری کم می‌باشند و ابعاد یا به عبارتی تعداد ویژگی‌های آن‌ها بسیار زیاد است، نیاز به ابزار تحلیلی سطح بالا همچون داده‌کاوی و الگوریتم‌های پیشرفته هوش مصنوعی و یادگیری ماشین احساس می‌شود. هدف اصلی و اساسی استخراج اطلاعاتی است که منجر به کشف الگوهای برای دسته‌بندی سرطان و شناسایی زیست‌نشانه‌هایی همچون پروتئین‌ها که به‌طور بالقوه ابزاری قدرتمند برای شناسایی، تشخیص و پیشگیری از بیماری‌ها به‌خصوص سرطان است. (Thomas et al., 2006)

¹³ Tumor suppressor

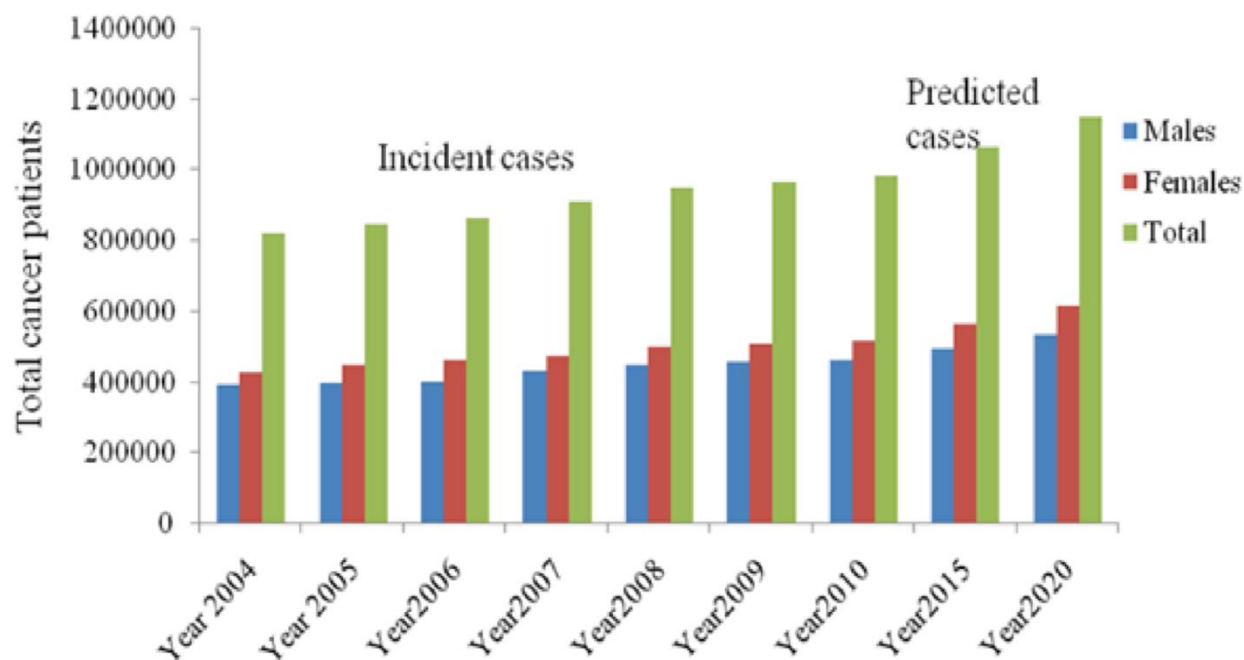
¹⁴ ontology

۱,۵ ضرورت مسئله

بنا بر آمار سازمان جهانی بهداشت سالانه ۸,۲ میلیون نفر انسان در سراسر جهان به دلیل سرطان می‌میرند که این مقدار ۱۳٪ از کل میزان مرگ‌ومیر جهان است. پیش‌بینی می‌شود که به این میزان تا ۷۰٪ تا دو دهه آینده اضافه شود. بیشترین میزان سرطان کشف‌شده به ترتیب مربوط به سرطان ریه (۱,۸ میلیون، ۱۳٪ از کل)، سرطان سینه (۱,۷ میلیون، ۱۲٪ از کل) و کلورکتال (۱,۴ میلیون نفر، ۹,۷٪ از کل) است. تخمین زده می‌شود تا سال ۲۰۲۵ این میزان به ۱۹,۳ میلیون نفر به دلیل رشد جمعیت افزایش یابد. (شکل ۱-۳) (Shukla et al., 2016)

تعداد بیماران سرطانی سال به سال رو به افزایش است و این خود یک معضل پزشکی است نه فقط از نظر بهداشت و درمان بلکه از نظر اقتصادی می‌تواند کشورها را تا حد ورشکستگی اقتصادی پیش ببرد. اگر سرطان‌ها در مرحله اول تشخیص داده شوند به‌طور کامل قابل‌معالجه هستند و اگر در مراحل دوم تشخیص داده شوند در حدود ۷۰ درصد شانس‌معالجه دارند و گر در مراحل سوم تشخیص داده شوند در حدود ۳۰ درصد شانس بهبودی دارند و اگر تشخیص در مرحله چهارم بوده باشد به‌طور حتم به بافت‌های دیگر گسترده شده است، شانس‌معالجه و بهبودی در حدود پنج درصد است که پنج سال ادامه حیات داشته باشد.

در سه دهه گذشته، محققین اطلاعات زیادی را درباره ژن‌ها و پروتئین‌ها و نقش آن‌ها در تولید سلول‌های طبیعی و سرطانی گزارش نموده‌اند. یکی از اکتشافات مهم آن‌ها، نقش ژن‌های جهش‌یافته در تولید سلول‌های سرطانی بوده است. عوامل محیطی باعث جهش‌های ژنتیکی می‌شوند در حال شناسایی هستند. با کمک از روش‌های مختلف مولکولی قادر هستیم که قدرت بیان ژن‌ها و پروتئین‌های معیوب را تعیین نماییم. حتی پیدا کردن زیست‌نشانگرهای جدیدی مانند پروتئین که شاخص یک نوع سرطان هستند در تشخیص زودرس و معالجه به‌موقع بیماری سرطان کمک‌های شایان توجهی را می‌نماید که این امر با به‌کارگیری علوم رایانه‌ای به‌خصوص داده‌کاوی ممکن می‌شود. پس از تعیین شکل فضایی پروتئین‌های معیوب می‌توان داروهای ضد سرطان جدیدی را ساخت که بتواند سلول‌های در حال سرطانی شدن را مورد هدف قرار بدهند تا از تولید و رشد آن‌ها به سلول‌های سرطانی جلوگیری شود. (پارسا، ۲۰۱۲)



شکل ۱-۳ پیش بینی نرخ رشد سرطان تا سال ۲۰۲۰ (Shukla et al., 2016)

۱,۶ تعاریف و اصطلاحات

قبل از ورود به فصل آتی جهت بررسی زوایای موضوع و درک مسئله نیاز به تعریف برخی از اصطلاحات عمده داریم که در اینجا به آن‌ها اشاره خواهیم کرد.

آنتی‌بادی^{۱۵}: آنتی‌بادی یا ایمونوگلوبولینها پروتئین‌هایی می‌باشند که توسط سیستم ایمنی بدن برای آمیختن با آنتی‌ژن‌های سلول‌های (مانند باکتری و ویروس) بیگانه جهت خنثی‌سازی آن‌ها استفاده می‌شوند.

آنتی‌ژن^{۱۶}: آنتی‌ژن‌ها مولکول‌های بیگانه هستند که موجب تولید آنتی‌بادی در ارگان می‌شوند.

پپتید: پپتیدها زنجیره‌های کوچکی از آمینواسیدها می‌باشند.

پلی‌پپتاید: یک زنجیره ابرمولکولی از آمینواسیدها که ممکن است یک پروتئین یا یکی اجزا پروتئین باشد.

آنزیم: پروتئین‌هایی می‌باشند که به‌عنوان کاتالیزور در فرآیندهای شیمیایی داخل بدن شرکت می‌کنند.

¹⁵ Antibody

¹⁶ antigen

جرم مولکولی (M):^{۱۷} جرم یک مولکول از یک ماده است. وزن مولکول عمدتاً به دالتون^{۱۸} بیان می‌شود. دالتون (Da) به‌عنوان یک دوازدهم جرم کربن-۱۲ (ایزوتوپ کربن) تعریف می‌شود و به‌طور تقریبی برابر است با: $1.66 \times 10^{-24} \text{g}$

نقطه ایزو الکتریک^{۱۹}: نقطه ایزو الکتریک یک پروتئین، میزان PH یک حلال است که در آن نقطه پروتئین از لحاظ بار الکتریکی در حالت خنثی است. PH، یا سطح هیدروژن، میزان اسیدی بودن یک حلال را اندازه‌گیری می‌کند. برای حلال‌های خنثی PH برابر هفت است.

اثرانگشت جرم-پپتاید^{۲۰}: مجموع جرم مولکولی پپتیدهایی که با تجزیه آنزیماتیک پروتئین تولید شده است. جرم پپتیدها به‌وسیله طیف جرمی مشخص می‌شود. در نتیجه اثرانگشت جرم-پپتاید برای شناسایی پروتئین‌ها با مقایسه این اثرانگشت با جرم پپت ایدهای حاصل از تجزیه توالی‌های پروتئین محاسبه شده در حالت تئوریک^{۲۱} در بانک داده‌ها است.

طیف جرمی^{۲۲}: طیف جرمی برداری است که شدت بار به جرم یک مورد یا تحت یک شرایط خاص را نمایش می‌دهد.

توصیف پروتئین‌های بیان شده^{۲۳}: یک بردار است که بیانگر شدت بار به جرم یک مورد خاص در میان موارد و شرایط متفاوت است.

¹⁷ Molecular mass

¹⁸ dalton

¹⁹ Isoelectric point

²⁰ Peptide-mass fingerprint

²¹ In silico approach

²² Mass spectrometry

²³ Protein expression profile

۱,۷ اختصارات

جدول ۱-۱ اختصارات

اختصار	واژه کامل
DNA	Deoxyribonucleic acid
RNA	Ribonucleic acid
MS	mass spectrometry
WHO	world health organization
NCBI	National Center for Biotechnology Information
DC	decision tree
SVM	support vector machine
ANN	artificial neural network
CV	cross-validation
RS	random sampling
PPI	protein-protein interaction
PPT	protein post translation
CNS	central neural system
MS-TOF	mass spectrometry Time Of Flight
SELDI-MS	surface enhanced laser mass spectrometry
MALDI-MS	matrix assisted laser ionization mass spectrometry
ESI-MS	electrospray ionization mass spectrometry

۱,۸ آرایش کلی گزارش

این تحقیق از چهار فصل تشکیل شده است که در فصل اول به بیان کلیات و مفاهیم و همچنین تشریح مسئله خواهیم پرداخت، ضمناً در پایان هر فصل یک جمع‌بندی و خلاصه‌ای از مطالب بیان شده آن فصل خواهیم داشت. در فصل دوم علم پروتئومیکس و کاربردهای علوم کامپیوتری در آنالیز داده‌های حاصل از طیف‌سنج جرمی را بررسی خواهیم کرد در فصل سوم به داده‌کاوی و مرور ادبیات داده‌کاوی در پروتئوم جهت دسته‌بندی و کشف زیست‌نشانه‌های سرطانی خواهیم داشت. در فصل آخر خلاصه‌ای از فصل‌های گذشته را بیان خواهیم کرد و فرصت‌ها و چالش‌های داده‌کاوی پروتئوم جهت تشخیص سرطان را بیان کرده و سپس داده‌ها در دسترس و منابع دسترس‌پذیر را معرفی خواهیم کرد.

۱,۹ خلاصه فصل

انسان تقریباً متشکل از چهار ابر مولکول کربوهیدرات‌ها، پروتئین‌ها، اسیدهای نوکلئوتید و لیپیدها شده است که در این میان پروتئین‌ها به دلیل فراوانی، رفتار پویا، تغییرات پس از ترجمه و مطابق نبودن میزان بیان آن‌ها با mRNA به دیگر ابرمولکول‌ها برای شناخت بیماری‌های پیچیده مانند سرطان برتری دارد. سرطان بیماری است که از تقسیم غیرعادی سلولی ایجاد می‌شود که به دلیل نبود پروتئین خاص، فعالیت زیاد یک پروتئین یا کم‌فعالیتی یک پروتئین ایجاد می‌شود. سرطان دارای انواع مختلفی از سارکوما، خون، لنفاوی و ... که سالانه جان میلیون‌ها انسان را می‌گیرد و مسئله سرطان می‌تواند موجب ضعیف شدن کشورها و حتی ورشکستگی آن‌ها نیز شود. هرگونه تغییر بالینی در بدن ما باعث تغییر در میزان، نوع و فعالیت‌های پروتئین‌های موجود در بافت و یا ارگان‌های ما خواهد داشت که نتیجتاً سرطان نیز این تغییر را در الگوی پروتئین‌های بدن ما می‌گذارد. به پروتئین‌های بیان‌شده در بافت یا ارگان ما در یک لحظه خاص پروتئوم گفته می‌شود و علم بررسی پروتئوم را پروتئومیکس می‌گویند. داده‌های حاصل از پروتئوم بسیار پیچیده بوده و از نظر ویژگی بسیار زیاد است بنابراین برای کشف تغییرات الگوی پروتئین‌ها که منجر به شناسایی زیست‌نشانه‌هایی مانند انواع پروتئین که در بیماری‌های مثل سرطان نقش کلیدی دارند و همچنین برای دسته‌بندی بیماری‌های همچون سرطان نیاز به تحلیل‌های سطح بالایی همچون داده‌کاوی و علوم هوش مصنوعی احساس می‌شود. با کمک نتایج حاصل از این تحلیل‌ها می‌توان شناخت بهتری نسبت به سرطان پیدا کرد و در نتیجه داروها و روش‌های درمانی مناسبی را جهت مقابله با این بیماری ساخت.

فصل دوم

پروتئومیکس و کاربردهای علوم کامپیوتری در داده

های حاصل از طیف سنج جرمی

۲,۱ مقدمه

در فصل قبل تعریفی از پروتئوم و علم پروتئومیکس ارائه کردیم. هدف از این فصل پاسخ به پرسش‌هایی زیر است:

- روش‌های پروتئومیکس برای تفسیر و توضیح بیماری‌ها چیست؟
- کدام نمونه مناسب است؟
- پروتئین‌های در نمونه چگونه تفکیک، کمیت شماری و توصیف می‌شوند؟
- و رویکردهای برخورد با داده‌های حاصل از طیف‌سنج جرمی چیست؟

۲,۲ انواع پروتئومیکس

تحقیقات پروتئومیکس شامل نگرش‌های متنوعی به مسائل می‌شود. طیف این نگرش از اندازه‌گیری فراوانی و مقایسه پروتئین در دو نمونه آزمایش و کنترلی تا بررسی برهم‌کنش پروتئین در یک شبکه است. در این بخش به معرفی این رویکردها و هدف آن‌ها می‌پردازیم. لازم است که مجدداً تأکید کنیم در این پژوهش تمرکز بر روی روش‌های پروتئومیکس شناختی و مقایسه‌ای با استفاده از طیف‌سنج جرمی که همان روش توصیف پروتئین‌هاست می‌پردازیم.

۲,۲,۱ پروتئومیکس شناختی و مقایسه‌ای

هدف از پروتئین کمی اندازه‌گیری مطلق و یا مقایسه‌ای فراوانی پروتئین‌ها است. مزیت رویکرد مقایسه‌ای استفاده در مطالعات کمی برای اندازه‌گیری تغییرات پروتئین‌های بیان‌شده در یک بافت یا سلول است. رویکرد متداول مقایسه دودسته بیمار و سالم، یا سلول‌های تحت معالجه با یک دارو و سلول‌های دیگر و یا موارد دیگر که به دودسته کنترل و غیر کنترل تقسیم می‌شوند. پروتئین‌های با کاهش یا افزایش فراوانی در سلول شناسایی و کمی شماری شده تا تغییرات در بیان این پروتئین به من منظور کنترل و پایش بیماری اندازه‌گیری شود. این مورد بسیار پراهمیت است که بیماری‌هایی از قبیل سرطان، دیابت و بیماری‌های قلبی دارای عوامل متعدد ژنتیکی بوده و مکانیسم‌های و کنش و واکنش‌های پیچیده‌ای را دارا می‌باشند. (Elo and Schwikowski, 2012) شکل ۱-۲ شمایی از جریان کاری فرآیند شناسایی زیست‌نشانگرها را نمایش می‌دهد.

۲,۲,۲ تغییرات پس از ترجمه پروتئین‌ها

پروتئین‌ها در پاسخ به پیام‌های داخل سلول و خارج از سلول دچار تغییرات پس از ترجمه می‌شوند. این تغییرات پس از ترجمه یا اختصاراً PTM تنظیم‌کننده مکانیزم‌های فرآیندهای سلولی است و شناسایی این تغییرات پیش‌نیاز درک عملکرد پروتئین‌ها در سلول است. تعدادی مختلفی از PTM شناسایی شده است که شامل اکستیل‌اسیون، متیل‌اسیون، فسفوریلاسیون و گلیکوزیله شدن است. در میان این‌ها گسترده‌ترین مطالعات بر فسفوریلاسیون صورت گرفته است به این دلیل که تقریباً در تمامی فرآیندهای سلولی این تغییرات پس از ترجمه اتفاق می‌افتد. بنابراین تغییرات پس از ترجمه در پروتئین‌ها می‌تواند به عنوان عاملی برای بررسی و شناسایی بیماری مورد استفاده قرار گیرد. هدف شناسایی پروتئین‌هایی که دچار تغییرات پس از ترجمه شده، شناسایی مکان این تغییرات در سلول، اندازه‌گیری فراوانی آن‌ها و تعیین میزان ارتباط چندین PTM با یکدیگر است. تعداد نرم‌افزارها و ابزارها برای استفاده به طور فزاینده‌ای در چند سال گذشته افزایش پیدا کرده است. به طور مثال PTMScout (<http://ptmscout.mit.edu>) وب‌سایتی برای مشاهده، بررسی و آنالیز داده‌های پروتئین حاصل از PTM که شامل بانک اطلاعاتی از آزمایش‌ها PTM تا با یکدیگر با ابزار مقایسه‌ای مجموعه داده است. به عنوان مثالی دیگر ابزار scan-x است که هدف آن پیش‌بینی مکان‌های فسفوریلاسیون با استفاده از طیف وسیعی از داده‌های عمومی فسفو پروتئومیکس است. مرجع (Gallego and Virshup, 2007) نگاهی مروری به این رویکرد داشته است.

۲,۲,۳ مکان‌یابی پروتئین‌ها^۱

پروتئین‌ها باید در مکان‌های به خصوصی از سلول حضور داشته باشند تا بتواند عملکرد مناسب داشته باشند و عدم قرارگیری در جای نامناسب می‌تواند موجب عملکرد نادرست سلول شود. همچنان که پروتئین‌هایی که موجب یک بیماری می‌شوند تمایل دارند تا در زیر سلول‌های همان بخش ظاهر شوند. با این تفاسیر نقش برداری از مکان پروتئین‌ها در زیر سلول‌ها می‌تواند منجر به شناسایی زیست‌نشانه‌هایی نوینی برای شناسایی بیماری‌ها و ساخت و توسعه داروهای جدید باشد. برای مکان‌یابی از ترکیب روش‌های کامپیوتری و آنالیزهای پروتئین مانند طیف‌سنج جرمی استفاده می‌شود. (Elo and Schwikowski, 2012)

۲,۲,۴ برهم‌کنش پروتئین‌ها^۲

پروتئین‌ها به ندرت به تنهایی فعالیتی می‌کنند و معمولاً در شبکه از ارتباطات پیوسته در فرآیندهای سلولی شرکت می‌کنند و هرگونه تغییر و اختلال در این ارتباطات پروتئین که اختصاراً PPI^۳ گفته می‌شود می‌تواند زمینه‌ساز بیماری‌ها همچون سرطان شود. با شناخت این تعاملات پویای پروتئینی تشخیص بیماری‌ها و ساخت داروهای درمانی تسهیل می‌شود از این رو کاربردهای شبکه‌های پروتئینی در زمینه شناخت عوامل ژنتیکی بیماری‌ها، شناخت ویژگی‌های شبکه‌های پروتئینی و ارتباط زیر شبکه‌های پروتئینی به بیماری‌ها رو به گسترده شدن است. همچنین از شبکه‌های پروتئین‌ها برای دسته‌بندی بیماری‌ها و به خصوص سرطان استفاده می‌شود. مرجع (Wang et al., 2015) مروری بر شناسایی نشانگرهای سرطانی با رویکرد شبکه در تعاملات زیست‌نشانه‌ها در سلول انجام داده است.

¹ Protein localization

² Protein interaction

³ Protein-protein interaction



شکل ۲-۲۱ جریان کاری مراحل شناسایی زیست‌نشانگرها (Elo and Schwikowski, 2012)

۲.۳ روش‌های نمونه‌گیری برای کشف نشانگرهای سرطانی

طراحی مناسب آزمایش امری حیاتی برای موفقیت یک مطالعه است. با طراحی ضعیف یک آزمایش امکان تعیین این‌که آیا یک مشاهده به‌درستی تغییرات زیستی را نمایش می‌دهد یا این تغییرات تنها مربوط به مسائل فنی است، وجود ندارد. هزینه‌های بالا آزمایش‌ها پروتئومیکی معمولاً منجر به طراحی ضعیف آزمایش‌ها می‌شود و در نتیجه منجر به نتایج ضعیف یا اشتباه در مطالعه می‌شود. (Elo and Schwikowski, 2012)

زیست‌نشانگرها که مولکول‌های زیستی (ابر مولکول‌های ذکرشده در فصل اول) بوده و اطلاعات بسیاری را در رابطه با بیمارها دارا می‌باشند و می‌توانند برای تشخیص، پیش‌بینی، ارزیابی ریسک بیماری‌ها مورد استفاده قرار بگیرند، میزان و نوع آن‌ها در روش‌های نمونه‌گیری مختلف، متفاوت است در نتیجه تشخیص و انتخاب درست این نمونه‌ها در نتایج آزمایش بسیار تأثیرگذار است و نظر به این‌که این امر موردی حیاتی و اثرگذار در مطالعات پروتئومیکس است در این بخش نگاهی هرچند گذار اما تحلیلی و انتقادی به این روش‌ها خواهیم داشت.

۲,۳,۱ مایعات زیست پذیر^۴

سروم (مایعی است که پس از حذف پروتئین‌های عامل لختگی خون از پلاسما باقی می‌ماند) و پلاسما (بخش مایع خون که در آن سلول‌های خونی شناور است) رایج‌ترین نمونه‌های بیمارستانی برای شناخت و اثبات زیست‌نشانگرها است به دلیل دسترسی و شناخته‌شده بودن در مراکز درمانی است. مایعات زیست پذیر دیگر شامل اور^۵ CSF، بزاق، مایع آسیت و ... هم به سرعت در حال شناخته شدن به عنوان منبعی برای یافتن زیست‌نشانگرها است. بیشتر نمونه‌های پیش آمده بیشتر برای مطالعات ابتدایی کشف زیست‌نشانگرها مناسب است و در عمل برای استفاده نهایی در مراکز درمانی به دلیل دشواری در نمونه‌گیری و هزینه‌بر بودن قابل استفاده نیستند. در مقابل، اوره به دلیل سهولت در دستیابی و شیوه غیرتهاجمی بودن (عدم آسیب به بدن برای به دست آوردن نمونه) و سرعت در پردازش آن بسیار مناسب است. پروتئوم اوره منبع غنی از زیست‌نشانگرهای بیماری برای آنالیز است.

مایعات زیست پذیر مجموعه‌ای پیچیده از پروتئین‌ها است که از تومورها و بافت‌های سالم ناشی می‌شود؛ بنابراین، ممکن است موجب عدم دقت در بیان پروتئین‌های مقایسه‌ای بین دو گروه شود. مقادیر زیاد پروتئین‌ها در مایعات زیست پذیر (از قبیل آلبومین و هموگلوبین) ممکن است نقش پروتئین‌های با میزان کم اما دخیلی در سرطان را کمرنگ کند. برای غلبه به این مشکل نیاز به فن‌های متنوع آزمایشگاهی است.

۲,۳,۲ بافت^۶

مقایسه بیان پروتئین‌های شناخته‌شده بین بافت تومور و بافت سالم مجاور می‌تواند اطلاعات ارزشمندی را برای تشخیص بسیاری از بیماری به خصوص سرطان فراهم کند. چندین مزیت در استفاده از بافت برای پیدا کردن زیست‌نشانگرها وجود دارد که شامل این حقیقت می‌شود که زیست‌نشانگرهای کشف شده به طور واضح منشأ تومور می‌باشند. به هیچ عنوان نشانگرهای تومور در سلول‌های تومور پنهان نمی‌شوند؛ بنابراین، تمرکز بر روی بافت محلی تومور در قیاس با مایعات زیست پذیر بسیار بیشتر است و در نتیجه شناسایی نشانگرهای منتخب ساده‌تر است. اغلب با استفاده از فن‌های LCM و تجزیه و تحلیل پروتئوم می‌توان منبع دقیق نشانگرهای پروتئینی داخل تومور را کشف کرد. با این وجود مطالعات پروتئوم به کمک نمونه‌ها بافتی دارای ایراداتی نیز است. اول اینکه روشی تهاجمی است، دسترسی به برخی از تومورها و بافت‌ها محدود است. به علاوه، بافت تومور ترکیبی ناهمگون از سلول‌های بدخیم همچنان که از بافت‌های همبند، بافت چربی و سلول‌های التهابی است که نمایانگر چالش‌های فنی برای دریافت و آنالیز نمونه و کمبود دقت در شناسایی نشانگرهای بالقوه است.

۲,۳,۳ سل لاین^۷

برخلاف بافت، سل لاین‌های تومور می‌تواند بیان‌کننده متقارن جمعیت سلول‌ها باشد که می‌شود به سادگی در آزمایشگاه کشت کرد، می‌تواند در مقدار بسیار زیاد در دسترس باشد و می‌تواند برای طیف وسیعی از مطالعات کشف نشانگرها مورد استفاده قرار بگیرد. همچنین آماده‌سازی بخش‌های زیر سلولی شامل غشای پلاسما، سیتوزولی یا بخش هسته سلول، ترشحات، اگزوموز از طریق کشت برای مطالعات پروتئومیکسی نسبتاً ساده است. سلول‌های تومور ممکن است در محیط آزمایشگاهی به منظور شبیه‌سازی شرایط آسیب دیدگی ژن‌ها، بیان زیاد ژن‌ها و پاسخ آن‌ها به روش‌های درمانی قبل از مطالعات پروتئومیکسی استفاده شود. با این وجود

⁴ biofluids

⁵ Central nervous system

⁶ Tissue

⁷ Cell line

استفاده از سل‌لاینها دارای چندین محدودیت است از جمله کمبود دسترسی به سل‌لاینها مناسب تومورها و بافت‌های مخاطی به‌عنوان گروه کنترل در مطالعات، تغییرات در ژن یا بیان پروتئین‌های شناسایی‌شده که ممکن است در شرایط کشت دوبعدی اتفاق بیافتد و نبود تعامل بین سلول‌های تومور با سلول‌های استرومال و یا ایمنی که شرایط محیط تومور در بدن را مشخص می‌کند. (Yang et al., 2015, Panis, 2015)

۲.۴ فن‌آوری‌های آنالیز پروتئومیکس

قدم کلیدی در پروتئومیکس جداسازی و نمایش مجموعه‌ای از پروتئین‌های موجود در یک نمونه است. آنالیز پروتئومیکس شامل فن‌آوری‌های متنوعی است که هرکدام دارای مزایا و معایب می‌باشند. در این قسمت نگاهی به این روش‌ها خواهیم داشت.

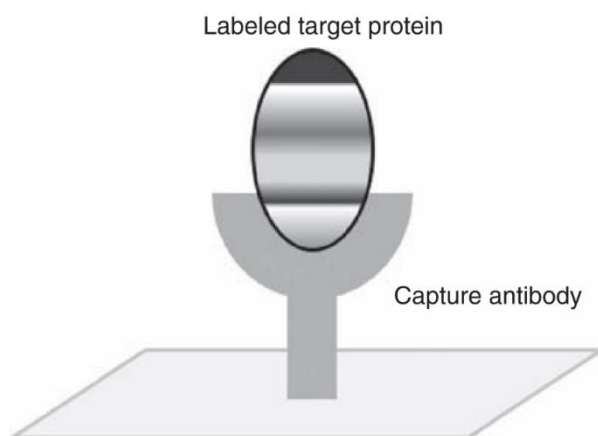
۲.۴.۱ فن‌آوری ریزآرایه‌های پروتئینی

ریزآرایه‌های پروتئینی می‌تواند به دودسته عمده: ریزآرایه‌های فاز مستقیم و ریزآرایه‌های فاز معکوس تقسیم شود. در ریزآرایه‌ها فاز مستقیم مولکول‌ها با یک الگوی ماتریسی در یک سطح کوچک بی‌حرکت قرار می‌گیرند. این مولکول‌های بی‌حرکت با آنالیت‌ها (نمونه‌ای مورد ارزیابی) مورد آزمایش جهت تعیین حضور یا فراوانی پروتئین‌ها موردنظر واکنش می‌دهند. این مولکول‌ها می‌تواند شامل آنتی‌بادی‌ها، پروتئین‌ها، تکه‌هایی از پروتئین‌ها، آنزیم‌ها و یا پپتیدها باشند که با توجه به هرکدام از آن‌ها به‌عنوان مولکول گیرنده استفاده شود فن‌ها و کاربردها در آنالیز متفاوت است. به‌طور مثال ریزآرایه‌های آنتی‌بادی شامل دو فن‌آوری ساندویچ ریزآرایه و فن‌آوری تک-آنتی‌بادی ریزآرایه است. در روش ساندویچ ریزآرایه، دو آنتی‌بادی تعیین‌کننده پروتئین برای پروتئین هدف مورد استفاده قرار می‌گیرد. آنتی‌بادی گیرنده بر روی سطح برای گیر انداختن پروتئین هدف در روی سطح بی‌حرکت قرار می‌گیرد. سپس آنتی‌بادی دوم که با فلورسنت برچسب‌گذاری شده به پروتئین هدف که توسط آنتی‌بادی اول که پروتئین را گیر انداخته است می‌چسبد. (شکل ۲-۲)

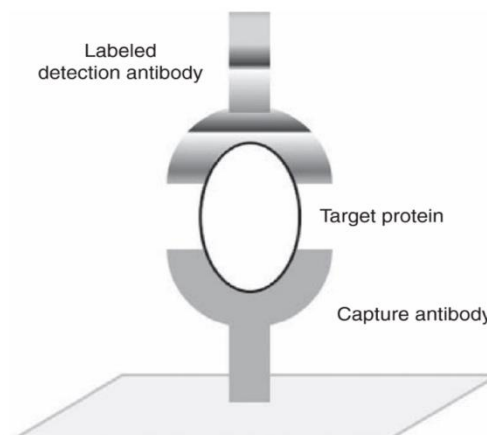
در روش آرایه‌های تک-آنتی‌بادی، نمونه مورد آزمایش جهت شناسایی مستقیم پروتئین‌ها توسط آنتی‌بادی‌های بی‌حرکت بر روی سطح برچسب‌گذاری می‌شوند. (شکل ۲-۳)

در ریزآرایه‌های فاز معکوس، آنالیت‌ها در یک سطح بی‌حرکت باقی می‌مانند. این مورد امکان ارزیابی هم‌زمان مقادیر کوچکی از بافت از نمونه‌های زیاد و متنوع را می‌دهد. بسته به طراحی ریزآرایه‌های پروتئین، چه آنالیت‌ها چه مولکول‌های گیرنده، به‌منظور شناسایی پروتئین‌ها برچسب‌گذاری می‌شوند. از روش‌های رایج برچسب‌گذاری شامل فلورسنت، نورتایی شیمیایی^۸ و رادیواکتیویته هست که روش فلورسنت پرتعدادترین آن‌ها است. (Dziuda, 2010)

⁸ chemiluminescence



شکل ۳-۲ روش آرایه‌های تک‌آنتی‌بادی (Dziuda, 2010)



شکل ۲۲-۲ روش ساندویچ ریز آرایه (Dziuda, 2010)

۲،۴،۲ ژل الکتروفورز دوبعدی پلی آکریل آمید (۲D-PAGE)

ژل الکتروفورز دوبعدی پروتئین‌ها را در بعد اول بر اساس بار آن‌ها و در بعد دوم بر اساس جرم آن‌ها جدا می‌کند. 2D-PAGE روش کمی پرتوان است که اجازه می‌دهد بیش از ۱۰ هزار پروتئین را با وضوح بالا تفکیک کرد. یکی از ایرادات این روش آن است که پروتئین‌های با مقادیر زیاد موجب نامفهوم سازی و پنهان ماندن پروتئین‌های با مقادیر کم شود. روش two-2D-DIGE (dimensional in-gel electrophoresis) یکی از انواع روش ژل الکتروفورز دوبعدی پلی آکریل آمید است. در این روش پروتئین‌ها دو نمونه متفاوت با کمک رنگ فلورسنت برچسب‌گذاری می‌شوند و سپس الکتروفورز می‌شوند. (Dubitzky et al., 2007)

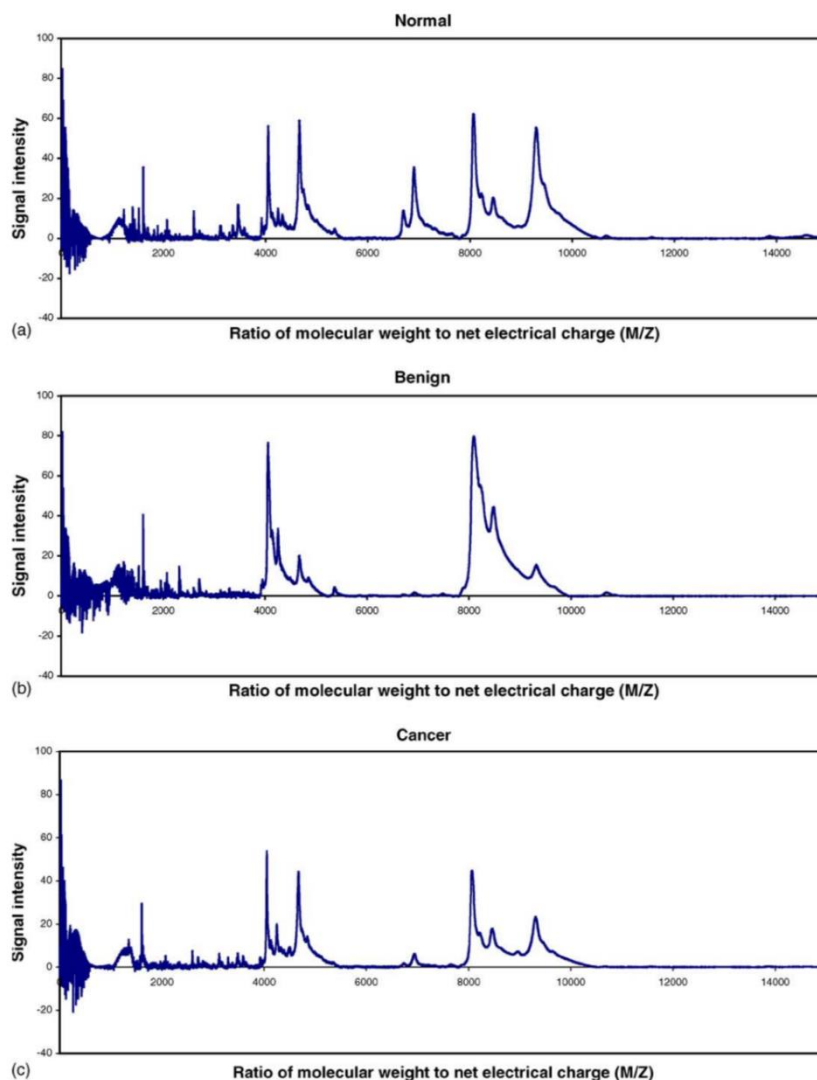
۲،۴،۳ طیف‌سنج جرمی^۹

طیف‌سنج جرمی یا به اختصار MS نقشی اساسی در شناسایی پروتئین‌ها و فرآیند پس از ترجمه آن‌ها بازی می‌کند. طیف‌سنج جرمی از سه بخش اساسی تشکیل شده است: (۱) منبع یونی که پروتئین‌ها را به گازهای یونی تبدیل می‌کند. (۲) ارزیاب جرمی که نسبت جرم به بار (m/z) یون‌ها را اندازه‌گیری می‌کند. (۳) و یک آشکارساز که تعداد یون‌های مشاهده‌شده برای یک مقدار مشخص m/z را شمارش می‌نماید. به‌طور عمده دو منبع یون‌ساز electrospray ionization (ESI) و matrix-assisted laser desorption/ionization (MALDI) موجود هستند. همچنان time-of-flight (TOF) به‌عنوان ابزار آنالیز در MALDI به‌کاربرده می‌شود. به‌طور مختصر نمونه‌های پروتئینی با مولکول‌های با چگشت ماتریسی مخلوط شده و سپس به لکه‌های بلورهای شده تبدیل می‌شوند و بر روی صفحات فلزی قرار می‌گیرند. سپس تابش‌های پرتاب‌شده از لیزر موجب پراکنده شدن و یونیزاسیون مخلوط می‌شود. پروتئین‌های یونیزه از اتاق یونی عبور می‌کنند و به آشکارساز برخورد می‌کنند. بنا بر ولتاژ بکار رفته و شتاب یون‌ها، میزان بار به جرم (m/z) هر کدام از یون‌ها را می‌توان اندازه گرفت و سپس نمایش داد. Surface-enhanced laser desorption time-of-flight (SELDI-TOF) گونه جدیدی از MALDI-TOF است. عنصر کلیدی در SELDI-TOF MS ترارش‌های پروتئینی با سطح شیمیائی برای به دام انداختن دسته‌ای از پروتئین‌ها در شرایط خاص است. MALDI MS و SELDI-

⁹ Mass Spectrometry

TOF MS فناوری‌های با حساسیت بالا هستند و برای شناسایی و تشخیص پروتئین‌ها در حجم بالا بسیار مناسب می‌باشند. (Dubitzky et al., 2007)

در پروتئومیکس رایج‌ترین روش برای شناسایی پروتئین‌ها با به‌کارگیری MS، رویکرد پائین به بالا است. در این رویکرد مولکول‌های مورد ارزیابی پپتیدها هستند که با فرآیند تجزیه آنزیماتیک پروتئین‌ها نمونه حاصل شده‌اند که مراحل آن را شرح دادیم. طیف ایجادشده از تکه‌های پپتیدها به tandem MS spectra شناخته می‌شود که نتایج حاصل از آن بصری سازی می‌شود که موجک‌های موجود در تصویر حاصل نشان‌دهنده آمینواسیدهای حاضر در پپتیدها است؛ اما در این رویکرد نتایج حاصل نمایش‌دهنده پپتیدها است و نیاز به مرحله دیگری برای شناسایی پپتیدها برای پیش‌بینی پروتئین‌ها موجود در نمونه است. این فرآیند می‌تواند با استفاده از نرم‌افزارهای جستجوی توالی در پایگاه داده‌ها مانند mascot استفاده شود که این روش در بخش‌های بعدی بررسی خواهد شد؛ اما رویکرد بعدی، از بالا به پائین است که MS مستقیماً برای ارزیابی پروتئین‌هایی سالم و تجزیه نشده استفاده می‌شود؛ اما این رویکرد نیاز به تجهیزات بیشتر و پیچیده‌تر و در نتیجه پرهزینه‌تر برای بکارگیری است که همین عوامل باعث عدم استقبال از این رویکرد است. (Swan et al., 2013) شکل ۲-۴ نمونه‌ای از طیف جرمی بصری سازی شده را نمایش می‌دهد.



شکل ۲-۴ طیف جرمی حاصل از سرم خون با فناوری SELDI

این طیف نمایش دهنده سه دسته سرطان خوشخیم، سرطان بدخیم و سالم را نشان می‌دهد. (Li et al., 2004)

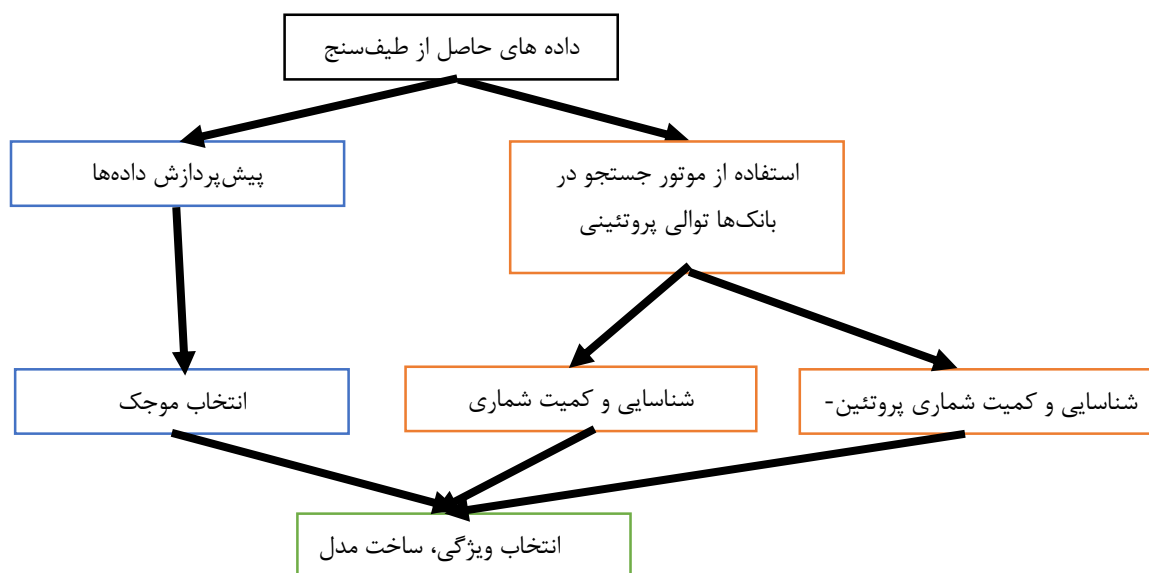
۲.۵ فرصت‌ها و چالش‌های آنالیز داده‌ها در طیف‌سنج جرمی

استفاده از طیف‌سنج جرمی فرصت و روزه‌های پیشرفت بسیاری را گشوده است ولی با این وجود دچار کاستی‌هایی نیز است. محور اساسی در شناسایی پروتئین‌ها با استفاده از طیف‌سنج جرمی، کشف زیست‌نشانگرهای بیماری‌ها است. این زیست‌نشانگرها می‌توانند برای اهداف مختلفی مانند شناسایی فرآیندهای بیولوژیکی، شناسایی و پیش‌بینی بیماری‌ها و میزان پیشرفت بیماری‌های شناخته‌شده در بدن استفاده شود. همچنان که برای بررسی تأثیرگذاری یک داروی خاص و توسعه و ساخت داروهای جدید بکار

می‌رود. فرصت دیگر حاصل از طیف‌سنج جرمی، نگرش به پروتئین‌ها از دیدگاه شبکه و مسیرهای ارتباطی آن‌ها است؛ اما با تمام این فرصت‌ها دارای محدودیت‌هایی نیز است. به‌طور مثال برای نمونه‌گیری‌های متعدد به دلیل محدودیت‌های زمانی و همچنین هزینه زیاد مناسب نیست. همچنین پیدا کردن پروتئین‌هایی با میزان کم در نمونه که در فرایند ایجاد بیماری مشارکت دارند دشوار است. مورد بعدی هم عدم قطعیت در شناسایی درست پروتئین‌ها بنا به توالی پپتیدی آن‌ها در بانک‌های اطلاعاتی است.

۲.۶ آنالیز داده‌های طیف‌سنج جرمی

داده‌های حاصل از طیف‌سنج جرمی مجموعه‌ای مोजک‌ها برحسب بار به جرم است که برای شناسایی پپتیدها است. ضرورت استفاده از روش‌های علوم کامپیوتری برای شناسایی پروتئین‌ها از مोजک‌های شناسایی‌شده و مقایسه نمونه‌ها امری اجتناب‌ناپذیر است. بدین منظور دو رویکرد متفاوت جهت برخورد با مोजک‌های حاصل از طیف‌سنج جرمی وجود دارد. رویکرد اول مستقیماً بر روی مोजک‌ها پردازش و عملیات داده‌کاوی را انجام می‌دهد که به آن انتخاب مोजک^{۱۰} گفته می‌شود و در رویکرد دوم استفاده از موتورهای جستجو در بانک‌های پروتئینی و توالی یابی شده جهت شناسایی و توصیف پروتئین‌ها است که این رویکرد نیز در پیوستگی استفاده از برچسب‌گذاری در حین آماده‌سازی نمونه‌ها برای آنالیز که تقریباً روشی آزمایشگاهی است استفاده می‌شود و یا از روش‌ها بدون برچسب‌گذاری و مبتنی بر روش‌های کامپیوتری جهت کمی شماری پروتئین‌ها استفاده می‌کند. شکل ۲-۵ جریان آنالیز داده‌های طیف‌سنج جرمی را به‌طور خلاصه بیان می‌کند.



شکل ۲-۵ جریان آنالیز داده‌های حاصل از طیف‌سنج جرمی

¹⁰ Peak picking

۲.۶.۱ انتخاب موجک

در انتخاب موجک داده‌های حاصل از طیف‌سنج جرمی فارغ از این که کدام پپتیدها و کدام پروتئین‌ها در نمونه حاضر هستند مورد پردازش قرار می‌گیرند در عوض موجک‌هایی با بیشترین سیگنال به‌عنوان زیست‌نشانگر کاندید می‌شوند. کاتاجاما و اریزیک دو ایراد این روش را برشمردند اول این که یک رویکرد مستقیم برای شناسایی پروتئین ارائه نمی‌کند و برای این منظور تحلیل‌های بیشتر نیاز است. دومین ایراد این است که پیش‌پردازش داده‌ها شامل نرمال‌سازی، همگام‌سازی موجک‌ها و کاهش نویز امری ضروری و اجتناب‌ناپذیر است. (Katajamaa and Orešič, 2005) درواقع بدون پیش‌پردازش داده‌ها امکان مقایسه موجک‌های کاندید بین نمونه‌ها وجود ندارد و دقت مدل‌های ارائه‌شده به‌شدت کاهش می‌یابد. در جدول ۱-۲ نمونه‌ای از داده‌های حاصل از طیف‌سنج جرمی که برای پردازش توسط روش‌های انتخاب موجک فراهم‌شده‌اند نمایش داده‌شده است. به این جدول ماتریس شدت M/Z نیز گفته می‌شود.

۲.۶.۲ موتورهای جستجو

در طی فرآیند طیف‌سنج جرمی، جرم پپتیدهای شناسایی‌شده حاضر در نمونه مورد تحلیل قرار می‌گیرد. جرم این پپتیدها هم‌راستا با جرم تکه‌هایشان جهت این که کدام پپتیدها حضور دارند و به کدام پروتئین‌ها متعلق می‌باشند مورد بررسی قرار می‌گیرد.

نرم‌افزارهای جستجو از قبیل Mascot برای بررسی احتمال حضور پروتئین‌ها توسعه داده‌شده است. این نرم‌افزار در ارتباط بانک توالی پروتئینی مانند UniportKB و NCBItr^{۱۱} و همچنان از بانک داده‌های خاص مانند sgn متعلق به TOMATO^{۱۲} و Tair متعلق به Arabidopsis^{۱۳} استفاده می‌کند. پپتیدهای حاضر در نمونه آنالیز شده با طیف‌سنج جرمی بررسی شده و با قیاس تناسب جرم به حجم پپتیدهای نمونه با پپتیدهای پروتئین‌های شناخته‌شده در بانک داده‌ای موردنظر شناسایی و پیش‌بینی می‌شود. نیلسون و همکاران معتقدند که این روش به‌طور کامل دقیق نیست به خاطر شباهت توالی برخی از پروتئین‌ها و نسبت کم پروتئین‌های توالی شده به کل پروتئین‌های موجود در بدن انسان که موجب عدم قطعیت و ضعف این روش می‌شود. (Neilson et al., 2011)

همان‌طور که ذکر کردیم شناسایی پروتئین با کمک موتورهای جستجو در پیوستگی استفاده از روش‌های برچسب‌گذاری و بدون برچسب‌گذاری است. روش‌های و رویکردهای برچسب‌گذاری متنوعی وجود دارد بنا به این که جز روش‌های آزمایشگاهی برای کمیت‌شماری پروتئین‌ها بوده و از موضوع بحث ما خارج است از بیان آن‌ها صرف‌نظر می‌کنیم؛ اما ضعف این روش شامل هزینه‌های بالای آزمایشگاهی، محدودیت در تعداد نمونه‌های که می‌توان آنالیز کرد و همچنین عدم سازگاری با برخی از نمونه‌های که برای تحلیل گرفته می‌شود.

¹¹ <http://www.ncbi.nlm.nih.gov/>

¹² <http://solgenomics.net>

¹³ <http://www.arabidopsis.org>

اما روش دیگر تحلیل داده‌های بدون استفاده از برچسب‌گذاری و روش‌های اضافی آزمایشگاهی است. دو روش بدون برچسب‌گذاری برای کمیت‌شماری پروتئین‌ها وجود دارد (۱) شدت سنجی سیگنال (۲) طیف‌شماری.

در شدت سنجی سیگنال از ناحیه زیر سطح منحنی (AUC) طیف‌ها موحک‌ها برای مقایسه مقدار پپتیدهای حاضر در نمونه استفاده می‌شود. در روش دوم طیف مشاهده‌شده برای پپتیدهای یک پروتئین جمع زده می‌شود. مرجع (Wong and Cagney, 2010) مروری به این رویکردها داشته است. تعداد زیادی نرم‌افزار منبع باز و تجاری بدون برچسب‌شمارشگر موجود است. به‌عنوان مثالی از نرم‌افزارهای رایگان برای روش AUC می‌توان از MSInspect و MSQuant نام برد و برای روش طیف‌شماری می‌توان از نرم‌افزارهای PepC و APEX نام برد که می‌تواند برای کمیت‌شماری روش‌های بدون برچسب‌گذاری استفاده شود سپس اعداد و کمیت‌های حاصل با استفاده از روش‌های کامپیوتری جهت شناسایی پروتئین‌ها مورد تحلیل قرار می‌گیرند. روش‌های کامپیوتری به دلیل عدم استفاده از روش‌های پیچیده آزمایشگاهی به‌سادگی قابل‌استفاده می‌باشند. در نتیجه انعطاف بیشتری در نمونه‌گیری و آماده‌سازی نمونه‌گیری وجود خواهد داشت. رویکردهای کمی مقادیر عددی با توجه به پروتئین شناسایی‌شده از طیف‌سنج جرمی تولید می‌کنند که این می‌تواند یک مزیت باشد برای زمانی که می‌خواهیم از روش‌های داده‌کاوی برای دسته‌بندی بیماری‌های همچون سرطان و یا شناسایی زیست‌نشانه‌های موثر در بیماری مثل سرطان با قیاس دودسته از نمونه‌های سالم و بیمار حاصل شود، استفاده کنیم. (Swan et al., 2013) جدولی از نرم‌افزارهای مورد استفاده در پروتئومیکس در پیوست ((الف)) آمده است. جدول ۲-۲ نمونه‌ای از پروتئین‌های شناسایی‌شده را در یک ماتریس نشان می‌دهد که به آن ماتریس بیان پروتئین^{۱۴} نیز گفته می‌شود.

جدول ۲-۱ ماتریس شدت M/Z (Dziuda, 2010)

در اینجا ستون اول بیانگر جرم به بار یک پپتید خاص و ردیف‌های ستون اول بیانگر نمونه و سایر مقادیر بیانگر شدت سیگنال یک پپتید در یک نمونه است.

m/z	Class 1					Class 2		***	Class J		
	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	***	***	***	***	Sample N-1	Sample N
500.0223	25.9588	26.6526	24.4485	29.7677	26.1579					24.4203	25.9005
500.4307	28.6677	29.4781	26.9261	33.0075	29.8984					27.7164	29.0940
500.8573	33.8643	35.1272	31.9417	38.9422	34.3449					31.4630	33.0437
501.2751	38.0995	40.2232	35.2459	43.2378	36.7521					33.3736	35.1533
501.6931	39.6762	42.3087	35.7506	44.4585	38.0613					34.1921	36.1327
502.1113	39.7707	42.6244	35.3323	44.2172	38.0231					34.1357	36.2976

15985.13	5.1260	4.9562	5.1565	4.4309	5.6946					4.5777	4.5680
15988.53	5.1059	4.9439	5.1494	4.4263	5.6690					4.5746	4.5636
15992.61	5.0613	4.9013	5.1072	4.4050	5.5761					4.5602	4.5554
15995.77	5.0525	4.8931	5.0944	4.4004	5.5560					4.5534	4.5514
15998.93	5.0297	4.8816	5.0752	4.3936	5.5237					4.5499	4.5437

¹⁴ Expression protein matrix

جدول ۲-۲ ماتریس بیان پروتئین

ستون اول بیان گر نام پروتئین‌ها، در سایر ستون‌ها ردیف اول نمایش گر نام نمونه و مقادیر جدول بیان پروتئین‌ها بر حسب M/Z												
	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	***	***	***	***	***	Sample N – 1	Sample N
Variable 1	50.6083	45.8562	56.6637	66.0191	48.7441						20.3157	22.0191
Variable 2	81.3635	80.8688	86.9904	99.9598	80.8768						35.7969	48.9452
Variable 3	30.7451	21.9715	28.7226	25.6417	25.2696						16.3164	15.8974
Variable 4	25.2859	24.6415	32.1971	30.4539	23.9356						32.3733	65.2203
Variable 5	75.9739	99.6320	95.6099	73.0236	50.5500						86.1237	94.3348
Variable 6	31.3223	61.3896	40.8125	30.7833	25.0216						43.7728	47.9578

Variable p-3	26.2785	36.7642	30.6385	30.8155	24.0845						40.8044	43.6642
Variable p-2	17.9427	19.6682	20.0322	20.0845	16.7964						34.6466	40.4123
Variable p-1	19.5918	20.5732	22.3328	22.2665	18.7692						19.7486	25.0679
Variable p	75.2334	90.1458	89.2247	75.6991	75.5829						87.5451	74.6747

۲.۷ خلاصه فصل

یکی از روش‌های پروتئومیکس شامل شناسایی، توصیف و کمیت شماری پروتئین‌ها است که گاهی نیز به پروتئومیکس قیاسی نیز شناخته می‌شود. روش دیگر بررسی تغییرات پس از ترجمه پروتئین‌ها در سلول است. روش دیگر مکان‌یابی پروتئین‌ها است که به بررسی مکان پروتئین‌ها در سلول می‌پردازد و در آخر برهم‌کنش پروتئین‌ها که به تعاملات پروتئین‌ها با یکدیگر در یک شبکه نگاه می‌کند. برای جمع‌آوری داده می‌توان از مایعات زیست پذیر، بافت‌ها و سل‌لاینها استفاده کرد که مایعات زیست پذیر با توجه به گردش آن‌ها در بدن و به تبع آن دارا بودن اطلاعات بسیار از بدن و راحتی جمع‌آوری نمونه روش مناسب و پیشنهادی بسیاری از تحقیقات است. برای آنالیز این نمونه‌ها از فناوری‌های تراشه‌های پروتئینی، ژل الکتروفورز دوبعدی پلی آکریل آمید (2D-PAGE) و طیف‌سنج جرمی می‌توان استفاده کرد که هرکدام از این فناوری‌ها دارای انواع مختلف و زیرمجموعه‌های خاص خود هستند. در این میان طیف‌سنج جرمی به دلیل توان بالایی در آنالیز انواع نمونه‌ها و توان بالا در آنالیز حجم بالایی از پروتئین‌ها و غیر جانب‌دارانه بودن یعنی در ابتدای آزمایش از پروتئین‌های بیان شده خبری نداریم به‌عنوان فناوری مرسوم توسط کلینیک‌ها و محققین استفاده می‌شود. در برخورد با داده‌های حاصل از طیف‌سنج جرمی دو رویکرد وجود دارد. این‌که مستقیماً روی موجک‌های حاصل پردازش صورت گیرد و یا این‌که ابتدا به شناسایی پروتئین‌ها اقدام کرده و سپس به پردازش آن‌ها به پردازیم که این ایراد رویکرد اول برای کشف زیست‌نشانگرها نیاز به تحلیل‌های اضافی است و همچنین به‌کارگیری پیش‌پردازش داده‌ها امری اجتناب‌ناپذیر است.

فصل سوم

داده‌کاوی و کاربرد آن در دسته‌بندی و کشف

زیست‌نشانگرهای سرطان

۳,۱ مقدمه

فرآیند داده‌کاوی پروتئوم شامل پیش‌پردازش داده‌ها (بیشتر برای زمانی که مستقیم بر روی مویک‌ها پردازش و مدل‌سازی می‌کنیم)، انتخاب ویژگی و کاهش بعد (به دو منظور افزایش دقت دسته‌بندی و کشف زیست‌نشانه‌ها که موجب تمایز دودسته می‌شوند)، ساخت مدل دسته‌بندی و در آخر ارزیابی عملکرد مدل است. در این فصل شش مدل دسته‌بندی درخت تصمیم (DC)^۱، شبکه‌های عصبی مصنوعی (ANN)^۲، جنگل تصادفی (RF)^۳، مدل‌های مبتنی بر قاعده^۴، بردارهای ماشین پشتیبان (SVM)^۵ و بیز ساده (NB)^۶ را مورد بررسی قرار داده و در آخر مروری بر تحقیقات پیشین خواهیم داشت.

۳,۲ داده‌کاوی و کاربرد آن در پروتئومیکس

داده‌کاوی رویکردی است که در علم و زمینه‌های تجاری برای استخراج مفاهیم معنادار و قابل استفاده از مجموعه داده‌های بزرگ و پیچیده استفاده می‌شود. فرآیند داده‌کاوی معمولاً تکرارپذیر است و داده‌های ناشناخته و دارای اطلاعات بالقوه به کمک ابزار تحلیلی پیچیده مورد کاوش و کشف قرار می‌گیرد. فرآیند کشف شامل پیدا کردن روابط و الگوها در داده‌های خام است که بشود از آن برای تصمیم‌گیری و آنالیزهای بیشتر استفاده کرد. تکامل شیوه‌های داده‌کاوی از پیشرفت در هوش مصنوعی، آمار و مدیریت انبار داده‌ها حاصل شده است. (Thomas et al., 2006) داده‌کاوی می‌تواند شامل گونه‌های نظارتی، یا همان مدل‌های دسته‌بندی^۷ باشد که در این گونه داده‌ها هر کدام دارای برچسب هستند که دسته آن داده را مشخص می‌کند، داده‌ها به دو دسته داده‌های آموزشی و تست تقسیم می‌شوند، با کمک داده‌های آموزشی مدل خود را می‌سازیم و سپس با داده‌های تست مدل خود را ارزیابی می‌کنیم در مقابل گونه غیر نظارتی که همان مدل‌های خوشه‌بندی^۸ است، داده‌ها فاقد برچسب هستند و بر اساس شباهت‌ها با کمک مدل‌های

¹ Decision tree

² Artificial neural network

³ Random forest

⁴ Rule-based classifiers

⁵ Support vector machine

⁶ Nave bayes

⁷ classification

⁸ clustering

خوشه بند در یک گروه قرار می گیرند که به این گروه ها خوشه گفته می شود. (Swan et al., 2013) در این بخش ما بر روی گونه های نظارتی و به تبع مدل های دسته بند تمرکز خواهیم کرد. داده کاوی با نظارت در پروتئومیکس توصیفی و مقایسه ای معمولاً با دو هدف دنبال می شود (۱) دسته بندی سرطان تا (بدخیم و خوش خیم، سرطانی و سالم، سرطان نوع A و سرطان نوع B و ...) (۲) کشف زیست نشانگرهای سرطانی که از مقایسه دسته تا مثلاً دسته سالم و سرطان، سرطان خوش خیم و بدخیم و ... به دست می آید که لازمه حصول به هدف دوم انتخاب ویژگی و کاهش بعد است. اگر ما در فرآیند داده کاوی مستقیماً به سراغ موجک ها برویم، پس از انتخاب ویژگی و ساخت مدل نیاز است تا بر روی موجک های منتخب تحلیل اضافی جهت روشن سازی اینکه هر کدام از این موجک های منتخب مربوط به کدام پروتئین است، صورت گیرد. داده کاوی داده های پروتئینی به دلیل کم بودن نمونه تا (سطرها) و زیاد بود ویژگی تا (ستون) بسیار چالش برانگیز است و نیاز به ابزارهای تحلیلی سطح بالا داده کاوی است. شکل شمایی از جریان مراحل داده کاوی بر داده های پروتئومی به منظور دسته بندی سرطان را نمایش می دهد. شکل ۳-۱ جریان کاری داده کاوی بر روی داده های حاصل از طیف سنج جرمی را زمانی که داده ها به صورت ماتریس بیان پروتئین باشند نمایش می دهد.

۳.۲.۱ پیش پردازش داده ها

داده های حاصل از MS دارای نویز بوده و هرگونه ساخت مدل دسته بند قبل از پیش پردازش ممکن است نتایج گمراه کننده ای را به بار آورد. بیشتر تحقیقات منتشر شده از نرم افزارهای که برای پیش پردازش داده ها توسعه داده شده اند استفاده کرده اند. به طور مثال نرم افزار ciphergen که برای پیش پردازش داده های حاصل از SELDI-MS استفاده می شود. این نرم افزار مکان و شدت هر کدام از پروتئین های در نمونه را یافته و فعالیت های پیش پردازش شامل کاهش خط مبنا^۱، همگام سازی موجک تا نرمال سازی شدت موجک تا و شناسایی موجک تا را انجام می دهد. کاهش خط مبنا نویزهای شیمیایی و الکترونیکی را پاک می کند. معمولاً کاهش خط مبنا در دو مرحله انجام می شود. اول تخمین خط مبنا با روش های پارامتریک و غیر پارامتریک می شود و سپس داده های اصلی از پارامتر تخمین زده شده کسر می شود. زمانی که کاهش خط مبنا تکمیل شد نوبت بعدی پیش پردازش نرمال سازی است. به دلیل اینکه موجک های یک طیف مقدار فراوانی نسبی یک پروتئین را توصیف می کنند برای مقایسه معنی دار بین طیف های مختلف (منظور نمونه تا) نرمال سازی صورت می گیرد.

۳.۲.۲ نفرین بعد

مهم ترین چالش در داده کاوی پروتئوم مسئله ابعاد داده هاست بدین معنی که تعداد نمونه ها کم و تعداد ویژگی ها (موجک ها) بسیار زیاد است و این درست برعکس داده های تجاری مانند خرده فروشی ها که تعداد ویژگی ها (ستون تا) کم ولی تعداد نمونه ها (سطرها) بسیار زیاد است. یک رویکرد مستقیم می تواند این باشد که هر کدام از موجک تا را به عنوان یک ویژگی در نظر بگیریم اما متأسفانه تعداد ویژگی در داده های پروتئوم بسیار زیاد بوده و می تواند بین ۳۰۰۰-۱۵۰۰۰ ویژگی باشد. این مشکل در داده کاوی به عنوان نفرین بعد یا مشکل ابعاد مطرح است. ابعاد بالا معمولاً بدین معناست که تعداد زیادی از ویژگی های بی استفاده و بی اثر باعث پنهان ماندن ویژگی های کلیدی یک مجموعه داده می شوند. ابعاد بالا علاوه بر اشغال حافظه و کاهش سرعت پردازش می تواند دقت و درستی مدل دسته بند را هم تحت تأثیر قرار دهد. (Wang et al., 2017)

غربال طیف داده ها به منظور شناسایی موجک ها، به عنوان استخراج ویژگی شناخته می شود. بر اساس این فرآیند هر گروهی از نقاط M/Z که داخل یک گروه قرار می گیرد به وسیله میانگین یا ماکزیمم داده های آن گروه توصیف می شود. در نتیجه این مقادیر نماینده

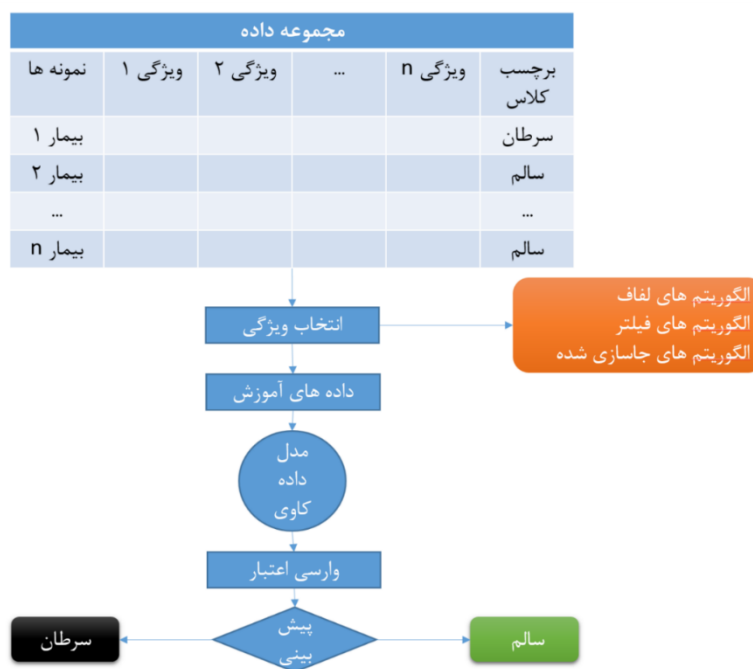
¹ Baseline reduction

به عنوان ویژگی (موجک) در نظر گرفته می شود. این گروه ها می توانند مستقل یا غیرمستقل، هم اندازه یا متغیر، پوشا یا غیر پوشا باشند. با تغییر اندازه گروه ها پژوهشگر می تواند فرآیند استخراج داده ها را بهینه کند.

اما در مقابل انتخاب و استفاده از ویژگی های به خصوص برای ساخت مدل موجب دقت و کارآمدی بیشتر دسته بندی می شود. فرآیند داده کاوی با کاهش ویژگی سرعت و سهولت بیشتر پیدا می کنند و نتایج تفسیرپذیر تر می شوند. همچنین برای کشف زیست نشانگرها کاندید که موجب تمایز دو دسته از هم می شوند استفاده از انتخاب ویژگی را اجتنابناپذیر می کند. پس پاک کردن ویژگی های غیر ضرور و بی استفاده امری ضروری است. البته باید به این نکته توجه داشت که کاهش ابعاد همیشه تضمین کننده نتایج موفقیت آمیز در انتخاب ویژگی نیست بنابراین لازم است تا متغیرهای انتخاب شده اعتبار سنجی شوند.

۳.۲.۳ روش های انتخاب ویژگی

به طور کلی روش های انتخاب ویژگی به سه دسته روش فیلتر^۲، روش لفاف^۳ و روش جاسازی شده^۴ طبقه بندی می شود. روش های فیلتر مستقل از الگوریتم دسته بندی بوده و به طور جداگانه داده ها رو مورد بررسی قرار می دهد. روش های لفاف مدل دسته بندی یادگیرنده بر مبنای زیرمجموعه ای از ویژگی های انتخاب شده مورد ارزیابی قرار می دهند. اگرچه این روش وابستگی بین متغیرها را بررسی می کند اما ریسک بیش برآزش را افزایش داده و همچنین از لحاظ محاسباتی بسیار سنگین است. روش های جاسازی شده زیر مجموعه ای بهینه از ویژگی های انتخاب شده را با ساخت انتخاب ویژگی داخل الگوریتم ارزیابی می کند. اگرچه این روش در قیاس با روش های فیلتر کمی محاسبات بیشتری را نیاز دارد اما در عوض در کنش و واکنش مستقیم با مدل دسته بندی می باشد. (Jagga and Gupta, 2015)



شکل ۱-۳ جریان کاری داده کاوی بر روی داده های طیف جرمی (Jagga and Gupta, 2015)

² Filter method

³ Wrapper method

⁴ Embedded method

۳,۲,۴ مدل‌های دسته‌بندی

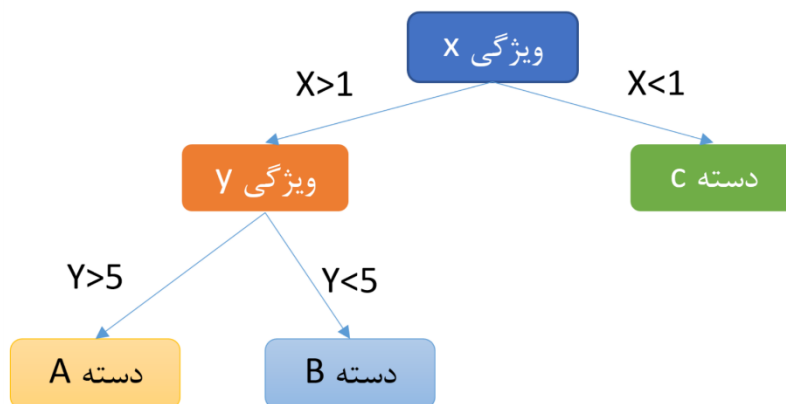
بیز ساده: بیز روشی برای دسته‌بندی پدیده‌ها، بر پایه احتمال وقوع یا عدم وقوع یک پدیده است و در نظریه احتمالات با اهمیت و کاربرد است. اگر برای فضای نمونه‌ای مفروضی بتوانیم چنان افرازی انتخاب کنیم که با دانستن این که کدام یک از پیشامدهای افراز شده رخ داده است، بخش مهمی از عدم اطمینان تقلیل میابد.

این قضیه از آن جهت مفید است که می‌توان از طریق آن احتمال یک پیشامد را با مشروط کردن نسبت به وقوع و یا عدم وقوع یک پیشامد دیگر محاسبه کرد. در بسیاری از حالت‌ها، محاسبه احتمال یک پیشامد به صورت مستقیم کاری دشوار است. با استفاده از این قضیه و مشروط کردن پیشامد موردنظر نسبت به پیشامد دیگر، می‌توان احتمال موردنظر را محاسبه کرد.

دسته‌بندی مبتنی بر قاعده: هم‌زمان با پیدایش علم داده‌کاوی در دهه ۹۰ الگوریتم‌های استخراج قوانین وابستگی از پایگاه داده‌ها نیز پا به عرصه گذاشت. نویسندگان زیادی در زمینه استخراج قوانین وابستگی در پایگاه داده‌ها بحث کرده‌اند رابرت. اس (۲۰۰۳) در مقاله خود اقدام به مقایسه الگوریتم‌های مهم استخراج قوانین وابستگی پرداخت است. در این مطالعه مزیت‌های و معایب سه الگوریتم مهم مورد استفاده در استخراج قوانین وابستگی یعنی *apriori*, *sampling* و *partitioning* پرداخته است.

اساساً ارتباط میان مجموعه اشیا (چیزها) وابستگی‌های جالب توجهی هستند که منجر به امکان آشکارسازی الگوهای مفید و قوانین وابستگی برای پشتیبانی تصمیم، پیش‌بینی‌های مالی، سیاست‌های بازاریابی، وقایع پزشکی و خیلی کاربردهای دیگر می‌شود. در حقیقت توجهات زیادی را در تحقیقات اخیر به خود جلب کرده است. تحلیل وابستگی‌ها یک حالت غیر نظارتی داده‌کاوی است که به جستجو برای یافتن ارتباط در مجموعه‌ها می‌پردازد. (تیمورپور، بابک؛ نجفی حیدر، ۱۳۹۴)

درخت تصمیم: درخت تصمیم یک مدل یادگیرنده است که در یک ساختار درخت مانند نمونه تا را تفکیک می‌کند. شکل یک درخت ساده تصمیم را نمایش می‌دهد که یک مجموعه داده را بر اساس مقادیر دو ویژگی به سه کلاس تقسیم کرده است. درخت‌های تصمیم ساده برای درک نحوه دسته‌بندی بسیار قابل فهم هستند. نمونه‌های جدید بر اساس میزان پیروی از هر یک از سه شاخه موجود بر اساس ویژگی‌هایشان دسته‌بندی می‌شوند. روش‌های از قبیل C4.5 با یک درخت خالی شروع کرده و مرتباً داده‌ها را تقسیم می‌کنند، شاخه‌های درخت را می‌سازند تا زمانی که تمام نمونه‌های یک شاخه به یک دسته خاص تعلق گیرد، برگ‌های درختان بر اساس معیارهای خاصی ساخته می‌شود. میزان خطا در شاخه‌های درختان به اندازه کافی کم است. (Han et al., 2011) (شکل ۳-۲)



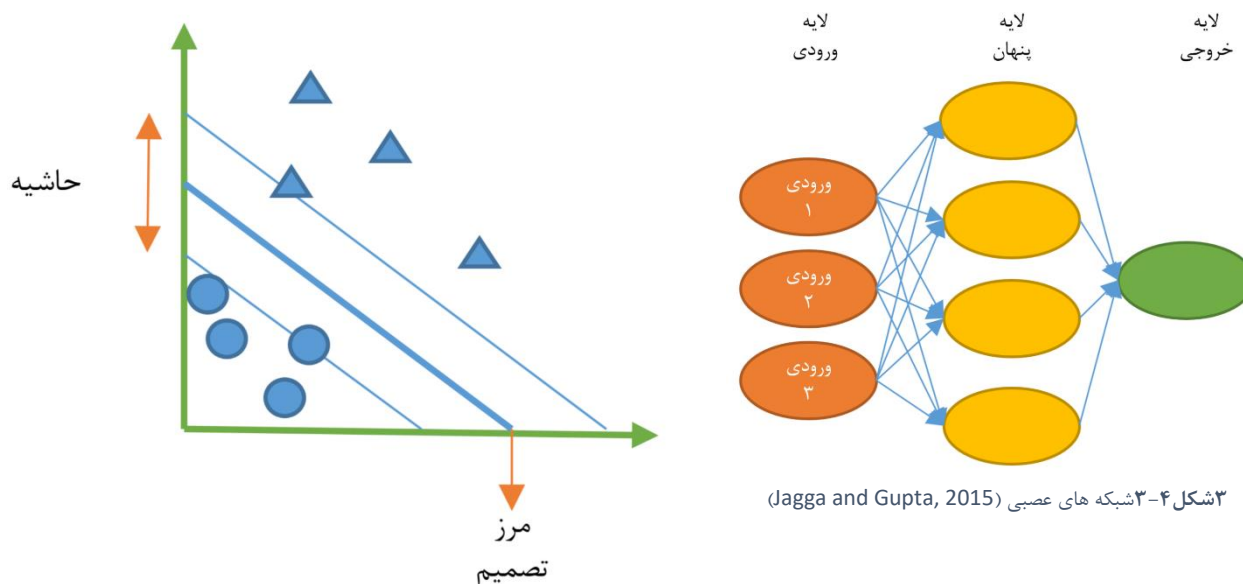
شکل ۳-۲ درخت تصمیم (Swan et al., 2013)

جنگل تصادفی: جنگل تصادفی بر مبنای مجموعه‌ای از درخت‌ها که بر اساس داده‌های آموزشی مدل شده‌اند ساخته می‌شود. هر کدام از درختان تصمیم به یک زیرمجموعه از ویژگی‌های نمونه تا دسترسی دارند و در آخر هنگام پیش‌بینی داده‌های تست، هر کدام از درختان یک دسته را برای داده موردنظر پیش‌بینی می‌کنند. هر کدام از دسته‌ها که بیشترین رأی را بیاورد به‌عنوان دسته آن داده موردنظر انتخاب می‌شود. (Meyfroidt et al., 2009)

بردارهای ماشین پشتیبان: این مدل به‌وسیله vapnik توسعه داده شد، شکل ... یک الگوریتم SVM را نمایش می‌دهد. این الگوریتم بر اساس مفهوم جدا پذیر بودن خطی داده‌ها را دسته‌بندی می‌کند. ویژگی به خصوصی که SVM را تعریف می‌کند عبارت‌اند از: (۱) تعیین معیاری که بهترین دسته‌بند خطی را بر اساس بیشینه کردن حاشیه تا تعریف می‌کند. (۲) شناسایی بردارهای پشتیبان که کم‌ترین تعداد نمونه داده‌های آموزشی نیاز است تا بهترین دسته‌بند خطی تعریف شود، به این دلیل که آن‌ها در مرز حاشیه تا قرار دارند. (۳) استفاده از کرنل‌ها برای انتقال ویژگی‌های اصلی به یک فضای بزرگ‌تر غیرخطی به‌منظور ایجاد جدا پذیری خطی است. الگوریتم SMO^۵ یکی از پرطرفدارترین الگوریتم‌های SVM است. (شکل ۳-۴)

شبکه‌های عصبی مصنوعی: شبکه‌های عصبی مصنوعی با الهام از ساختار و عملکرد مغز انسان توسعه داده شده‌اند. این شبکه از مجموعه‌ای از عناصر محاسبه‌گر (نرون) که بر اساس یک الگوی درون ارتباطی گستره به هم متصل شده‌اند. در کل هر نرون یک متغیر از یک دسته‌بند خطی است، اما با آمیختگی این نرون‌ها یک مدل دسته‌بندی پیچیده غیرخطی ساخته می‌شود که می‌تواند برای حل مسائل پیچیده استفاده شود. شکل ۳-۳ یک مدل شبکه عصبی مصنوعی را نمایش می‌دهد. (Swan et al., 2013).

⁵ Sequential minimal optimization



شکل ۳-۴ شبکه های عصبی (Jagga and Gupta, 2015)

شکل ۳-۳ بردارهای ماشین پشتیبان (Swan et al., 2013)

جدول ۳-۱ مقایسه مدل های دسته بندی (Swan et al., 2013)

روش	مزایا و معایب	سرعت یادگیری	سهولت تفسیر پذیری
پیز ساده	مزایا: قابلیت ساده و سریع در یکارگیری، مناسب برای مجموعه داده های با داده های کم شده، معایب: فرض مستقل بودن ویژگی ها از هم	۱	۴
درخت تصمیم	خروجی این الگوریتم به سادگی قابل تفسیر است اما بستگی به نوع الگوریتم مورد استفاده و پیچیدگی درخت ساخته شده دارد، همچنان مناسب برای مجموعه داده ها با داده کم شده	۲	۱
جنگل تصادفی	روشی کارآمد برای مجموعه داده های بزرگ اگرچه در مقابل outliera حساس نیست	۴	۳
دسته بندی های مبتنی بر قاعده	قواعد تولید شده به سادگی قابل خواندن است، مناسب برای کشف زیست نشانگرهای پنهان، اما امکان بیش برازش نیز دارد	۳	۱
بردارهای ماشین پشتیبان	استفاده از کرنل برای فراگیری توابع پیچیده، با این وجود بسیار کند بوده و چندین پارامتر نیز توسط کاربر باید تعریف شود	۵	۵
شبکه های عصبی مصنوعی	نتایج خروجی قابلیت خواندن ندارد و آموزش مدل ممکن است بسیار آهسته صورت گیرد.	۵	۵

۳.۲.۵ واریسی اعتبار^۶

از چالش های بزرگ داده کاوی و ساخت مدل بر روی داده های پزشکی اعتبار سنجی مدل بنا شده بر اساس داده های آموزشی است. بدین منظور داده تا به دو بخش داده های آموزش و داده های تست تقسیم می شوند و مدل بر اساس داده های آموزشی توسعه پیدا

^۶ Cross-validation

می‌کند سپس به کمک داده‌های تست مدل اعتبار سنجی می‌شود. یکی از رویکردهای ساده استفاده از نمونه‌برداری تصادفی RS^7 برای تقسیم داده‌ها به دو دسته داده‌های آموزشی و تست است؛ اما به دلیل این که فرآیندهای نمونه‌گیری پزشکی دشوار، پیچیده و هزینه‌بر است معمولاً تعداد نمونه بسیار کم است و از طرفی تعداد داده‌های آموزشی بر روی عملکرد و دقت مدل بسیار تأثیرگذار است. برای غلبه بر این مشکل از استفاده از روش‌های پیچیده‌تر نمونه‌گیری از مجموعه داده‌هاست که یکی از این روش‌ها واری اعتبار است. در این روش مجموعه داده به دو بخش داده‌های آموزش و تست تقسیم می‌شود و سپس مدل بر اساس داده‌های آموزش توسعه پیدا کرد و با داده‌های تست ارزیابی می‌شود. این رویه چند بار تکرار می‌شود و در آخر میزان خطا میانگین خطاهای به‌دست‌آمده از تکرار مراحل قبل است. یکی از رویکردهای رایج استفاده از الگوریتم واری اعتبار K -Fold است. در این روش مجموعه داده به k زیرمجموعه تقسیم شده و هر بار از یکی از این زیرمجموعه‌ها به‌عنوان داده تست و از $k-1$ داده دیگر به‌عنوان داده آموزش استفاده می‌شود این رویه انقدر ادامه می‌یابد تا تمام زیرمجموعه‌ها یک بار به‌عنوان داده تست انتخاب شوند. روش دیگر استفاده از الگوریتم واری اعتبار یکی-بیرون^۸ است. این روش دقیقاً مشابه رویکرد k -Fold است تنها با این تفاوت که تعداد k برابر برابر با تعداد نمونه‌هاست اما این روش حجم محاسبات را بالا می‌برد. (Thomas et al., 2006)

۳,۲,۶ بررسی عملکرد مدل

آخرین مرحله از فرآیند داده‌کاوی، مرحله بررسی عملکرد است. در این قسمت به چند روش برای ارزیابی عملکرد مدل‌های دسته‌بند در مواجهه با داده‌های تست می‌پردازیم.

دقت ۹: دقت یک مدل بر اساس نسبت میزان داده‌هایی که به‌درستی دسته‌بندی شده‌اند به مجموع داده‌هاست. معیار دقت با توجه به این اگر یک دسته به‌خصوص به‌طور معناداری بیشتر از سایر دسته‌ها باشد می‌تواند گمراه‌کننده باشد.

حساسیت و ویژگی ۱۰: برای دودسته حاصل از مدل چهار خروجی: درست-مثبت، درست-منفی، غلط-مثبت و غلط-منفی امکان‌پذیر است. حساسیت نسبت تعداد نمونه‌های مثبت که به‌درستی دسته‌بندی شده‌اند به کل نمونه‌های مثبت است؛ اما ویژگی در داده‌های پزشکی و تشخیصی احتمال یک شخص سالم به اشتباه به‌عنوان بیمار دسته‌بندی شود اس

منحنی ROC. این منحنی میزان حساسیت و ویژگی یک مدل تصمیم را بر اساس تمام ترشه‌لدهای^{۱۱} ممکن تصمیم بیان می‌کند. این منحنی یک تصویر کلی از عملکرد مدل به ما می‌دهد همچنان که می‌تواند برای انتخاب بهینه سطح بهینه تصمیمات که منجر به افزایش دقت مدل دسته‌بند می‌شود استفاده کرد. همچنان می‌توان از این منحنی برای مقایسه عملکرد چندین مدل دسته‌بند استفاده کرد. شکل ۵-۳ یک نمونه از این منحنی را که برای مقایسه عملکرد دو مدل استفاده‌شده است را نمایش می‌دهد. درمجموع هر چه سطح زیر نمودار یک مدل بیشتر باشد آن مدل بهتر است.

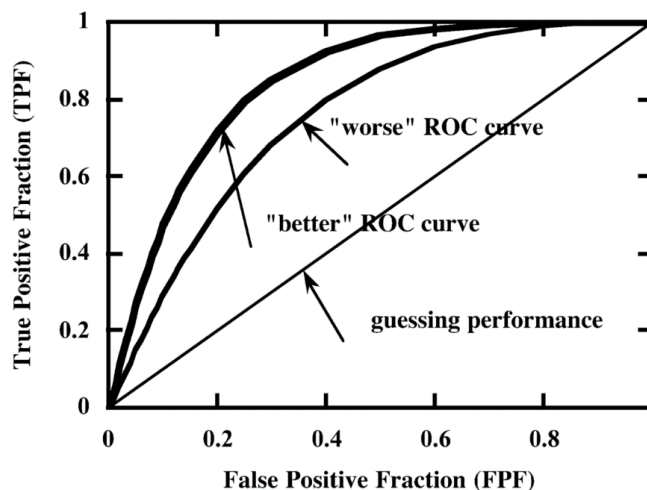
⁷ Random sampling

⁸ Leave-one-out cross validation

⁹ Accuracy

¹⁰ Sensitivity and specificity

¹¹ threshold



شکل ۵-۳ نمونه‌ای از منحنی ROC (Thomas et al., 2006)

۳.۳ مروری بر ادبیات کاربرد داده کاوی در دسته بندی و تشخیص زیست نشانگرهای سرطان

برای بررسی تحقیقات گذشته ما چند معیار برای انتخاب مقالات داشتیم. اول این که این مقالات برای ارائه مدل های دسته بندی و یا کشف زیست نشانگرها چه پردازش مستقیم بر روی مویک و چه پردازش بر روی پروتئین های شناسایی شده به کمک روش های های بدون برچسب گذاری و برچسب گذاری باشد. دوم این که فرآیند ساخت مدل بر روی داده های حاصل از طیف سنج جرمی به دلیل توان بالا این روش در کشف زیست نشانگرها و میزان بالا پروتئین های استخراجی که به طور معمول بین ۳۰۰۰ تا ۱۵۰۰۰ و بیش از آن است. بدین منظور از فهرست ارائه شده توسط دو مرجع (Swan et al., 2013) و (Jagga and Gupta, 2015) استفاده کردیم. اما یکی از مشکلات اساسی این مراجع، بیشتر مقالات فهرست شده مربوط به سال های قبل از ۲۰۱۱ بود. برای رفع این مشکل مجدداً در موتور جستجوی گوگل اسکولار (<https://scholar.google.com>)، پایگاه داده ساینس دایرکت (<http://www.sciencedirect.com>) و همچنین بانک مقالات NCBI^{۱۲} (<https://www.ncbi.nlm.nih.gov>) کلید واژه های "cancer"، "proteomic mass spectrometry"، "machine learning"، "data mining" را جستجو کرده و پس از یافتن مقالات مناسب با مقالات قبلی ادغام کرده و بخش مهم این مقالات که مواد و روش های تحقیق و بخش نتیجه گیری و/استنتاج بود مورد بررسی و نقد قرار دادیم.

۳.۳.۱ بررسی مقالات و تحقیقات صورت گرفته

در داده کاوی پروتئوم همان طور که اشاره شده دو هدف دسته بندی سرطان ها و کشف زیست نشانگرهای سرطان ها دنبال می شود. بیشتر تحقیقات صورت گرفته در این حوزه ارائه پردازش بروی مویک و به عبارت جدول داده های ماتریس شدت صورت گرفته است. برخی از این تحقیقات فقط به ارائه مدل های دسته بندی پرداخته اند مانند (Htike and Win, 2015) که از مدل ترکیبی

¹² National Center for Biotechnology Information

درخت لجستیک برای دسته‌بندی ۸۰ نمونه بیمار سرطان لوزالمعده و ۱۰۱ نمونه فرد سالم انجام داده‌اند. بیشترین مدل‌های استفاده‌شده در تحقیقات اخیر استفاده مدل الگوریتم‌های بردار ماشین پشتیبان و درخت تصمیم است و این به دلیل توانایی بالای این دو الگوریتم در برخورد با داده‌های گمشده است. نکته‌ای که قابل توجه است استفاده از انتخاب ویژگی برای کاهش بعد است که منجر به افزایش دقت معنادار دقت دسته‌بند فارغ از هدف تحقیق که خواه برای کشف زیست‌نشانگر یا خواه صرفاً برای دسته‌بندی باشد. تحقیق (Guan et al., 2009) یک مثال عالی برای تأیید این موضوع است. این تحقیق برای دسته‌بندی ۳۵ نمونه سرطان تخمدان و ۳۵ نمونه کنترلی است و برای دسته‌بندی از الگوریتم بردار ماشین پشتیبان استفاده کرده است. زمانی که از این الگوریتم استفاده کرده است به دقت ۸۳,۳٪ با واری اعتبار یکی-بیرون رسیده است. در نوبت دیگر از یک الگوریتم انتخاب ویژگی مبتنی بر ماشین پشتیبان همراه مدل استفاده کرده است و این بار به دقت ۹۷,۲٪ رسیده است که این مقدار تفاوت معناداری را نمایش می‌دهد. همچنان این تحقیق برای نمونه‌گیری برای انتخاب داده‌های تست مقایسه‌ای بین دو روش تقسیم تصادفی داده‌ها به دو دسته تست و آموزشی و روش واری اعتبار از الگوریتم یکی-بیرون استفاده شد پرداخته است که به‌وضوح برتری و دقت روش واری اعتبار اثبات شده است. البته باید به این نکته توجه داشت که نتایج این تحقیق بر روی پردازش یک مجموعه داده به دست آمده و نتیجه‌ای قطعی و اثبات شده نیست. همچنان لازم به ذکر است که در مرحله نمونه‌برداری برای بررسی عملکرد تقریباً نزدیک به تمام تحقیقات از روش‌های واری اعتبار استفاده کرده‌اند. در روش‌های جمع‌آوری نمونه مطلوب‌ترین نمونه‌ها برای کشف زیست‌نشانگرها، مایعات زیست پذیر نمونه‌های مربوط به خون است به دلیل این که خون در بدن انسان توسط قلب پمپاژ شده و در سرتاسر بدن به گردش درمی‌آید و به همین دلیل می‌تواند حاوی اطلاعات ارزشمندی جهت شناخت و توصیف نشانگرهای سرطانی باشد. در پیش‌پردازش داده‌ها بیشتر پژوهش‌ها از نرم‌افزار آماده استفاده کرده‌اند که فهرستی از مهم‌ترین نرم‌افزارهای مورد استفاده در پیوست ((الف)) آمده است. البته برخی تحقیقات مانند (Ushijima et al., 2007) از روش‌های خاصی برای پیش‌پردازش داده‌ها استفاده کرده‌اند به‌طور مثال برای نرمال‌سازی داده‌ها از فرمول زیر استفاده کرده‌اند به‌طوری که V_{max} نشان دهنده شدیدترین سیگنال مشاهده شده و V_{min} نشان دهنده ضعیف‌ترین سیگنال مشاهده شده است. همچنین V_i موجک موردنظر ما است که قصد نرمال‌سازی آن را داریم. اعداد به‌دست‌آمده در بازه [0-1] است.

$$NV_i = \frac{V_i - V_{min}}{V_{max} - V_{min}},$$

برای انتخاب موجک استفاده از الگوریتم نزدیک‌ترین همسایه با $K=10$ برای عرض پنجره‌های شناسایی موجک‌ها استفاده شده است.

برای تشخیص موجک‌های مشترک میان سوژه‌ها موردنظر از فرمول زیر که با میانگین هسته گوسی با مراکز موجک‌ها ساخته شده است استفاده کرده‌اند.

$$A(x) = \frac{1}{N_G} \sum_{i=1}^{N_G} \sum_j \exp \left[-\frac{(x - p_{i,j})^2}{(\sigma p_{i,j})^2} \right]$$

در این فرمول N_G بیانگر اندازه هر گروه، $p_{i,j}$ بیانگر مقدار m/z از i امین سوژه و j امین موجک است و σ بیانگر پارامتر محاسبه شده برای پهنای موجک‌ها است.

همچنین در این تحقیق مقایسه‌ای بین روش دسته‌بند الگوریتم *AdaBoost* و مدل دسته‌بند بردار ماشین پشتیبان با الگوریتم *SVM-RFE* صورت گرفته است که نتایج آن‌ها بر روی مجموعه داده‌هایشان برتری اندک *AdaBoost* را نشان می‌دهد اگرچه این اختلاف جزئی است اما به دلیل اینکه الگوریتم *SVM-RFE* محاسبات پیچیده‌تری و زمان بری را انجام می‌دهد الگوریتم *AdaBoost* را برتر می‌دانند.

در جدول ۲-۳ تحقیقات صورت گرفته از سال ۲۰۰۲ تا سال ۲۰۰۱۵ را به‌طور خلاصه آورده‌ایم. لازم به ذکر است که ارزیابی عملکرد مدل‌های دسته‌بند بین تحقیقات یکسان نبوده است. به‌گونه‌ای که برخی از معیارهای دقت و برخی از حساسیت و ویژگی استفاده کرده‌اند و بعضاً در برخی تحقیقات از معیارهای *ROC* و اندازه‌گیری سطح زیر منحنی استفاده شده است.

جدول ۲-۳ مرور تحقیقات گذشته دسته‌بندی سران و کشف زیست‌نشانگر

مرجع	دقت عملکرد دسته‌بند و زیست‌نشانگرهای شناسایی شده	روش‌های شناسایی، داده‌کاوی و ارزیابی	هدف پژوهش و مجموعه داده تا
(Sinues et al., 2015)	دقت دسته بند بالا ۹۰٪ و کشف هشت زیست نشانگر	استفاده از مدل SVM برای دسته بندی و انتخاب ویژگی	دسته بندی و کشف زیست نشانگرهای سرطان سینه با داده های بزاق دهان، ۱۴۱ نمونه بیمار و ۱۱ نمونه سالم، آنالیز داده تا SESI-MS
(Htike and Win, 2015)	دقت عملکرد دسته بند: ۷۴،۰۳۳۱٪	روش ترکیبی لوجستیک درخت برای مدل سازی، الگوریتم RELIEF برای انتخاب ویژگی، استفاده از روش TOP-HAT برای کاهش خط مبنا، استفاده از رویکرد یکی بیرون برای اعتبار سنجی	هدف دسته بندی سرطان لوزالمعده ۸۰ نمونه سرطان لوزالمعده در مقابل ۱۰۱ نمونه سالم،
(Guan et al., 2009)	بهترین دقت ۸۳٪ با استفاده از SVM و LOO-CV، رسیدن به دقت ۹۷،۲٪ با ترکیب مدل SVM و انتخاب ویژگی بر مبنای SVM، شناسایی ۳۸ زیست نشانگر با چهار رویکرد متفاوت	شناسایی با نرم افزار mzMine (v0.60)، SVMs با روش های مرتبط انتخاب ویژگی، LOO-CV، 12-FOLD CV، 52-20-split validation	دسته بندی و کشف زیست نشانگرهای سرطان تخمدان، ۳۷ بیمار بیمارمان مبتلا به سرطان تخمدان پاپیلری سروز و ۳۵ نمونه کنترلی
(Vlahou et al., 2003)	۸۵٪ دقت دسته بندی	دسته بندی بر اساس درخت-رگرسیون (CART)، اعتبار سنجی 10-FOLD، برای دسته بندی موجک استفاده از نرم افزار CIPHERGEN Systems، نرم افزار الگوهای زیست نشانگر (BPS)	دسته بندی سرطان بدخیم، خوشخیم و سالم در سرطان تخمدان، ۴۴ سرطان بدخیم، ۶۱ سرطان خوشخیم و ۳۴ نمونه سالم
(Le et al., 2005)	۸۵٪ دقت و پیدا کردن چند زیست نشانگر	استفاده از نرم افزار mascot و استفاده از C.SVM با استفاده از اعتبار سنجی یکی-بیرون	دسته بندی نمونه های پروستات از ۱۹ بیمار با متاستاز استخوان و ۱۹ بیمار فاقد آن
(Adam et al., 2002)	۹۰٪ از داده های تست به درستی دسته بندی شدند.	داده تست شامل: ۱۵ کنترل، ۱۵ خوشخیم و ۳۰ پروستات. شناسایی	دسته بندی سرطان پروستات و کنترل، ۹۷ نمونه کنترلی، ۹۲

نمونه سرطان خوشخیم، ۱۹۷ سرطان پروستات	موجک با نرم افزار CIPHERGEN SELDI software و همگام سازی موجک با الگوریتم peakminer، استفاده از درخت تصمیم.		
شناسایی زیست نشانگر برای سرطان پروستات، ۱۷۹ سرطان ادرنوکازسینوما و ۷۴ سرطان خوشخیم	پیش پردازش: گروه بندی، تصحیح خط مبنا و نرمال سازی با استفاده از نرم افزار TOFWorks	دقت ۸۷٫۹٪، ۲۶ موجک به عنوان زیست نشانگرهای محتمل شناسایی شد.	(Oh et al., 2009)
شناسایی زیست نشانگرهای سرطان سر و گردن، پنج مجموعه از چهار نمونه به همراه نمونه کنترل برای هر مجموعه	استفاده از نرم افزارهای iTRQ و ProteinPilot. استفاده از مدل بیز ساده و اعتبار سنجی 3-FOLD	سه تا از زیست نشانگرها شناخته شد.	(Ralhan et al., 2008)
توسعه پنل ^{۱۳} برای شناسایی زیست نشانگرهای سرطان سینه، ۴۰ نمونه پلازما خون از سرطان سینه و ۴۰ نمونه از سالم	شناسایی پروتئین تا به صورت بدون برچسب گذاری با کمک نرم افزار Eli Lilly، استفاده از مدل شبکه های عصبی مصنوعی	۸۵٪ دقت در داده های تست، دوتا از بهترین پنج پنل پروتئین شناخته شد که شامل هفت پروتئین شد.	(Zhang and Chen, 2009)
کشف زیست نشانگرها، ۶۵ نمونه از بیماران مبتلا به سرطان سینه، سپس نمونه گیری مجدد از آنها بعد از مصرف چهار هفته ای docetaxel 75 mg/m2	استفاده از نرم افزار SpecAlign برای کاهش خط مبنا، همگام سازی موجک، استفاده از مدل AdBoost برای دسته بندی و 5- FOLD برای CV	۶ زیست نشانگر شناخته شد	(Ushijima et al., 2007)
شناسایی زیست نشانگرها، ۱۳۲ بیمار با سرطان لنفوم B-Cell و ۷۵ نمونه کنترلی، داده های آنالیز شده SELDI-TOF-MS از سرم خون	استفاده از درخت تصمیم برای مدل سازی و استفاده از نرم افزار CIPHERGEN برای پیش پردازش داده تا	نه موجک به عنوان زیست نشانگر بالقوه انتخاب در نظر گرفته شد، چهار موجک برای پیش بینی پاسخ بیماران به درمان های استاندارد، حساسیت ۹۴٪، ویژگی ۹۴٪ در ۸۵ نمونه از مجموعه تست، حساسیت ۹۴٪ و ویژگی ۹۲٪ با ۶۶ نمونه تست.	(Zhang et al., 2007)
کشف زیست نشانگرهای سرطان معد، ۷۹ نمونه از بیماران مبتلا به سرطان معد و ۳۳ نمونه از افراد فاقد سرطان که ۱۰ نفر از آنها دارای التهاب معد بودند.	استفاده از مدل دسته بند مبتنی بر SVM و برای واری اعتبار از الگوریتم 10-FOLD استفاده شده است.	۹ نشانگر یافت شده با ۸۹٪ دقت در عملکرد	(Cohen et al., 2011)
شناسایی زیست نشانگرها که موجب متاستاز استخوانی در سرطان سینه می شود، نمونه از ۱۱۱ زن مبتلا به سرطان سینه به دو گروه تقسیم شده که گروه اول شامل ۴۱ نفر متاستاز استخوان و ۳۶ نفر دیگر فاقد آن، گروه دیگر	استفاده از جنگل تصادفی برای مدل دسته بند، ساخت درخت ۱۰۰۰ مرتبه تکرار شده و هر بار در هر گره سه ویژگی مورد آزمایش قرار می گرفت	شناسایی ۱۳ زیست نشانگر، قدرت دسته بند با حساسیت ۹۱٪ و ویژگی ۹۳٪	(Washam et al., 2013)

(منظور تعدادی از پروتئین ها که با هم کار می کنند و ظاهر می شوند) Panel^{۱۳}

شامل ۱۷ نفر متاستاز استخوان و ۱۷ نفر دیگر فاقد آن			
--	--	--	--

۳,۴ خلاصه فصل

در این فصل به بررسی مراحل داده کاوی بر روی پروتئوم پرداختیم. پیش پردازش داده ها معمولاً برای مواردی است که مستقیماً پردازش را بر روی موجک ها انجام می دهیم. فعالیت های که در پیش پردازش استفاده می شود شامل کاهش خ مبنا، کاهش نویز و نرمال سازی و انتخاب موجک هاست. شش مدل درخت تصمیم، بردارهای ماشین پشتیبان، جنگل تصادفی، مدل های مبتنی بر قاعده، بیز ساده و شبکه های عصبی مصنوعی را معرفی کردیم و ضعف و قدرت آن ها را بیان کردیم. برای واری اعتبار روش های یکی-بیرون و K-FOLD را معرفی کردیم. سپس معیارهای ارزیابی مدل مانند ویژگی، حساسیت، دقت و همچنین ROC را شرح دادیم. در مرحله آخر به مرور ادبیات کاربرد داده کاوی در دسته بندی و کشف زیست نشانگرهای سرطانی پرداختیم. مشاهده کردیم که بیشترین مدل های استفاده شده بردارهای ماشین پشتیبان و درخت تصمیم بودند. همچنین استفاده از روش های انتخاب ویژگی علاوه بر این که برای کشف زیست نشانگرها ضروری است موجب دقت و سرعت عملکرد مدل دسته بند می شود.

فصل چهارم

جمع‌بندی و نتیجه‌گیری

۴.۱ مقدمه

در این فصل ابتدا به مرور فصل‌های گذشته خواهیم پرداخت و سپس به بیان چالش‌ها و فرصت‌های مطالعات سرطان با بکارگیری داده‌کاوی پروتئوم خواهیم پرداخت و پس از آن مجموعه داده‌های خود را شرح خواهیم و تعدادی بانک‌داده مناسب برای تحقیق معرفی خواهیم کرد.

۴.۲ مروری بر فصل‌های گذشته

در فصل اول به بیان مفاهیم و کلیات پرداختیم. دانستیم که انسان تقریباً متشکل از چهار ابرمولکول کربوهیدرات‌ها، پروتئین‌ها، اسیدهای نوکلئوتید و لیپیدها است و در این بین پروتئین‌ها به دلیل رفتار پویا در بدن انسان، تغییرات پس از ترجمه و همچنین معادل نبودن میزان تولیدشان مطابق با mRNA گزینه‌ای مناسب‌تر برای مطالعات بیماری‌ها در سطح مولکولی است. سرطان را بیماری ژنتیکی در سطح مولکولی تعریف کردیم که از تکثیر غیرعادی سلول‌ها به دلیل کم‌کاری، پرکاری و یا عدم حضور پروتئینی خاص در یک بافت یا ارگان از بدن بروز می‌کند؛ بنابراین میزان پروتئین‌های بیان‌شده، الگو و رفتار آن‌ها می‌تواند معیاری مناسب برای تشخیص بیماری‌ها و کشف زیست‌نشانگرها - ابرمولکول‌ها - مرتبط با بیماری‌ها به‌خصوص در سرطان باشد. به مجموع پروتئین‌های بیان‌شده و الگوی آن‌ها در یک لحظه خاص در سلول یا یک بافت پروتئوم گفته می‌شود و علم بررسی پروتئوم را پروتئومیکس می‌گویند. در علم پروتئومیکس کشف این تغییرات و کشف این زیست‌نشانگرها می‌تواند منجر به ساخت داروهای جدید و همچنین ابداع روش‌های درمانی جدید نیز گردد. بیان کردیم، به دلیل این که نمونه‌های بیمارستانی معمولاً تعداد نمونه‌هایشان کم اما ویژگی و متغیرهای بسیاری دارند برای تحلیل و دسته‌بندی آن‌ها به‌خصوص جهت کشف زیست‌نشانگرها نیاز به بهره‌گیری از علوم مانند داده‌کاوی، یادگیری ماشین و هوش مصنوعی اجتناب‌ناپذیر است. سپس در پایان فصل اول برخی از اصطلاحات مورد نیاز برای ورود به دنیای پروتئومیکس را تعریف کردیم و پس از آن آرایش کلی گزارش خود را ارائه کردیم.

در فصل دوم ابتدا به سراغ روش‌ها و رویکردهای پروتئومیکس به حل مسائل پرداختیم. دانستیم که پروتئومیکس از روش‌های شناسایی، توصیف و کمیت شماری پروتئین‌ها برای مقایسه و دسته‌بندی دو دسته‌بیمار و دسته‌کنترل جهت دسته‌بندی سرطان‌ها و کشف زیست‌نشانگرها استفاده می‌کند. در رویکرد تغییرات پس از ترجمه به اتفاقات و فعالیت‌های پروتئین‌ها پس از ترجمه در سلول و بافت را مشاهده و اندازه‌گیری می‌کند. در مکان‌یابی پروتئین‌ها به بررسی مکانی پروتئین‌ها بنا به شرایط بالینی مختلف

انسان می پردازد و در برهم کنش پروتئین ها، تعاملات پروتئین ها را در یک شبکه ارتباطی با یکدیگر می سنجد؛ که هدف ما در این تحقیق پروتئئومیکس کمی و مقایسه ای بود. برای دریافت نمونه از بیماران سه نوع نمونه گیری مایعات زیست پذیر، بافت ها و سل لاین ها وجود دارد که با وجود اینکه بافت ها اطلاعات ارزشمندی را در خود دارند اما به دلیل مشکلات نمونه برداری و تهاجمی بودن آن ها توصیه نمی شود ولی در مقابل مایعات زیست پذیر هم به دلیل گردهش در بدن اطلاعات متنوعی در خود دارند و به خصوص نمونه های ادرار به دلیل غیرتهاجمی بودن نمونه های مناسب برای پروتئومیکس است. فناوری های بکار برده شده در پروتئومیکس را به سه دسته عمده فناوری ریزآرایه های پروتئینی، ژل الکتروفورز دوبعدی پلی آکریل آمید (2D-PAGE) و طیف سنج جرمی تقسیم بندی کردیم و عنوان نمودیم که طیف سنج جرمی به دلیل سهولت و توان بالا در آنالیز حجم عظیمی از پروتئین ها نسبت به روش ها دیگر بیشتر مورد توجه محققین قرار گرفته است و نکته مهم دیگر اینکه غیر جانب دارانه است یعنی برخلاف فناوری های دیگر که از قبل فرض بر وجود تعداد مشخصی پروتئین در نمونه گرفته می شود و برای آنتی بادی جهت به دام انداختن پروتئین ها ساخته می شود، است و تعداد و نوع پروتئین ها از قبل مشخص نیست. طیف سنج خود به انواع مختلف از لحاظ فناوری لیزر تقسیم می شود اما دو رویکرد عمده در رابطه با آن وجود دارد رویکرد اول از پائین به بالا هست که ابتدا پروتئین ها را تجزیه به پپتیدها و شاخه های آمینواسیدهای آن ها می کنند و سپس آنالیز می شود و بعد از آنالیز با فن آوری های خاصی مجدداً نوع پروتئین ها حاصل از آنالیز شناسایی می کنند در رویکرد دوم یا از بالا به پائین پروتئین ها بدون تجزیه شدن به قطعات کوچک تر خود مستقیماً تفکیک، کمیت شماری و شناسایی می شوند که رویکرد دوم به دلیل مشکلات آزمایشگاهی پیچیده و هزینه های بیشتر مرود توجه محققین نبوده است. ما هم تمرکز خود را به داده های حاصل از رویکرد اول گذاشتیم و دانستیم نتایج حاصل از رویکرد اول به دو صورت عمده مورد پردازش قرار می گیرد. اول این که مستقیماً به سراغ موجک ها حاصل برویم که این امر نیاز به پیش پردازش داده ها و همچنین آنالیزهای اضافه تر برای شناسایی پروتئین ها دارد و در روش دوم استفاده از موتورهای جستجو بانک توالی اطلاعات پروتئینی است که در پیوستگی استفاده از روش ها برچسب گذاری که روشی آزمایشگاهی است و یا غیر برچسب گذاری که روشی غیر آزمایشگاهی و مبتنی بر علوم کامپیوتری است.

در فصل سوم به سراغ مراحل داده کاوی بر روی پروتئوم جهت دسته بندی و کشف زیست نشانگرها رفتیم و بیان کردیم زمانی که مستقیماً به سراغ موجک ها برویم نیاز به پیش پردازش داده ها از قبیل کاهش خط مبنا، نرمال سازی داده ها و حذف نویزها داریم که برای این مسئله نرم افزارهای کامپیوتری مناسب و همچنین روش های مختلف وجود دارد. از اهمیت کاهش بعد در داده های پروتئینی گفتیم و بیان کردیم برای کشف زیست نشانگرها انتخاب ویژگی یک امر ضروری است و برای دسته بندی سرطان ها هم موجب افزایش سرعت و دقت دسته بند ما می شود سپس به سراغ مدل های پر استفاده دسته بند در تحقیقات چند سال اخیر رفتیم و شش مدل: درخت تصمیم، بیز ساده، مدل مبتنی بر قاعده، بردارهای ماشین پشتیبان، جنگل تصادفی و شبکه های عصبی را تشریح کردیم و در یک جدول به طور خلاصه آن ها را با یکدیگر مقایسه نمودیم. سپس به سراغ روش های نمونه گیر از داده ها رفتیم و بیان کردیم به دلیل دشواری در نمونه گیر داده های پزشکی تعداد سطرهای داده ها بسیار کم است و نیاز به روش های پیچیده نمونه برداری جهت داده تست و داده آزمایش برای ساخت و ارزیابی مدل داریم. روش واری اعتبار را معرفی کردیم و دو نمونه K-FOLD و یکی بیرون را که از روش های واری اعتبار است شرح دادیم. در قسمت بعدی به بیان معیارهای ارزیابی مدل دسته بند خود پرداختیم و معیارهای دقت، ویژگی، حساسیت و منحنی ROC را معرفی نمودیم. در بخش آخر مروری انتقادی به تحقیقات و پژوهش های صورت گرفته داشتیم و آن ها را در یک جدول خلاصه بندی کردیم. دانستیم که پرکاربردترین مدل دسته بند الگوریتم بردار ماشین های پشتیبان و درخت تصمیم است به دلیل قدرت و دقت بالای دسته بندی که البته نشان دادیم در همه تحقیقات این نکته لزوماً صدق نمی کند. همچنین لزوم استفاده از انتخاب ویژگی را با نمایش تأثیرگذاری آن در یکی از تحقیقات نشان دادیم و

برای پیش‌پردازش داده‌ها علاوه بر استفاده از نرم‌افزارهای کامپیوتری استفاده از روش‌ها ابتکاری و جدید نیز صورت می‌گیرد که استفاده آن در یکی از تحقیقات را بیان کردیم. نشان دادیم برخی تحقیقات علاوه بر دسته‌بندی، در کشف زیست‌نشانگرها به دنبال یافتن پنل‌ها هستند. پنل‌ها به مجموعه‌ای از پروتئین‌ها گفته می‌شود که با یکدیگر در سلول فعالیت می‌کنند و با یکدیگر نیز ظاهر می‌شوند.

۴,۳ بررسی چالش‌ها و پیشنهاد فرصت‌ها

تعداد زیادی چالش همچنان در طراحی آزمایش‌ها، تحلیل داده‌ها، استانداردسازی روش‌ها و ارائه و توزیع داده‌های حاصل از پروتئومیکس در بانک داده‌ها موجود است. اثرات پروتئومیکس در تحقیقات پزشکی روزبه‌روز در حال بیشتر شدن است و ابزارهای پروتئومیکس به‌طور گسترده‌ای برای شناسایی و درمان بیماری‌ها بکار برده می‌شود. نقطه کلیدی در این بین، شناسایی زیست‌نشانگرها است که عدم وجود یا تغییر در فراوانی آن‌ها در حالت زیستی سلول‌ها تأثیرگذار است. با توجه تحقیقات قبلی در حوزه‌های کشف زیست‌نشانگرها هنوز هم نیاز به مطالعات بیشتر جهت شناسایی هرچه بیشتر زیست‌نشانگرها وجود دارد. به‌طور مثال می‌توان از روش دسته‌بندی مبتنی بر قاعده که در بررسی و تحلیل ریزآرایه‌ها در بیوانفورماتیک استفاده می‌شود بهره برد. بدین‌صورت که شناسایی شبکه‌ای از زیست‌نشانگرها که در یک قاعده خاص باهم ظاهر می‌شوند. همچنین از روش‌های علوم ژئومیکس مانند آنالیز مجموعه‌ای از ژن‌ها مشخص در علم پروتئومیکس استفاده کرد. بدین‌صورت که دسته‌ای از پروتئین‌های کاندید را در دسته‌های کنترل و بیمار شناسایی و کمیت‌شماری و مقایسه کنیم. همچنان بحث پیش‌پردازش و کاهش ابعاد در پروتئومیکس جدی است. با وجود تلاش‌های صورت گرفته اما همچنان در بسیاری از تحقیقات نتایج مناسب نیست. به دلیل این‌که تعداد نمونه‌ها در داده‌های پروتئومیکس کم است دقت دسته‌بندی‌های گزارش‌شده می‌تواند گمراه‌کننده باشد و زمانی که این داده‌ها در دنیای واقعی با مقادیر نمونه‌های بسیار زیاد استفاده شود نرخ خطاها بسیار بیشتر می‌شود به همین دلیل توصیه می‌شود تا به بحث پیش‌پردازش داده‌ها و روش‌های جدید کاهش ابعاد در علم پروتئومیکس پرداخته شود. همچنین جای خالی ارتباط سایر بیماری‌ها به‌طور مثال چاقی و یا دیابت با یک سرطان خاص یا مجموعه‌ای از سرطان‌ها به‌شدت احساس می‌شود.

۴,۴ داده‌های تحقیق

در فصل دوم به طور مفصل از جنس داده‌های پروتئومیکس و نمایش آنها صحبت کردیم. داده‌های در دسترس ما مجموعه داده‌های با همکاری بخش آزمایشگاه پروتئومیکس دانشگاه کالیفرنیا جنوبی است که شامل ۲۰ نمونه مبتلا به سرطان لوزالمعده، ۲۰ نمونه عفونت مزمن لوزالمعده و ۲۰ نمونه از اشخاص سالم است. نمونه‌گیری از پلاسما خون است و این داده‌ها با فن‌آوری ESI-MS آنالیز شده است. همچنین برای توصیف و شناسایی و کمیت‌شماری پروتئین‌ها از روش بدون برچسب و با استفاده از نرم‌افزار sequest حاصل شده است. شکل ۴-۱ شمایی از داده‌ها را که در یک فایل Excel است را نمایش می‌دهد.

همچنین بانک‌های داده‌ای مناسب بسیار زیادی برای داده‌ها پروتئینی موجود است مانند:

- اطلس پروتئینی انسان (<http://www.proteinatlas.org>) که داده‌های پروتئین آن با فن‌آوری ریز آرایه‌ها جمع‌آوری شده است. بانک داده‌ای شامل ۴۴ بافت نرمال و ۲۰ بافت سرطانی، ۵۶ سل لاین و سطح بیان ترنسکریپت‌ها است.
- پایگاه داده مرجع پروتئین انسانی (<http://www.hprd.org>): که اطلاعات مناسبی در رابطه با تعاملات پروتئینی در یک شبکه دارد.
- مرکز ملی اطلاعات زیست‌فن‌آوری (<https://www.ncbi.nlm.nih.gov>): اطلاعات، مقالات و کتب بسیاری در رابطه با ژنومیکس و پروتئومیکس دارا است.
- بانک داده پروتئین (<http://www.rcsb.org/pdb/home/home.do>)
- بانک داده شناسایی پروتئین‌ها (<https://www.ebi.ac.uk>)
- Proteopedia (<http://proteopedia.org>)
- UniProt (<http://www.uniprot.org>)

A		B	C		D	E	F	G	H	I
Gene	Protein	Protein description	Peptide	charge	A1	A2	A9	B2		
1 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	GWMNDPNGLWYDEK	2	63976296	24769764	40046584	5498		
3 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	GWM[+16]NDPNGLW	2	6197049	2523306	12923109	104		
4 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	VFWYEPSSQK	2	8723820	8625410	7146495	496		
5 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	WIMTAAK	2	1086354	1498491	6836681	602		
6 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	WIM[+16]TAAK	2	177862	12706056	851060	777		
7 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	IEIYSSDDLK	2	190592160	118888592	136130736	1684		
8 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	IEIYSSDDLKSWK	3	742601	1987192	2445724	477		
9 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	SSMSLVRK	2	210993	484660	3281242	54		
10 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	KFSLNTEYQANPETELIN	3	8995346	22224170	23039502	563		
11 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	FSLNTEYQANPETELINL	3	41705792	41687748	46593780	4379		
12 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	AEPLNISNAGPWSR	2	5289871	6530384	16779556	239		
13 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	GLEDPEEYLR	2	522095232	784777280	649137216	56294		
14 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	MGFEVSASSFFLDR	2	4875871	5925656	7153578	575		
15 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	MGFEVSASSFFLDR	3	75067	1966584	2582508	198		
16 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	M[+16]GFEVSASSFFLDR	2	1367760	1	1			
17 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	M[+16]GFEVSASSFFLDR	3	3904164	9913573	19360182	628		
18 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	MGFEVSASSFFLDRGNL	3	1511042	5894734	4381281	235		
19 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	ENPYFTNR	2	247100000	292735456	193823904	7914		
20 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	MSVNNQPFK	2	4375857	81305216	2918126	2939		
21 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	M[+16]SVNNQPFK	2	4616198	5216906	1125524	1895		
22 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	MSVNNQPFKSENDLSY	3	3358167	3965633	3892572	429		
23 SUC2	P00724	Invertase 2 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) GN=SUC2 PE=1 SV=1	SENDLSYK	2	153155616	173682416	130623936	6286		
SUC2 Total					1274133118	1607309228	1311073297	92276		
25 AKAP13	T12802	A-kinase anchor protein 13 OS=Homo sapiens GN=AKAP13 PE=1 SV=2	STPSLP[+57]M[+16]V	3	383171360	396120864	449970880	10815		
26 AKAP13	T12802	A-kinase anchor protein 13 OS=Homo sapiens GN=AKAP13 PE=1 SV=2	SVSIQNTIGVGNDENM	3	25716	45718	14628	21		
AKAP13 Total					383197076	396166582	449985508	10818		
28 ECD	O95905	Protein edcynoseless homolog OS=Homo sapiens GN=ECD PE=1 SV=1	EEKEQNYDLTEVSESM[4	29705026	8415698	4520073	747		
ECD Total					29705026	8415698	4520073	747		

شکل ۴-۱ نمونه‌ای از داده‌های آزمایشگاه پروتئومیکس دانشگاه کالیفرنیا جنوبی

۴,۵ خلاصه فصل

در این فصل ابتدا مروری بر مطالعات فصل‌های گذشته داشتیم و خلاصه‌ای از نتیجه‌گیری‌های خود را ارائه کردیم. سپس به بیان چالش‌ها و فرصت‌ها پرداختیم و از لزوم پرداختن به مسئله کاهش بعد و نگاه کردن به سرطان از رویکردهای دیگر صحبت به میان آوردیم و چند موضوع برای مطالعات آتی پیشنهاد دادیم. بعد از آن به معرفی داده‌های در دسترس خود و چند بانک اطلاعاتی جامع و مناسب که برای تحقیقات آتی مناسب است را معرفی کردیم.

مراجع

مراجع فارسی

۱. تیمورپور، بابک؛ نجفی حیدر، ۱۳۹۴، "داده کاوی با R به همراه متن کاوی و تحلیل شبکه های اجتماعی". ویرایش ۱، تهران: مرکز تحقیقات و توسعه سازمان اتکا.
۲. علی پور، محمد. ۱۳۸۴، "ارائه روشی برای تشخیص بیماری سرطان به کمک داده کاوی پروتئوم انسان"، پایان نامه کارشناسی ارشد مهندسی برق-الکترونیک، دانشکده فنی مهندسی، دانشگاه تربیت معلم سبزواری
۳. پارسا، ن. ۲۰۱۲. اساس سلولی و مولکولی سرطان در انسان، مقاله مروری. مجله سلول و بافت، 2, (Cell & Tissue Journal), 365-376.
۴. شیردل، س. ا.، عالمی، م. & خلیفه، خ. ۲۰۱۳. مباحث نوین در زیست شناسی: پروتئومیکس و نانوزیست فناوری. زیست فناوری دانشگاه تربیت مدرس، ۴، ۹-۲۲

مراجع انگلیسی

5. ADAM, B.-L., QU, Y., DAVIS, J. W., WARD, M. D., CLEMENTS, M. A., CAZARES, L. H., SEMMES, O. J., SCHELLHAMMER, P. F., YASUI, Y. & FENG, Z. 2002. Serum protein fingerprinting coupled with a pattern-matching algorithm distinguishes prostate cancer from benign prostate hyperplasia and healthy men. *Cancer research*, 62, 3609-3614.
6. COHEN, M., YOSSEF, R., EREZ, T., KUGEL, A., WELT, M., KARPASAS, M. M., BONES, J., RUDD, P. M., TAIEB, J. & BOISSIN, H. 2011. Serum apolipoproteins CI and C-III are reduced in stomach cancer patients: results from MALDI-based peptidome and immuno-based clinical assays. *PloS one*, 6, e14540.
7. DUBITZKY, W., GRANZOW, M. & BERRAR, D. P. 2007. *Fundamentals of data mining in genomics and proteomics*, Springer Science & Business Media.
8. DZIUDA, D. M. 2010. *Data mining for genomics and proteomics: analysis of gene and protein expression data*, John Wiley & Sons.
9. ELO, L. L. & SCHWIKOWSKI, B. 2012. Mining proteomic data for biomedical research. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2, 1-13.

10. GALLEG0, M. & VIRSHUP, D. M. 2007. Post-translational modifications regulate the ticking of the circadian clock. *Nature Reviews Molecular Cell Biology*, 8, 139-148.
11. GARRETT, R. & GRISHAM, C. 2005. Biochemistry: Belmont, CA. CA: Thomson Brooks/Cole, 32005.
12. GUAN, W., ZHOU, M., HAMPTON, C. Y., BENIGNO, B. B., WALKER, L. D., GRAY, A., MCDONALD, J. F. & FERNÁNDEZ, F. M. 2009. Ovarian cancer detection from metabolomic liquid chromatography/mass spectrometry data by support vector machines. *BMC bioinformatics*, 10, 259.
13. HAN, J., PEI, J. & KAMBER, M. 2011. *Data mining: concepts and techniques*, Elsevier.
14. HTIKE, Z. Z. & WIN, S. L. 2015. Premalignant Pancreatic Cancer Diagnosis Using Proteomic Pattern Analysis. *Journal of Medical and Bioengineering Vol*, 4.
15. JAGGA, Z. & GUPTA, D. 2015. Machine learning for biomarker identification in cancer research—developments toward its clinical application. *Personalized Medicine*, 12, 371-387.
16. JOERGER, A. C. & FERSHT, A. R. 2016. The p53 pathway: Origins, inactivation in cancer, and emerging therapeutic approaches. *Annual review of biochemistry*, 85, 375-404.
17. KATAJAMAA, M. & OREŠIČ, M. 2005. Processing methods for differential analysis of LC/MS profile data. *BMC bioinformatics*, 6, 179.
18. LAM, S., JIMENEZ, C. & BOVEN, E. 2014. Breast cancer classification by proteomic technologies: current state of knowledge. *Cancer treatment reviews*, 40, 129-138.
19. LE, L., CHI, K., TYLDESLEY, S., FLIBOTTE, S., DIAMOND, D. L., KUZYK, M. A. & SADAR, M. D. 2005. Identification of serum amyloid A as a biomarker to distinguish prostate cancer patients with bone lesions. *Clinical chemistry*, 51, 695-707.
20. LI, L., TANG, H., WU, Z., GONG, J., GRUIDL, M., ZOU, J., TOCKMAN, M. & CLARK, R. A. 2004. Data mining techniques for cancer detection using serum proteomic profiling. *Artificial intelligence in medicine*, 32, 71-83.
21. MEYFROIDT, G., GÜIZA, F., RAMON, J. & BRUYNOOGHE, M. 2009. Machine learning techniques to examine large patient databases. *Best Practice & Research Clinical Anaesthesiology*, 23-۱۲۷ , ۱۴۳
22. NEILSON, K. A., ALI, N. A., MURALIDHARAN, S., MIRZAEI, M., MARIANI, M., ASSADOURIAN, G., LEE, A., VAN SLUYTER, S. C. & HAYNES, P. A. 2011. Less label, more free: approaches in label-free quantitative mass spectrometry. *Proteomics*, 11, 535-553.
23. OH, J. H., LOTAN, Y., GURNANI, P., ROSENBLATT, K. P. & GAO, J. 2009. Prostate cancer biomarker discovery using high performance mass spectral serum profiling. *Computer methods and programs in biomedicine*, 96, 33-41.
24. OTTO, T. & SICINSKI, P. 2017. Cell cycle proteins as promising targets in cancer therapy. *Nature Reviews Cancer*, 17, 93-115.
25. PANIS, C. 2015. Proteomic Tools for Cancer Research: Updating the Oncoproteomics. *Journal of Proteomics & Bioinformatics*, 1.
26. RALHAN, R., DESOUZA, L. V., MATTA, A., TRIPATHI, S. C., GHANNY, S., GUPTA, S. D., BAHADUR, S. & SIU, K. M. 2008. Discovery and verification of head-and-neck cancer biomarkers by differential protein expression analysis using iTRAQ labeling, multidimensional liquid chromatography, and tandem mass spectrometry. *Molecular & Cellular Proteomics*, 7, 1162-1173.
27. SHUKLA, S., GUPTA, D. L. & PRASAD, B. R. 2016. Comparative Study of Recent Trends on Cancer Disease Prediction using Data Mining Techniques. *International Journal of Database Theory and Application*. ۱۱۸-۱۰۷ , ۹ ,
28. SINUES, P. M.-L., LANDONI, E., MICELI, R., DIBARI, V. F., DUGO, M., AGRESTI, R., TAGLIABUE, E., CRISTONI, S. & ORLANDI, R. 2015. Secondary electrospray ionization-mass spectrometry and a novel statistical bioinformatic approach identifies a cancer-related profile in exhaled breath of breast cancer patients: a pilot study. *Journal of breath research*, 9, 031001.

29. SOSA, M. S., BRAGADO, P. & AGUIRRE-GHISO, J. A. 2014. Mechanisms of disseminated cancer cell dormancy: an awakening field. *Nature Reviews Cancer*, 14, 611-622.
30. SWAN, A. L., MOBASHERI, A., ALLAWAY, D., LIDDELL, S. & BACARDIT, J. 2013. Application of machine learning to proteomics data: classification and biomarker identification in postgenomics biology. *Omics: a journal of integrative biology*, 17, 595-610.
31. THOMAS, A., TOURASSI, G. D., ELMAGHRABY, A. S., VALDES, R. & JORTANI, S. A. 2006. Data mining in proteomic mass spectrometry. *Clinical Proteomics*, 2, 13.
32. USHIJIMA, M., MIYATA, S., EGUCHI, S., KAWAKITA, M., YOSHIMOTO, M., IWASE, T., AKIYAMA, F., SAKAMOTO, G., NAGASAKI, K. & MIKI, Y. 2007. Common peak approach using mass spectrometry data sets for predicting the effects of anticancer drugs on breast cancer. *Cancer informatics*, 3, 285.
33. VLAHOU, A., SCHORGE, J. O., GREGORY, B. W. & COLEMAN, R. L. 2003. Diagnosis of ovarian cancer using decision tree classification of mass spectral data. *BioMed Research International*, 2003, 308-314.
34. WANG, J., YUE, S., YU, X. & WANG, Y. 2017. An efficient data reduction method and its application to cluster analysis. *Neurocomputing*, 238, 234-244.
35. WANG, J., ZUO, Y., MAN, Y.-G., AVITAL, I., STOJADINOVIC, A., LIU, M., YANG, X., VARGHESE, R. S., TADESSE, M. G. & RESSOM, H. W. 2015. Pathway and network approaches for identification of cancer signature markers from omics data. *Journal of Cancer*, 6, 54.
36. WASHAM, C. L., BYRUM, S. D., LEITZEL, K., ALI, S. M., TACKETT, A. J., GADDY, D., SUNDERMANN, S. E., LIPTON, A. & SUVA, L. J. 2013. Identification of PTHrP (12-48) as a plasma biomarker associated with breast cancer bone metastasis. *Cancer Epidemiology and Prevention Biomarkers*, 22, 972-983.
37. WATSON, J. D., BAKER, T., BELL, S., GANN, A., LEVINE, M. & LOSICK, R. 2003. *Molecular biology of the gene*, Pearson/Benjamin Cummings.
38. WONG, J. W. & CAGNEY, G. 2010. An overview of label-free quantitation methods in proteomics by mass spectrometry. *Proteome bioinformatics*, 273-283.
39. YANG, J., ROY, R., JEDINAK, A. & MOSES, M. A. 2015. Mining the human proteome: biomarker discovery for human cancer and metastases. *The Cancer Journal*. ۲۳۶-۲۴۷, ۲۱ ,
40. ZHANG, F. & CHEN, J. Y. A neural network approach to multi-biomarker panel development based on LC/MS/MS proteomics profiles: A case study in breast cancer. *Computer-Based Medical Systems*, 2009. CBMS 2009. 22nd IEEE International Symposium on, 2009. IEEE, 1-6.
41. ZHANG, X., WANG, B., ZHANG, X.-S., LI, Z.-M., GUAN, Z.-Z. & JIANG, W.-Q. 2007. Serum diagnosis of diffuse large B-cell lymphomas and further identification of response to therapy using SELDI-TOF-MS and tree analysis patterning. *BMC cancer*, 7, 235.

پیوست‌ها

پیوست ((الف))

جدول نرم‌افزارهای مورد استفاده در داده‌کاوی داده‌های طیف سنج جرمی

نرم‌افزار	کاربرد	رایگان بودن	آدرس وبسایت
mzMine	شناسایی موجک، برچسب‌گذاری، ایزوتوپ زدایی	بله	http://mzmine.github.io/
موتور جستجوی توالی در بانک‌داده‌ها			
Mascot	تعیین پروتئین‌های موجود در نمونه همچنین شامل empAI برای شناسایی پروتئین‌ها	خیر	http://www.matrixscience.com
Sequest	تعیین پروتئین‌های حاضر در نمونه	خیر	http://fields.scripps.edu/sequest
X!Tandem	تعیین پروتئین‌های حاضر در نمونه	بله	http://www.thegpm.org/tandem
شناسایی پروتئین‌ها بدون برچسب‌گذاری			
emPAI	بدون برچسب که در Mascot هست	بله	http://www.matrixscience.com
PepC	بدون برچسب	بله	http://sashimi.svn sourceforge.net/viewvc/ sashimi/trunk/trans_ proteomic_pipeline/ src/Quantitation/Pepc
APEX	بدون برچسب‌گذاری، کمیت شماری قطعی	بله	http://pfgrc.jcvi.org/index.php/ bioinformatics/apex.html
انتخاب ویژگی / یادگیری ماشین			
WEKA	متدهای متنوع برای انتخاب ویژگی و دسته- بندی	بله	http://www.cs.waikato.ac.nz/ml/ weka
کاربری‌های مختلف			
پکیج‌های مختلف برای R:	xcms MassSpecWavelet Bioconductor packages	بله	http://www.r-project.org/ http://bioconductor.org/packages/ release/bioc/html/xcms.html http://bioconductor.org/packages/ release/bioc/vignettes/ MassSpecWavelet/inst/doc/ MassSpecWavelet.pdf

پیوست ((ب))

فهرست واژگان انگلیسی به فارسی و بالعکس

فارسی	English
ابرمولکول	macromolecule
اثر-انگشت پپتیدی	Peptide-mass fingerprint
انتخاب مویک	Peak picking
آمیلاز بزاق	Salivary amylase
آنتی بادی	Antibody
آنتی ژن	antigen
بردارهای ماشین پشتیبان	Support vector machine
بیز ساده	Nave bayes
پسابش	dehydration synthesis
پلی پپتاید	polypeptide
پنل	Panel
ترجمه	translation
تعاملات بین پروتئین ها	Protein-protein interaction
تغییرات پس از ترجمه پروتئین ها	Post translational modification(PTM)
تنظیم کننده های مثبت	oncogenes
توصیف پروتئین های بیان شده	Protein expression profile
جرم مولکولی	Molecular mass
جنگل تصادفی	Random forest
جهش ژنتیکی	mutation
خودکشی سلولی	apoptosis
خوشه بندی	clustering
دالتون	dalton
درخت تصمیم	Decision tree
دسته بند مبتنی بر قاعده	Rule-based classifiers
دسته بندی	classification
رگ زایی	angiogenesis
رویکرد کامپیوتری در شناخت پروتئین	In silico approach
سرطان خون	Leukemia
سرکوب کننده های تومور	Tumor suppressor
سل لاین	Cell line

Artificial neural network	شبکه های عصبی مصنوعی
Mass spectrometry	طیف سنج جرمی
baseline reduction	کاهش خط مبنا
encode	کدگذاری
Expression protein matrix	ماتریس بیان پروتئین
Tissue	ماهیچه
biofluid	مایعات زیست پذیر
metastasis	متاستاز
Protein localization	مکان یابی پروتئین ها
Isoelectric point	نقطه ایزوالکتریک
Random sampling	نمونه برداری تصادفی
chemiluminescence	نورتابی شیمیایی
ontology	هستی شناسی پروتئین
Cross-validation	وارسی اعتبار

Abstract

Cancer is a very complex disease that occurs at the molecular level in the body. Cancers affected the form of cells, the patterns and the number of biomarkers, especially proteins that are in the tissue and the cell. In this context, proteomics, which examines the patterns and expression of proteins in the tissue and the cell, can be one of the keys to solving the cancer problem. Proteomics includes a variety of approaches and technologies in which the Expression Profiling proteomics approach and mass spectrometry technology can be used in conjunction with the use of high-level analytical methods such as data mining in the classification of cancer and the discovery of relevant biomarkers. In this report, we will discuss the approaches of proteomics and the technologies that used in it, and then focus the studies on expression profiling with the data obtained from the mass spectrometry technology. Subsequently, we will look at data mining and how to use data mining on data obtained from a mass spectrometer for the categorization and discovery of biomarkers, and, finally, we will have a review of past research.

Keys: data mining, mass spectrometry, cancer, biomarkers



Tarbiat Modares University

Faculty of Engineering

Information Technology Engineering Department

Seminar report

Application of proteome data mining in Cancer classification and biomarker discovery

Student

Rasoul norouzi

Supervisor

Dr.amir albadvi

July 2017