

کاربرد روش کاهش بعد چندعاملی (MDR) در شناسایی مدل‌های چند لوکوسی دخیل در رخداد بیماری بهجت

انوشیروان کاظم‌نژاد^۱، پریسا ریاحی^۲، شایان مصطفایی^۳

تاریخ دریافت: ۱۳۹۹/۱۰/۱۸

تاریخ پذیرش: ۱۴۰۰/۰۸/۰۸

چکیده:

رگرسیون لجستیک به دلیل تنگی و وجود متغیرهای رگرسیونی جداکننده‌های کامل و حجم بالایی محاسبات در اثرات متقابل مرتبه بالا، دقت کافی برای تشخیص اثرات اصلی و متقابل بین جهش‌های ژنتیکی را در مراتب خیلی بالا ندارد. الگوریتم کاهش بعد چندعاملی به‌عنوان یک الگوریتم توانمند برای شناسایی اثرات متقابل مراتب بالا در ساختارهای ابربعد محسوب می‌شود. در این تحقیق با استفاده از اطلاعات ۷۴۸ مورد بیمار مبتلا به بیماری بهجت که به مرکز تحقیقات روماتولوژی، بیمارستان شریعتی تهران مراجعه کرده بودند و ۷۷۶ شاهد سالم، برای شناسایی اثرات متقابل بین پلی‌مورفیسم‌های ژن ERAPI دخیل در رخداد بیماری بهجت از الگوریتم کاهش بعد چندعاملی استفاده شده است. محاسبات با استفاده از نرم‌افزار mdr 3.0.2 انجام گرفته است. مدل‌های حاصل از الگوریتم کاهش بعد چندعاملی با دقت متعادل بالای ۰٫۶ دخیل در افزایش ریسک بیماری بهجت تعیین شده‌اند. الگوریتم کاهش بعد چندعاملی توان و سرعت بالایی در محاسبه اثرات متقابل پلی‌مورفیسم‌ها یا جهش‌های ژنتیکی و شناسایی اثرات متقابل مهم و معنی‌دار دارد.

واژه‌های کلیدی: الگوریتم کاهش بعد چندعاملی، بیماری بهجت، اثرات متقابل ژن-ژن

۱ مقدمه

از پیشامد موردبررسی است. به‌طورکلی تحلیل‌های محاسباتی برای اپیستازیس مشکل بوده و دارای محدودیت‌های است زیرا این ساختار بسیار پیچیده است. روش کاهش بعد چندعاملی اولین بار توسط ریچی و همکارانش برای به دست آوردن اثرات متقابل بین ژن‌های دخیل بیماری سرطان پستان غیرارثی معرفی شده است [۱]. اولین مزیت کاهش بعد چندعاملی سرعت عمل آن در تشخیص و شناسایی هم‌زمان چند لوکوس است. مزیت دیگر آن ناپارامتری بودن آن است که این یک تفاوت اصلی در مقایسه با روش‌های پارامتری آماری قدیمی مثل مدل‌های خطی تعمیم‌یافته است. مثلاً در رگرسیون لجستیک، با وارد شدن هر اثر اصلی به مدل، تعداد اثرات متقابل نیز به‌صورت نمایی افزایش می‌یابد که این یک مشکل اساسی است. مزیت سوم آن، بدون پذیره بودن این روش است که این در بیمارهایی مثل سرطان سینه که در آن مد وراثت، نامعلوم و تقریباً پیچیده است، مهم است. مزیت چهارم، حداقل شدن مقدار مثبت کاذب در آزمون‌های چندگانه است با استفاده از اعتبارسنجی متقاطع^۶ در انتخاب بهترین مدل با اثرات متقابل است.

محاسبه اثرات متقابل ژن-ژن و ژن با محیط در بیماری‌های پیچیده چندعاملی و شناسایی اثرات متقابل مهم و معنی‌دار، یکی از مهم‌ترین و چالش‌برانگیزترین زمینه‌های مطالعاتی دانشمندان ژنتیک و آمار است. در محاسبه اثرات متقابل در مراتب بالاتر به‌ویژه در ساختارهای ابربعد، استفاده از روش مرسوم کلاسیک رگرسیون به دلیل یک سری محدودیت‌ها و بار محاسباتی نسبتاً بالا، امکان‌پذیر نیست. رگرسیون لجستیک به دلیل تنگی^۴ و وجود متغیرهای رگرسیونی جداکننده‌های^۵ کامل، دقت کافی برای تشخیص اثرات اصلی و متقابل بین ژن‌ها را در مراتب بالاتر ندارد، چراکه برآورد پارامترها نامعتبر است. در حضور متغیرهای جداکننده نیز در صورتی که متغیر توضیحی بزرگ‌تر از یک مقدار T باشد، یک رابطه بین یک متغیر پیشگو و یکی از متغیرهای توضیحی می‌شود. کاربرد الگوریتم کاهش بعد چندعاملی برای شناسایی اثرات متقابل بین ژن‌ها یا اپیستازیس^۶ در مورد هر فنوتیپی

^۱ استاد آمار زیستی، گروه آمار زیستی، دانشکده علوم پزشکی، دانشگاه تربیت مدرس، تهران، ایران، (نویسنده مسئول: kazem_an@modares.ac.ir)

^۲ دانش‌آموخته رشته آمار زیستی، گروه آمار زیستی، دانشکده علوم پزشکی، دانشگاه تربیت مدرس، تهران، ایران

^۳ استادیار آمار زیستی، گروه آمار زیستی، دانشکده بهداشت، دانشگاه علوم پزشکی کرمانشاه، کرمانشاه، ایران

^۴Sparsity

^۵Separation

^۶Epistasis

^۷Cross-validation, CV

۱۰۲ آماده‌سازی DNA و ژنوتیپ کردن پلی‌مورفیسیم‌ها

پروتکل مطالعه توسط کمیته اخلاقی دانشگاه علوم پزشکی تهران (کمیته اخلاقی ۹۱-۰۴-۴۱-۱۹۳۸۰-۲۹۶۳۷۱) تصویب شده است. نمونه‌های خون استخراج شده از آزمودنی‌ها با استفاده از پروتکل‌ها و تکنیک‌های استاندارد، ژنومیک DNA از آن‌ها استخراج شده است و سپس با استفاده از MGB-TaqMan Allelic Discrimination method^۸ و با استفاده از StepOnePlus Real-Time PCR System^۹ ژنوتایپ‌ها و آلل‌های پلی‌مورفیسیم‌های این ژن به دست می‌آیند. همه پلی‌مورفیسیم‌ها از کنترل کیفیت آماری برخوردار هستند [۵].

۲۰۲ روش کاهش بعد چندعاملی

روش کاهش بعد چندعاملی ابتدا همه ترکیبات دوعاملی را در نظر گرفته و سپس بهترین مدل دوعاملی (مدل با اثرات متقابل مرتبه دو) را با کمترین خطای پیش‌بینی در بین همه مدل‌های دوعاملی انتخاب می‌کند. سپس این فرایند در بین همه مرتبه‌های بالاتر ترکیبات عاملی یا اثرات متقابل تکرار می‌شود و در هر مرحله "بهترین" مدل انتخاب می‌شود. خطای طبقه‌بندی بر روی مجموعه آموزش و خطای پیش‌بینی بر روی مجموعه آزمون محاسبه می‌شود. در چنین مدل‌های به دلیل پیچیدگی و وجود اثرات مقابل زیاد، پدیده بیش‌برازشی بسیار رایج است. اعتبارسنجی متقاطع در مدل بندی تاحدودی قادر است مانع از رخداد این پدیده جلوگیری نماید. در اعتبارسنجی متقاطع مدل‌ها به دنبال مدلی هستیم که در مقایسه با دیگر مدل‌ها، سازگاری^{۱۰} اعتبارسنجی متقاطع آن بیشتر و خطای پیش‌بینی برای آن مدل کمتر باشد در صورتی که چند مدل با چنین شرایطی وجود داشته باشد آنگاه برای انتخاب بهترین مدل، از مدل صرفه‌جو^{۱۱} استفاده می‌شود. با این وجود هرچه تعداد لوسی‌ها افزایش پیدا می‌کند، خطای طبقه‌بندی همواره کاهش می‌یابد و این به دلیل بیش‌برازش مشاهده شده در خطای طبقه‌بندی است. در انتخاب مدل ضروری است تا خطای پیش‌بینی در نظر گرفته شود گرفته شود.

اعتبارسنجی متقاطع، مهم‌ترین مرحله در الگوریتم‌های یادگیری ماشین نظیر کاهش بعد چندعاملی است. چراکه توان پیش‌بینی یک مدل را ارزیابی کرده و از بیش‌برازشی تاحدودی جلوگیری می‌کند. بررسی‌های صورت گرفته نشان داده است که کاهش بازه‌های تعداد اعتبارسنجی متقاطع از ۱۰ به ۵، توان آزمون را کاهش نمی‌دهد و

بیماری بهجت یک بیماری التهابی مزمن سیستمیک است در این بیماری درگیری‌های چشمی که مکرراً اتفاق می‌افتد و کاهش بینایی را در پی دارد، زخم‌های دهانی، زخم‌های تناسلی و اختلالات پوستی شایع‌ترین تظاهرات بالینی هستند. این بیماری در کل دنیا وجود دارد، ولی عمدتاً در کشورهایی که در امتداد جاده ابریشم، از ژاپن تا خاورمیانه و حوضه مدیترانه شایع است. اگرچه علت بیماری هنوز شناخته شده نیست اما همانند بسیاری از بیماری‌های سندروم خودایمنی و خودالتهابی، عوامل ژنتیکی در این بیماری دخیل هستند. بیماری بهجت برای نخستین بار توسط هولوسی بهجت^۸ توصیف شده است [۲]. علاوه بر اثرات چشم و پوست، افت دهان و ناحیه تناسلی معیار اصلی تشخیص بیماری بهجت هستند. اگرچه بهجت به‌طور خاص در آمریکا (۵/۲ در ۱۰۰۰۰۰) و در اروپا (۲/۴ در ۱۰۰۰۰۰) شایع نیست. با توجه به میزان شیوع بهجت در کشورهای مجاور جاده ابریشم از جمله ایران، که شیوع آن ۸۰ نفر در هر ۱۰۰۰۰۰ است باعث شده است که مدیریت آن، یک مشکل برای متخصصان روماتیسم ایران و آسیا محسوب شود [۳، ۴]. لذا این تحقیق با هدف شناسایی اثرات متقابل پلی‌مورفیسیم‌های ژن ERAP1 دخیل در رخداد بیماری بهجت با کاربرد روش کاهش بعد چندعاملی است.

۲ روش کار

این پژوهش یک مطالعه درون رایانه‌ای است در این مطالعه از اطلاعات ۷۴۸ مورد بیمار مبتلا به بیماری بهجت و ۷۷۶ شاهد سالم استفاده شده است. اطلاعات افراد مورد از بیماران بهجتی مراجعه‌کننده به مرکز تحقیقات روماتولوژی، بیمارستان شریعتی تهران با روش نمونه‌گیری تصادفی ساده انتخاب شده و مورد مطالعه قرار گرفته‌اند. تمام تشخیص‌ها با معیارهای بین‌المللی بیماری بهجت (ICBD) تأیید شده است. گروه شاهد شامل ۷۷۶ فرد سالم، بدون سابقه خانوادگی و یا علائم بالینی هر نوع اختلالات روماتیسمی یا سایر اختلالات خودایمنی و همسان شده به لحاظ سن و جنسیت و نژاد با گروه مورد بیمار هستند [۵]. همه محاسبات با استفاده از نرم‌افزار mdr 3.0.2 انجام گرفته است. میزان خطای نوع اول به منظور سطح معنی‌داری آماری ۰/۰۵ در نظر گرفته شده است.

⁸Hulusi Behcet

⁹Applied Biosystems, Foster City, CA, USA

¹⁰Consistency

¹¹Statistical parsimony

در بین همه مرتبه‌های بالاتر ترکیبات عاملی تکرار می‌شود و در هر مرحله «بهترین» مدل انتخاب می‌شود. در بین همه «مدل‌های نهایی» مدلی که خطای پیش‌بینی را حداقل و سازگاری اعتبارسنجی متقابل را حداکثر می‌کند، انتخاب می‌شود. اگر چندین مدل، خطای پیش‌بینی و سازگاری اعتبارسنجی متقابل (تعداد دفعاتی که مدل به دست آمده از بین ۹۱۰ داده‌ها یکسان باشند) یکسانی داشته باشد، مدل با کوچک‌ترین تعداد پلی‌مورفیسم‌ها انتخاب می‌شود [۷].

۳ نتایج

اطلاعات دموگرافیک شرکت‌کنندگان میانگین سنی بیماران بهجتی 40.26 ± 10.88 و در گروه کنترل سالم 38.88 ± 11.54 سال بود (با p -مقدار برابر با ۰/۰۷۶). از بین ۷۴۸ بیمار بهجتی، ۴۴۸ نفر (۵۹/۸٪) و از ۷۷۶ نفر فرد کنترل سالم، ۴۷۶ نفر (۶۱/۳٪) مرد بودند. همچنین جدول ۱ فراوانی الل-های ژن ERAP1 را نشان می‌دهد.

با توجه به جدول ۲، پلی‌مورفیسم rs27044 به‌عنوان اثر اصلی با بیشترین مقدار دقت متعادل 0.5432 در مقایسه با ۱۰ پلی‌مورفیسم دیگر شناسایی شد. این پلی‌مورفیسم دارای بیشترین سازگاری اعتبارسنجی را در بین مدل‌های تک لوکوسی داشته است (۷۱/۰). مدل دولوکوسی به دست آمده از این روش با دقت متعادل ۵۶/۴۱، و همچنین سازگاری اعتبارسنجی متقاطع ۶۱/۰ به‌عنوان بهترین مدل دولوکوسی گزارش می‌شود. مدل سه لوکوسی گزارش شده شامل پلی‌مورفیسم‌های rs30187، rs469876 و rs13167972 می‌باشد. که دقت متعادل برای این مدل ۵۸/۴۴ است. مدل چهار لوکوسی $rs27434 \times rs28096 \times rs469876 \times rs13167972$ با دقت متعادل ۶۰/۷۸ به‌عنوان بهترین مدل در مرتبه چهار گزارش می‌شود. در اثرات متقابل مرتبه پنج پلی‌مورفیسم‌های rs27044، rs26653، rs469876، rs28096 و rs13167972 با دقت متعادل ۶۲/۳ و میزان از سازگاری اعتبارسنجی متقاطع ۳۱/۰ به‌عنوان بهترین مدل پنج لوکوسی از بین مدل‌های دیگر انتخاب شده‌اند.

در مرتبه شش، پلی‌مورفیسم‌های rs106547، rs30187، rs26653، rs469876، rs28096 و rs13167972 با دقت متعادل ۶۴/۱۴ و میزان سازگاری اعتبارسنجی متقاطع ۲۱/۰ بهترین مدل شش لوکوسی بین پلی‌مورفیسم‌های ژن ERAP1 را به دست دادند. پلی‌مورفیسم‌های rs1065407، rs30187، rs2287987، rs26653، rs469876 و rs28096 با دقت متعادل 0.6631 و سازگاری اعتبارسنجی متقاطع ۶۱/۰ به‌عنوان بهترین مدل هفت لوکوسی معرفی شده است. مدل هشت لوکوسی $rs1065407 \times rs2287987 \times rs30187 \times rs27044 \times rs26653 \times rs469876$

زمان محاسباتی را نصف کرده و برای تعداد مشاهدات کم پیشنهاد می‌شود [۶].

روش اجرای الگوریتم کاهش بعد چندعاملی در شکل ۱ نشان داده شده است که در آن هر سلول چندعاملی در فضای n بعدی به‌عنوان «ریسک بالا» یا «ریسک پایین» نام‌گذاری می‌شود و خطای پیش‌بینی هر مدل برآورد شده است. مراحل اجرای کاهش بعد چندعاملی:

۱. در گام اول، یک مجموعه n تایی ورینت‌های ژنتیکی گسسته از ترکیب تمام عوامل انتخاب می‌کنیم.

۲. در گام دوم، تمام n عامل و کلاس‌های چندعاملی یا سلول‌ها در فضای n بعدی نمایش داده می‌شوند. سپس نسبت تعداد موارد به تعداد شاهدها در درون هر کلاس چندعاملی برآورد می‌شود.

۳. هر سلول در فضای n بعدی به دو دسته «ریسک بالا»، اگر نسبت موردها به شاهدها برابر باشد یا بیشتر از مقدار آستانه تعیین شده (به‌عنوان مثال ۰/۱) باشد، یا «ریسک پایین» اگر نسبت محاسبه شده از مقدار آستانه بزرگ‌تر نباشد منتسب می‌شود. بدین طریق، یک مدل برای موردها و شاهدها با ادغام خانه‌های با ریسک بالا به یک گروه و خانه‌های با ریسک پایین به یک گروه دیگر ایجاد می‌شود. بدین طریق ما مدل Π بعدی را به مدل تک‌بعدی تبدیل می‌کنیم (یعنی اینکه یک متغیر با دو کلاس چندعاملی داریم ریسک بالا، ریسک پایین).

۴. در مرحله چهارم، خطای پیش‌بینی شده هر مدل به‌وسیله اعتبارسنجی متقابل ۱۰-تایی برآورد می‌شود. به این صورت که داده‌ها به‌صورت تصادفی به ۱۰ مجموعه مساوی تقسیم می‌کنیم. مدل کاهش بعد چندعاملی را برای هر ۹/۱۰ داده‌ها اجرا می‌کنیم و سپس به‌منظور پیش‌بینی وضعیت بیماری (متغیر پاسخ) ۱/۱۰ داده‌ها از مطالعه کنار گذاشته می‌شوند. نسبت داده‌هایی که موجب پیش‌بینی نادرست شده‌اند، برآورد خطای پیش‌بینی است. برای کاهش احتمال برآوردهای ضعیف خطای پیش‌بینی که به علت تقسیم تصادفی مجموعه داده‌ها است، اعتبارسنجی متقابل ۱۰-برابری، را ۱۰ مرتبه تکرار می‌کنیم و سپس متوسط مقدار خطاهای پیش‌بینی را محاسبه می‌کنیم.

روش کاهش بعد چندعاملی در ابتدا همه ترکیبات چندعاملی را در نظر گرفته و «بهترین» مدل چندعاملی را با کمترین خطای پیش‌بینی در بین همه مدل‌های هر مرتبه انتخاب می‌کند. سپس این فرایند

اعتبارسنجی متقاطع ۵/۱۰ و مدل نه لوکوسی $rs1065407 \times rs28096 \times rs469876 \times rs26653 \times rs27044 \times rs2287987$ با دقت متعادل ۰/۶۷۳۱ و سازگاری $rs13167972 \times rs30187 \times rs1748078 \times rs13167972$ با دقت متعادل ۰/۶۷۸۹ و دقت متعادل ۴/۱۰ به‌عنوان بهترین مدل‌ها به‌دست‌آمده‌اند. مدل‌های ده، یازده و دوازده لوکوسی آمده در جدول به ترتیب با دقت متعادل

۰/۶۸۲۱، ۰/۶۸۴۸ و ۰/۶۸۴۸ میزان سازگاری اعتبارسنجی متقاطع ۵/۱۰، ۱۰/۱۰ و ۱۰/۱۰ به‌عنوان بهترین مدل‌ها انتخاب می‌شوند (جدول ۲). ضرایب جینی اندازه و جهت ارتباط بین پلی‌مورفیسم‌ها در قالب یک نمودار دایره‌ای در شکل ۲ نمایش می‌دهد. پلی‌مورفیسم‌های $rs2287987$ و $rs27044$ بزرگ‌ترین ضرایب را در بین پلی‌مورفیسم‌های دیگر داشته‌اند.

جدول ۱. فراوانی الل‌های ژن ERAP1

نسبت بخت (فاصله اطمینان)	p-مقدار	فراوانی الل کوچک (%)		الل‌ها	مکان پلی‌مورفیسم در کروموزوم پنج	SNP
		شاهد	مورد			
۱۲۰ (۱/۰۳, ۱/۳۹)	۰/۱۸	۳۲/۵	۳۶/۶	G < T	۳۷۹,۷۷۶,۹۶	rs1065407
۰/۹۷ (۰/۸۳, ۱/۱۴)	۰/۷۴	۲۹/۱	۲۸/۵	G < C	۱۴۸,۷۸۳,۹۶	rs27044
۱۲۵ (۱/۰۰, ۱/۵۶)	۰/۵۲	۱۰/۳	۱۲/۶	T < C	۱۶۲,۷۸۳,۹۶	rs17482078
۱۲۷ (۱/۰۱, ۱/۵۹)	۰/۳۹	۱۰/۱	۱۲/۵	T < C	۵۰۶,۷۸۶,۹۶	rs10050860
۱۰۲ (۰/۸۸, ۱/۱۸)	۰/۸۲	۳۹/۷	۴۰/۱	T < C	۶۲۷,۷۸۸,۹۶	rs30187
۱۲۷ (۱/۰۱, ۱/۵۹)	۰/۴۰	۱۰/۲	۱۲/۵	C < T	۸۳۲,۷۹۳,۹۶	rs2287987
۱۰۰ (۰/۷۹, ۱/۲۶)	۰/۸۸	۹/۸	۹/۸	T < C	۸۴۰,۷۹۳,۹۶	rs27895
۰/۸۵ (۰/۷۱, ۱/۰۱)	۰/۵۹	۲۲/۸	۲۰/۱	C < T	۱۳۳,۷۹۵,۹۶	rs26618
۱۰۲ (۰/۸۹, ۱/۱۸)	۰/۷۵	۳۹/۷	۴۰/۲	C < G	۵۴۷,۸۰۳,۹۶	rs26653
۰/۸۱ (۰/۵۰, ۱/۳۲)	۰/۴۰	۲/۴	۱/۸	T < C	۷۶۱,۸۰۳,۹۶	rs3734016
۰/۹۸ (۰/۷۷, ۱/۲۵)	۰/۸۸	۹/۸	۹/۸	A < G	۸۹۲,۸۰۳,۹۶	rs72773968

است. از مزایای روش کاهش بعد چندعاملی می‌توان به سرعت عمل آن در تشخیص و شناسایی هم‌زمان چند لوکوس نام برد. محاسبه اثرات متقابل در حضور تعداد پلی‌مورفیسم‌های بالا و همچنین حجم نمونه بالا نیازمند یک الگوریتم توانمند با قدرت تشخیص بالاست که در عین حال سرعت پردازش بالایی داشته باشد. مزیت دیگر آن ناپارامتری بودن و بدون مدل بودن آن است که این یک تفاوت اصلی در مقایسه با روش‌های پارامتری آماری قدیمی مثل مدل‌های خطی تعمیم‌یافته است. این مدل علاوه بر مزیت‌های ذکر شده یک مزیت دیگر نیز دارد که این مدل را به‌طور برجسته‌ای از دیگر مدل‌ها متمایز می‌کند، حداقل شدن مقدار مثبت کاذب در آزمون‌های چندگانه است که این به دلیل استراتژی اعتبارسنجی متقابل در انتخاب بهترین مدل است [۱].

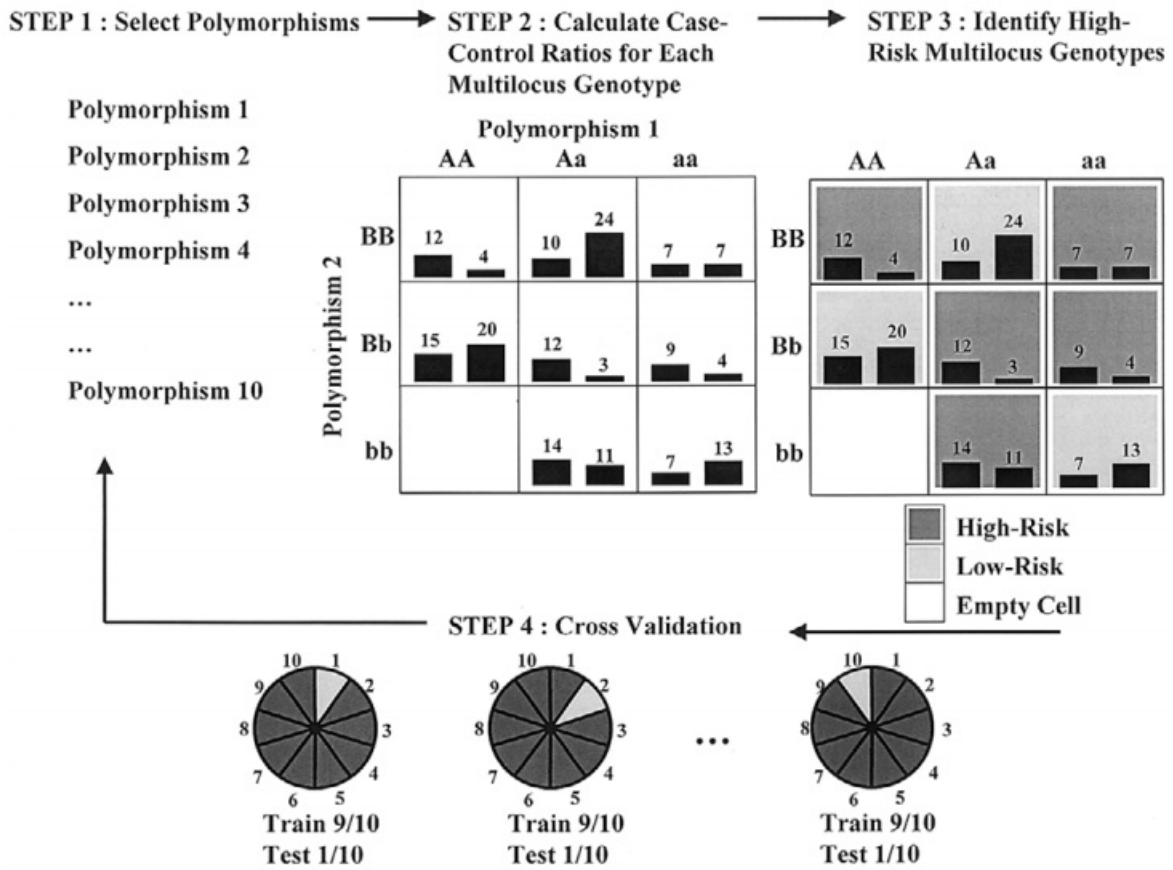
در این نمودار ضرایب اثرات متقابل بین هر جفت پلی‌مورفیسم بر روی خطی که این دو را به هم وصل می‌کند نوشته شده است. همچنین علامت ضریب نشان از جهت اثر دارد، به‌طوری‌که ضرایب با علامت منفی دارای اثر کاهنده و مثبت دارای اثر فزاینده‌اند. اندازه اثر اصلی هر پلی‌مورفیسم در کنار اسم آن گزارش شده است. رنگ نارنجی نشانگر رابطه فزاینده با درجه کم، رنگ طلایی نشانگر نبود و یا رابطه فزاینده بسیار کم، رنگ سبز نشانگر رابطه کاهنده با درجه کم و رنگ آبی نشانگر رابطه کاهنده با درجه بیشتر از گروه قبلی است.

۴ بحث و نتیجه‌گیری

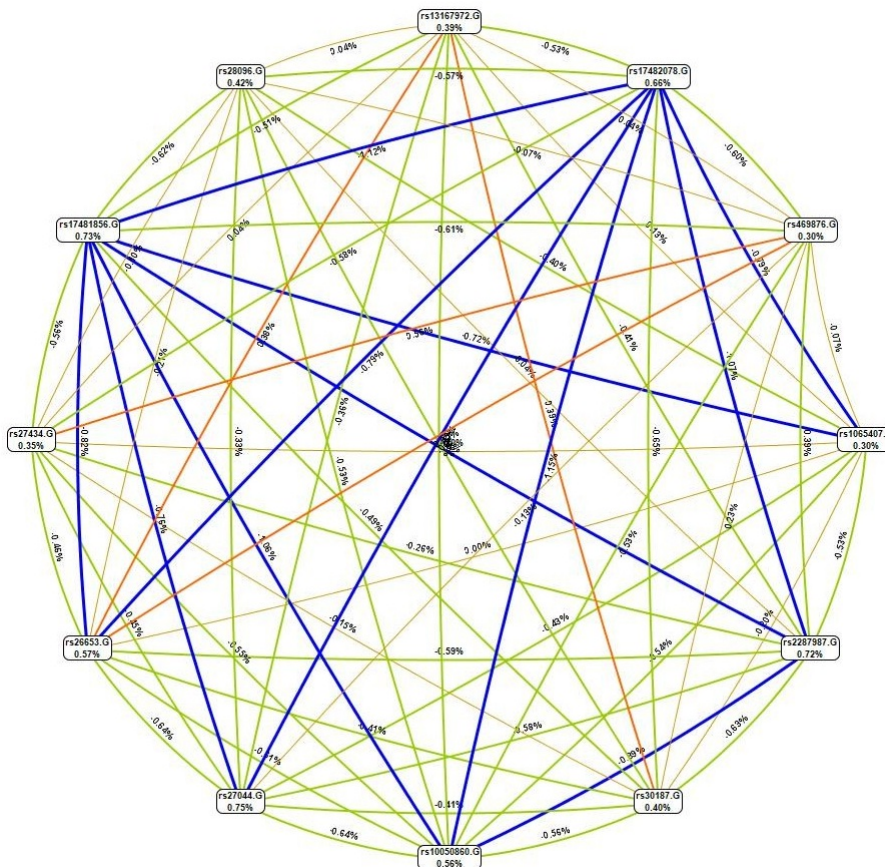
روش کاهش بعد چندعاملی یکی از ابزارهای قوی معرفی شده برای محاسبه اپیستازیس (اثر متقابل ژن-ژن) در حوزه علم آمار ژنتیک

جدول ۲. نتایج الگوریتم کاهش بعد چندعاملی برای پلی مورفیزم‌های ژن ERAP1

نتایج اعتبار سنجی متقاطع							مدل
ویژگی	حساسیت	دقت متعادل	CV	ویژگی	حساسیت	دقت متعادل	
۰٫۴۲۴۲	۰٫۶۶۲۲	۰٫۵۴۳۲	۷٫۱	۰٫۴۴۱۱	۰٫۵۴۸۵	۰٫۴۹۴۸	rs27044
۰٫۶۵۳۲	۰٫۴۷۴۹	۰٫۵۶۴۱	۶٫۱	۰٫۵۵۸۹	۰٫۴۰۸۰	۰٫۴۸۳۵	rs27434 × rs469876
۰٫۷۲۷۳	۰٫۴۴۱۵	۰٫۵۸۴۴	۶٫۱	۰٫۶۱۶۲	۰٫۳۸۸۰	۰٫۵۰۲۱	rs30187 × rs469876 × rs13167972
۰٫۷۴۰۷	۰٫۴۷۴۹	۰٫۶۰۷۸	۵٫۱	۰٫۶۱۹۵	۰٫۴۳۱۴	۰٫۵۲۵۵	rs27434 × rs28096 × rs469876 × rs13167972
۰٫۷۶۷۷	۰٫۴۷۸۲	۰٫۶۲۳۰	۳٫۱	۰٫۵۸۹۲	۰٫۳۶۷۹	۰٫۴۷۸۶	rs27044 × rs26653 × rs469876 × rs28096 × rs13167972
۰٫۷۴۰۷	۰٫۵۴۱۸	۰٫۶۴۱۴	۲٫۱	۰٫۵۲۵۳	۰٫۴۰۱۳	۰٫۴۶۳۳	rs1065407 × rs30187 × rs26653 × rs469876 × rs28096 × rs13167972
۰٫۷۶۴۳	۰٫۵۶۱۹	۰٫۶۶۳۱	۶٫۱	۰٫۵۴۸۸	۰٫۴۲۱۴	۰٫۴۸۵۱	rs1065407 × rs30187 × rs2287987 × rs26653 × rs469876 × rs28096 × rs13167972
۰٫۷۵۷۶	۰٫۵۸۸۶	۰٫۶۷۳۱	۵٫۱	۰٫۵۵۵۶	۰٫۴۳۸۱	۰٫۴۹۶۸	rs1065407 × rs2287987 × rs30187 × rs27044 × rs26653 × rs469876 × rs28096 × rs13167972
۰٫۷۵۴۲	۰٫۶۰۵۴	۰٫۶۷۹۸	۴٫۱	۰٫۵۶۵۷	۰٫۴۴۴۸	۰٫۵۰۵۲	rs1065407 × rs2287987 × rs27044 × rs26653 × rs469876 × rs28096 × rs13167972 × rs1748078 × rs30187
۰٫۷۵۴۲	۰٫۶۱۲۰	۰٫۶۸۲۱	۵٫۱	۰٫۵۵۵۶	۰٫۴۵۱۵	۰٫۵۰۳۵	rs1065407 × rs2287987 × rs30187 × rs27044 × rs26653 × rs27434 × rs469876 × rs28096 × rs13167972 × rs17482087
۰٫۷۵۷۶	۰٫۶۱۲۰	۰٫۶۸۴۸	۱۰٫۱	۰٫۵۴۵۵	۰٫۴۶۴۹	۰٫۵۰۵۲	rs1065407 × rs2287987 × rs30187 × rs10050860 × rs27044 × rs26653 × rs27434 × rs469876 × rs28096 × rs13167972 × rs17482078
۰٫۷۵۷۶	۰٫۶۱۲۰	۰٫۶۸۴۸	۱۰٫۱	۰٫۵۴۵۵	۰٫۴۶۴۹	۰٫۵۰۵۲	rs1065407 × rs2287987 × rs30187 × rs10050860 × rs27044 × rs26653 × rs27434 × rs469876 × rs28096 × rs13167972 × rs17482078 × rs17481856



شکل ۱. مراحل انجام الگوریتم کاهش بعد چندعاملی



شکل ۲. نمودار دایره‌ای نتایج MDR

کوریت‌های مهم اجرا شود، تعیین مدل نهایی سخت می‌شود. به‌عنوان مثال، اگر در تحلیل کاهش بعد چندعاملی معلوم شود که مدل نهایی، شامل ۴ پلی‌مورفیسم است، کاملاً مشخص نمی‌شود که آیا مدل نهایی یک مدل با اثر متقابل ۴ طرفه است، یا دارای دو اثر متقابل دوطرفه جداگانه، و یا دو اثر اصلی و یک اثر متقابل دوطرفه است [۷]. به عبارتی در مدل کاهش بعد چندعاملی تعیین اینکه مدل نهایی به‌دست‌آمده فقط دارای یک اثر متقابل است یا این پلی‌مورفیسم‌ها اثرات متقابل با مرتبه کوچک‌تری نیز دارند یا خیر. همچنین تفسیر نتایج در مواقعی که تعداد متغیرها تقریباً زیاد و حجم نمونه کم است سخت بوده و در این صورت توان پیش‌بینی مدل کم است (در حجم نمونه کوچک، نتایج مثبت کاذب و منفی کاذب به دست می‌آید) [۱]. در این مطالعه، مدل کاهش بعد چندعاملی در محاسبه و تشخیص مهم‌ترین اثرات متقابل بین پلی‌مورفیسم‌های ژن ERAP1 مرتبط با رخداد بیماری بهجت از دقت و سرعت خوبی برخوردار بود. این مدل علیرغم کاستی‌ها و معایب در حضور داده‌های گمشده و ناهمگونی‌های ژنتیکی، از توان بالایی در محاسبه و به دست آوردن مدل‌های چند لوکوسی ژنتیکی دارد.

در کنار مزیت‌ها و توانمندی‌های مدل کاهش بعد چندعاملی، این روش مشابه سایر روش‌های دیگر از معایب و کاستی‌هایی همراه است. یکی از عمده مشکلات روش کاهش بعد چندعاملی خالی بودن خانه‌های جدول توافقی ساخته‌شده در مدل‌هایی با ابعاد بالاست که در این صورت قادر نیست که این خانه‌ها را به‌عنوان خانه با ریسک بالا (فزاینده) و یا با ریسک پایین (کاهنده) نام‌گذاری کرد که این ممکن است در محاسبه درست و دقیق اثرات متقابل ابعاد بالا دچار مشکل شود. مشکل دیگر آن است که نسبت موردها به شاهد‌ها نزدیک به یک باشد که در این صورت نام‌گذاری خانه‌های جدول کار مشکلی خواهد بود [۸]. از عمده مشکلات این مدل دودویی بودن متغیر پاسخ است و نمی‌توان مقیاس‌های دیگری برای متغیر پاسخ در نظر گرفت. همچنین این مدل قادر نیست تا اثر متغیرهای مخدوش‌گر و یا هیچ کوریتی را تعدیل کند [۹]. در هنگام استفاده از روش کاهش بعد چندعاملی باید در نظر داشت که اگر تعداد عامل‌ها، بیش از ۱۵ باشد، محاسبات پیچیده‌تر شده و زمان محاسباتی بیشتر می‌شود و توان مدل در حضور داده‌های گمشده، خطای ژنوتیپی، ناهمگونی ژنتیکی و فنوکپی کم است [۱۰]. اگر کاهش بعد چندعاملی در حضور اثرات اصلی یا

مراجع

- [1] Ritchie, M.D., Hahn, L.W., Roodi, N., Bailey, L.R., Dupont, W.D., Parl, F.F. and Moore, J.H. (2001). Multifactor-dimensionality reduction reveals high-order interactions among estrogen-metabolism genes in sporadic breast cancer. *The American Journal of Human Genetics*, **69(1)**, 138–147.
- [2] Mahr, A., Belarbi, L., Wechsler, B., Jeanneret, D., Dhote, R., Fain, O., Lhote, F., Ramanoelina, J., Coste, J., and Guillevin, L. (2008). Population-based prevalence study of Behçet's disease: differences by ethnic origin and low variation by age at immigration. *Arthritis & Rheumatism: Official Journal of the American College of Rheumatology*, **58(12)**, 3951–3959.
- [3] Davatchi, F., Shahram, F., Chams-Davatchi, C., Shams, H., Nadji, A., Akhlaghi, M., Faezi, T., Ghodsi, Z., Larimi, R., Ashofteh, F., and Sadeghi Abdollahi, B. (2010). Behçet's disease in Iran: analysis of 6500 cases. *International journal of rheumatic diseases*, **13(4)**, 367–373.
- [4] International Team for the Revision of the International Criteria for Behçet's Disease (ITR-ICBD). (2014). The International Criteria for Behçet's Disease (ICBD): a collaborative study of 27 countries on the sensitivity and specificity of the new criteria, *Journal of the European Academy of Dermatology and Venereology*, **28(3)**, 338–347.
- [5] Mahmoudi, M., Jamshidi, A.R., Amirzargar, A.A., Farhadi, E., Nourijelyani, K., Fallahi, S., Oraei, M., Noori, S., and Nicknam, M.H. (2012). Association between endoplasmic reticulum aminopeptidase-1 (ERAP-1) and susceptibility to ankylosing spondylitis in Iran, *Iranian Journal of Allergy, Asthma and Immunology*, **11**, 294–300.
- [6] Mostafaei, S. and Riahi, P. (2020). Interaction Effects of Plasma Vitamins A, E, D, B9, and B12 and Tobacco Exposure in Urothelial Bladder Cancer: A Multifactor Dimensionality Reduction Analysis. *Nutrition and Cancer*, 1-2.

- [7] Coffey, C.S., Hebert, P. R., Ritchie, M. D., Krumholz, H. M., Gaziano, J. M., Ridker, P. M., Brown, N. J., Vaughan, D. E., and Moore, J. H. (2004). An application of conditional logistic regression and multifactor dimensionality reduction for detecting gene-gene interactions on risk of myocardial infarction: the importance of model validation. *BMC bioinformatics*, **5(1)**, 49.
- [8] Park, M.Y. and Hastie, T. (2008). Penalized logistic regression for detecting gene interactions. *Biostatistics*, **9(1)**, 30–50.
- [9] Riahi, P., Kazemnejad, A., Mostafaei, S., Meguro, A., Mizuki, N., Ashraf-Ganjouei, A., Javinani, A., Faezi, S.T., Shahram, F., and Mahmoudi, M. (2020). ERAP1 polymorphisms interactions and their association with Behçet's disease susceptibility: Application of Model-Based Multifactor Dimension Reduction Algorithm (MB-MDR), *PloS one*, **15(2)**, e0227997.
- [10] Calle, M.L., Urrea Gales, V., Malats i Riera, N., Van Steen, K. (2008) MB-MDR: model-based multifactor dimensionality reduction for detecting interactions in high-dimensional genomic data. *Technical Report 24*, Department of Systems Biology, Universitat de Vic, Vic,: Spain; available at <http://repositori.uvic.cat/xmlui/handle/10854/408>.

Application of multifactor dimensionality reduction (MDR) algorithm in detecting the n-locus models related to Behcet's disease

A. Kazemnejad ¹ , P. Riyahi ² , S. Mostafaei ³

Abstract:

Due to the sparsity and separation and a large amount of calculations in high-order interactions, Logistic regression is not accurate enough to detect the main and interaction effects between genetic markers at very high orders. The multifactorial dimension reduction algorithm is considered a powerful algorithm for identifying high-order interactions in high dimensional structures. In this study, information of 748 patients with Behcet's disease who were referred to the Rheumatology Research Center, Shariati Hospital, Tehran, and 776 healthy controls were used to identify the interaction effects between ERAP1 gene polymorphisms involved in the occurrence of Behcet's disease using the multifactor dimensionality reduction algorithm. Data analysis was performed using MDR 3.0.2 software. The models obtained from the multifactorial dimensional reduction algorithm with balanced accuracy above 0.6 have been determined to increase the risk of Behcet's disease. The multi-factor reduction algorithm has high power and speed in calculating the interaction effects of polymorphisms or genetic mutations and identifying important interactions.

Keywords: Multifactor dimensionality reduction algorithm, Behcet's disease, Gene-gen interaction.

¹Professor of Biostatistics, department of Biostatistics, faculty of Medical Sciences, Tarbiat modares University, Tehran, Iran.

Corresponding author, kazem_an@modares.ac.ir

²Graduated student of Biostatistics, department of Biostatistics, faculty of Medical Sciences, Tarbiat modares University, Tehran, Iran.

³Assistant professor of Biostatistics, department of Biostatistics, faculty of Health, Kermanshah University of Medical Sciences, Krmanshah, Iran.