



ارائه مدل یادگیری عمیق CNN با هدف بهبود دقت پیش بینی خطر بیماری پرفشاری خون بر مبنای SNP

سیدعلی لاجوردی^۱، مهرداد کارگری^۲، مریم السادات دانشپور^۳، مهدی اکبرزاده^۴

۱-دانشگاه تربیت مدرس - دانشکده صنایع و سیستم ها

۲-دانشگاه تربیت مدرس - دانشکده صنایع و سیستم ها

۳-دانشگاه علوم پزشکی شهید بهشتی - پژوهشکده علوم غدد درون ریز و متابولیسم

۴-دانشگاه علوم پزشکی شهید بهشتی - پژوهشکده علوم غدد درون ریز و متابولیسم

چکیده

فشارخون عامل قابل اصلاح برای بیماری‌های قلبی عروقی مانند بیماری اسکمییک قلب است که از دلایلی اصلی مرگ در سراسر جهان است که با نام کشنده خاموش معرفی می شود. از این رو، روش هایی که می توانند به طور دقیق خطر ابتلا به فشار خون بالا را در سنین پایین آشکار کنند، ضروری هستند. مدل‌های معمول پیش بینی خطر بیماری در درجه اول مبتنی بر عامل‌های شیوه زندگی است. اخیراً، در نظر گرفتن خطر عامل‌های ژنتیکی از جمله SNP‌های مرتبط با بیماری همراه با مدل سازی خطر، دقت پیش‌بینی فردی بیماری‌ها را بهبود بخ شیده است. SNP یک تغییر ژنتیکی کوچک در DNA است و رایج‌ترین نوع تنوع ژنتیکی است که در میان افراد رخ می‌دهد. مدل‌های مختلفی که از نشانگرهای ژنومی جهت پیش‌بینی بیماری فشار خون استفاده می‌کنند شامل چهار رویکرد آماری، فرا تحلیل، یادگیری ماشین و مدل بالینی هستند. مهم‌ترین مطالعه در این مدل‌ها تعداد بالای SNP‌های ورودی و ارتباط این SNP‌ها با یکدیگر است. در مطالعه حاضر از روش مبتنی بر یادگیری عمیق با رویکرد CNN برای استفاده از SNP‌های متعدد فشار خون در یک مطالعه کوهورت طولی استفاده و برای مقایسه نیز PRS با استفاده از دو نرم افزار plink و gcta64 محاسبه شده است. در ابتدا داده‌های ژنومی تبدیل به تصویر شده و وارد مدل CNN می‌شود که لایه‌های آن شامل لایه پیش‌پردازش، لایه ادغام و لایه کاملاً متصل و در نهایت لایه خروجی می‌باشد. داده‌های تحقیق شامل سه بخش داده ژنومی، سن و داده‌های طولی فنوتیپ فشارخون مبتنی بر مطالعه ژنتیک قلبی - متابولیک تهران است. برای مقایسه کارایی مدل از AUC استفاده شده است که مدل CNN با مقدار AUC 0.87، عملکرد بهتری نسبت به PRS و همچنین آخرین مدل‌های ارائه شده در ادبیات موضوع را نشان می‌دهد.

کلمات کلیدی: نمره خطر پلی ژنتیک، یادگیری ماشین، بیماری فشار خون، مارکر ژنتیکی، شبکه عصبی پیچشی



مقدمه

کشنده خاموش^۱ که عامل اصلی در بیماری‌های قلبی عروقی مانند بیماری اسکیمیک قلب و سکته مغزی است، فشارخون است که از علل مهم مرگ و بیماری در سراسر جهان می‌باشد. فشارخون و پرفشاری خون؛ هر دو از صفات چندژنی به حساب می‌آیند، که تعداد زیادی از مسیرهای متابولیکی را شامل می‌شوند (Martinez-Ríos et al. 2021; K. T. Mills, Stefanescu, and He 2020; Zhou et al. 2017). از این رو، روش‌هایی که می‌توانند به طور دقیق خطر ابتلا به فشار خون بالا را در سنین پایین آشکار کنند، ضروری هستند.

از طرف دیگر یکی از مهم‌ترین اهداف سلامت شخصی، بهبود دقت پیش‌بینی بیماری با برر سی انواع ژنتیکی است. بنابراین، روش‌های مختلف بالینی و ژنتیکی برای پیش‌بینی بیماری استفاده می‌شود. با این حال، مسئله مهم با این روش‌ها، تعداد بالای متغیرهای ورودی به عنوان نشانگرهای ژنتیکی با تعداد کمی نمونه است. روشی که می‌تواند بر این مشکل غلبه کند، یادگیری ماشینی است.

ژنوم انسان و تفاوت ژنومی

ژنوم انسان شامل تمام اطلاعات مورد نیاز برای حیات انسان است، که از ۳ میلیارد جفت باز در ۲۳ جفت کروموزوم به نام DNA تشکیل شده است. این مجموعه حاوی تمام اطلاعات لازم برای تعیین صفاتی است که فرد را می‌سازد. کدگذاری اختصاصی درون ژن‌ها از طریق چگونگی ترکیب چهار باز خاص به نام‌های آدنین، تیمین، گوانین و سیتوزین (A, T, G, C) انجام شده است.

در هنگام مقایسه ژنوم دو فرد، تفاوت‌های بسیاری (به لحاظ تعداد) دیده می‌شود که این تفاوت‌ها حاصل تغییرات ژنومی مانند چند پلی مورفیسم تک نوکلئوتیدی^۲ و تغییرات بزرگ‌تر مانند حذف^۳، جایگزینی^۴ و تغییر در تعداد کپی‌های ژنی^۵ هستند. هر کدام از این تغییرات همراه با عوامل محیطی می‌توانند منجر به بروز یک فنوتیپ متفاوت در فرد شوند.

SNP یک تغییر ژنتیکی کوچک در DNA است و رایج‌ترین نوع تنوع ژنتیکی است که در میان افراد رخ می‌دهد. SNP زمانی رخ می‌دهد که باز یک نوکلئوتید به عنوان مثال نوکلئوتید دارای باز A جایگزین یکی دیگر از سه باز نوکلئوتید (مثلاً G, C, یا T) می‌شود. البته این تغییرات می‌توانند بدون اثر (بدون تغییر در فنوتیپ)، اثر غیر بیماری‌زا (تغییر در فنوتیپ)، اثر ضد سلامت (دیابت، سرطان، بیماری‌های قلبی، بیماری‌های تنگن و هموفیلی) و اثر نهفته (فرضا تغییرات در ژنوم مناطق تنظیم‌کننده ژن به خودی خود مضر نیست، و تنها تحت شرایط خاصی مانند استعداد ابتلا به سرطان ریه و تغییر پاسخ فرد به یک درمان آشکار می‌شود) باشد (Filshtein et al. 2019).

¹ Silent killer

² Single nucleotide variation

³ Deletion

⁴ Insertion

⁵ Copy number variation



مطالعات گسترده ارتباط ژنومی⁶ (GWAS) و امتیاز خطر چند ژنی⁷ (PRS)

جهت برر سی ارتباط بین ژنوتیپ و فنوتیپ از برر سی توارث⁸ در خانواده از آزمون های لینکاژ یا پیو ستگی ژنتیکی⁹ استفاده می شد. این روش ها در بررسی بیماری هایی که یک ژن عامل بروز آنها بود، بسیار موفق بودند. اما این یافته ها در بیماری های چندعاملی پیچیده تعمیم پذیر نمی باشند. در مطالعات ارتباط ژنوم گسترده کشف نشانگرهای زیستی، شناسایی افراد دارای ژن بیماری های خاص و شناسایی نوع واکنش به درمان به طور وسیع تری صورت گرفته است (M. C. Mills and Rahal 2019).

این مطالعات اغلب از نوع مورد شاهدی انجام می گردد؛ بدین معنا که گروهی از بیماران در مقابل افراد سالم انتخاب می شوند و تمامی این افراد برای قسمت عمده چند مورفیسیم های تک نوکلئوتیدی شناخته شده تعیین ژنوتیپ می شوند، سپس میزان فراوانی الی هر کدام از این نشانگرها در گروه های مورد و شاهد مقایسه می شود. اساس این آزمون ها محاسبه میزان اثر¹⁰ بر مبنای محاسبه نسبت شانس¹¹ است بدین معنا که حضور یا عدم حضور یک آلل چقدر به حضور یا عدم حضور فنوتیپ مورد بررسی در یک جامعه ی مشخص و معلوم، مربوط است (Visscher et al. 2012).

امتیاز چند ژنی¹² (PGS) یا امتیاز خطر چند ژنی تخمینی از م سئولیت ژنتیکی یک فرد در برابر یک صفت یا بیماری است که با توجه به مشخصات ژنوتیپ و داده های مطالعه ارتباط گسترده ژنوم محاسبه می شود. در حالی که PRS های فعلی معمولاً تنها بخش کوچکی از واریانس صفت را توضیح می دهند، ارتباط آنها با بزرگترین عامل ایجاد تنوع فنوتیپی منجر به کاربرد معمول PRS ها در تحقیقات زیست پزشکی شده است (Choi, Mak, and O'Reilly 2020).

بیماری پرفشاری خون¹³ (HTN)

فشارخون در دو نوبت سنجیده می شود. عدد بالاتر که در زمان انقباض قلب محاسبه می شود، مرتبط با فشار سیستولی¹⁴ (SBP) است که معمولاً به فشارخونی که در سرخرگ ها است گفته می شود. فشارخون دیاستولی¹⁵ (DBP) به فشارخون ورید گفته می شود که در زمان استراحت قلب مابین دو نبض است. علل بیماری پرفشاری خون هنوز شناخته نشده است. پرفشاری خون بر اساس SBP بزرگتر از ۱۴۰ میلیمتر جیوه و DBP بزرگتر از ۹۰ میلیمتر جیوه تعریف می شود (K. T. Mills et al. 2016; Niu et al. 2021).

مدل های خطر بیماری

مدل های معمول خطر بیماری بر اساس همه گیر شناسی (با قدرت پیش بینی محدود) در درجه اول مبتنی بر خطر عامل های شیوه زندگی مانند وجود سابقه خانوادگی مثبت می باشند. اخیراً، در نظر گرفتن خطر عامل های ژنتیکی از جمله SNP های

⁶ Genome Wide Association Study (GWAS)

⁷ Polygenic Risk Score (PRS)

⁸ Inheritance

⁹ Genetic linkage

¹⁰ Effect size

¹¹ Odd ratio

¹² Genetic Risk Score

¹³ Hypertension

¹⁴ Systolic Blood Pressure

¹⁵ Diastolic Blood Pressure



مرتبط با بیماری و یا مربوط به فنوتیپ همراه با مدل سازی خطر، دقت پیش‌بینی فردی بیماری‌ها را بهبود بخشد است (Mosley et al. 2020).

در حال حاضر، توسعه مدل‌های خطر ژنتیکی بر روی دست‌یابی به قدرت پیش‌بینی دقیق برای شناخت افراد در خطر به شیوه‌ای قدرتمند تمرکز دارد. مطالعات GWA، SNPها را با توجه به ارتباط آن‌ها با یک بیماری/ فنوتیپ در سطح جمعیت تعریف می‌کند (Abraham and Inouye 2015).

در این تحقیق ابتدا مدل‌های ژنومی پیش‌بینی‌کننده صفات بر روی و در ادامه مدل‌های مختلف پیش‌بینی بیماری فشار خون مرور می‌گردد. این روش‌ها شامل چهار رویکرد آماری، فراتحلیل، یادگیری ماشین و مدل بالینی می‌باشند.

رویکرد آماری

رویکرد آماری، رویکرد غالب در تحلیل ارتباطات گسترده ژنومی است. اولین مدل آماری GRAMMAR¹⁶ می‌باشد (Aulchenko, De Koning, and Haley 2007). در این مدل داده‌ها تحت مدل ترکیبی¹⁷ مورد تجزیه و تحلیل قرار می‌گیرند.

پس از مدل GRAMMAR، مدل GCTA ارائه گردید (Yang et al. 2017). در روش GCTA برخلاف تجزیه و تحلیل روابط تک SNP، به اثرات همه SNPs به‌عنوان اثرات تصادفی با یک مدل خطی آمیخته¹⁸ نگاه می‌شود. افزایش عملکرد این مدل در مدل ACTA ارائه شده است (Gray, Stewart, and Tenesa 2012). با ظهور GPUهای سریع، الگوریتم جهت اجرای موازی تغییر کرده و مدل REACTA ارائه گردید (Cebamano et al. 2014). بر اساس SNPهای معرفی شده در GWASهای مرتبط با فشار خون، مدل PRS دارای AUC 0.804 می‌باشد (Vaura et al. 2021).

رویکرد فراتحلیل

در روش‌های فراتحلیل¹⁹ از مقادیر مختلفی مانند p-value و اندازه نمونه و ادغام مطالعات GWAS با یک فنوتیپ استفاده می‌شود. مهم‌ترین مدل در این رویکرد مدل METAL (Willer, Li, and Abecasis 2010) می‌باشد که امکان ادغام مطالعات مختلف در حوزه GWAS را بر اساس اندازه نمونه، واریانس داده، مقدار p-value و اندازه تاثیر ایجاد می‌کند. برای این رویکرد بسته‌های نرم‌افزاری مختلفی مانند GWAMA و MetABEL ارائه شده است (Evangelou and Ioannidis 2013).

رویکرد یادگیری ماشین

در بررسی‌های مشترک چند متغیر ژنتیکی و محیطی، چالش‌های متعددی وجود دارد. اول، در GWAS، ژنوتیپ‌ها تا یک میلیون SNP در چند هزار نفر تعیین می‌شود، که منجر به مشکل کوچک بودن نمونه و بزرگ بودن تعداد بعد

¹⁶ Genome wide Rapid Association Using Mixed Model and Regression

¹⁷ Mixed Model

¹⁸ Mixed Liner Model (MLM)

¹⁹ Meta-Analysis



می شود. دوم، هنگامی که تعداد زیادی از SNP های ژنوتیپ در یک مقیاس ژنومی هستند، باید عدم تعادل ارتباطی²⁰ بین SNP ها (که منجر به متغیرهای همبسته) شود، مورد توجه قرار گیرد (Niu et al. 2021; Szymczak et al. 2009).

رویکرد بالینی

در رویکرد بالینی علاوه بر عوامل ژنتیکی، عوامل بالینی نیز در بیماری فشار خون در نظر گرفته می شوند. عوامل مختلفی که در تحقیقات بالینی در نظر گرفته می شود شامل: سابقه خانوادگی، سن، جنس، مصرف نمک، مصرف الکل، مصرف دخانیات، کم تحرکی و چاقی از عوامل خطر بیماری فشاری خون بالا به شمار می رود. وو و همکاران بر مبنای ملاک های سنتی پیش بینی فشار خون تحقیقی را مبتنی بر روش های یادگیری ماشین ارائه داده اند که دارای AUC 0.757 برای پیش بینی فشار خون است (Wu et al. 2020). همچنین در تحقیق دیگری با استفاده عوامل بالینی سن، مصرف سیگار، سابقه خانوادگی، فعالیت فیزیکی و شاخص جرمی بدن²¹ دقت 0.854 و با اضافه کردن پارامترهای ژنتیکی دقت 0.871 بدست آمده است (Niu et al. 2021).

اهداف تحقیق

چالش های مختلفی در مدل سازی خطر مبتنی بر ژنوم و عوامل بالینی وجود دارد. اول، در GWAS، حدود یک میلیون ژنوتیپ SNP در چند هزار نفر خطای اندازه نمونه کوچک نسبت به تعداد صفات را ایجاد می کند. دوم، زمانی که تعداد زیادی از SNP های ژنوتیپ در یک منطقه ژنومی یک سان هستند، عدم تعادل پیوندی²² (LD) بین SNP ها (که منجر به متغیرهای همبسته می شود) باید در نظر گرفته شود. سوم، افزایش دقت پیش بینی در سنین پایین تر، حیاتی ترین موضوع است (Wu et al. 2020).

در نتیجه، مطالعه حاضر از روش های مبتنی بر یادگیری عمیق با رویکرد CNN برای استفاده از SNP های متعدد فشار خون در یک مطالعه کوهورت طولی استفاده می کند.

²⁰ Linkage Disequilibrium (LD)

²¹ Body Mass Index (BMI)

²² Linkage Disequilibrium



روش شناسی تحقیق

در این مقاله ابتدا با روش امتیاز خطر چند ژنی^{۲۳} (PRS) و در ادامه با استفاده از شبکه عصبی پیچشی^{۲۴} (CNN) خطر پرفشاری خون مدل سازی می شود و با توجه به اینکه پیش بینی رخداد بیماری در سنین پایین مهم است صرفاً جنسیت و سن به همراه SNP و بدون علائم بالینی وارد مدل شده است.

امتیاز خطر چند ژنی

از آنجایی که ترکیب ژنتیکی یک فرد از بدو تولد تا حد زیادی ثابت است، اطلاعات ژنتیکی این پتانسیل را دارد که به عنوان یک پیش بینی کننده اولیه خطر عمل کند. بیماری فشارخون تحت تأثیر چندین گونه ژنتیکی با اندازه های اثر کوچک^{۲۵} قرار می گیرد، بنابراین پیش بینی خطر معنی دار، بر روی تأثیر انبوه این گونه های چندگانه از طریق محاسبه یک متریک واحد که نشان دهنده خطر ژنتیکی کلی یک فرد است، ضروری است. ابتدا یک امتیاز خطر ژنتیکی^{۲۶} (GRS) ساده محاسبه شد که یک مجموع ساده از تعداد آلل های خطر (معمولاً از چند SNP از GWAS، که گاهی بر اساس اندازه های اثر وزن هستند) است که در هر فرد وجود دارد (Padmanabhan and Dominiczak 2020). اخیراً، با توجه به اینکه SNP هایی که آستانه اهمیت بسیار سختگیرانه برای ارتباط گسترده ژنومی را برآورده نمی کنند نیز می توانند پیش بینی کننده بیماری باشند، طیف وسیع تری از SNP ها، از هزاران تا میلیون ها SNP برای ایجاد یک GRS بهبود یافته به عنوان PRS استفاده شده اند (Torkamani, Wineinger, and Topol 2018). باید تأکید کرد اندازه خطر ارائه شده توسط PRS یک محدوده احتمالی است و این با اطلاعات خطر از نشانگرهای ژنتیکی اختلالات تک ژنی متفاوت است.

برای محاسبه PRS از دو نرم افزار plink و gcta64 استفاده می گردد. کد اجرا به صورت شکل ۱ می باشد. لازم به ذکر است اجرای کد روی فایل های به فرمت ped که داده های ژنومی را ذخیره کرده است، در بستر لینوکس اجرا می گردد.

```
plink.exe --ped test.ped --map test.map --noweb --make-bed --out test
gcta64 --bfile test --make-grm --out test --thread-num 10 --reml --pheno test.phen
gcta64 --reml --grm test --pheno test.phen --reml-pred-rand --out test
gcta64 --bfile test --blup-snp test.indi.blp --out test
plink --score test.snp.blp --noweb --bfile test --pheno test.phen --out test
gcta64 --reml --grm test --pheno test.phen --out test --cvblup
```

شکل (۱): کد shell محاسبه PRS

ضمناً برای اعتبار سنجی مدل ساخته شده، داده ها به دو دسته آزمایش^{۲۷} و یادگیرنده^{۲۸} تقسیم می شوند و اندازه PRS به

²³ Polygenic Risk Score

²⁴ Convolutional Neural Network

²⁵ Small effect

²⁶ Genetic Risk Score

²⁷ Test

²⁸ Train

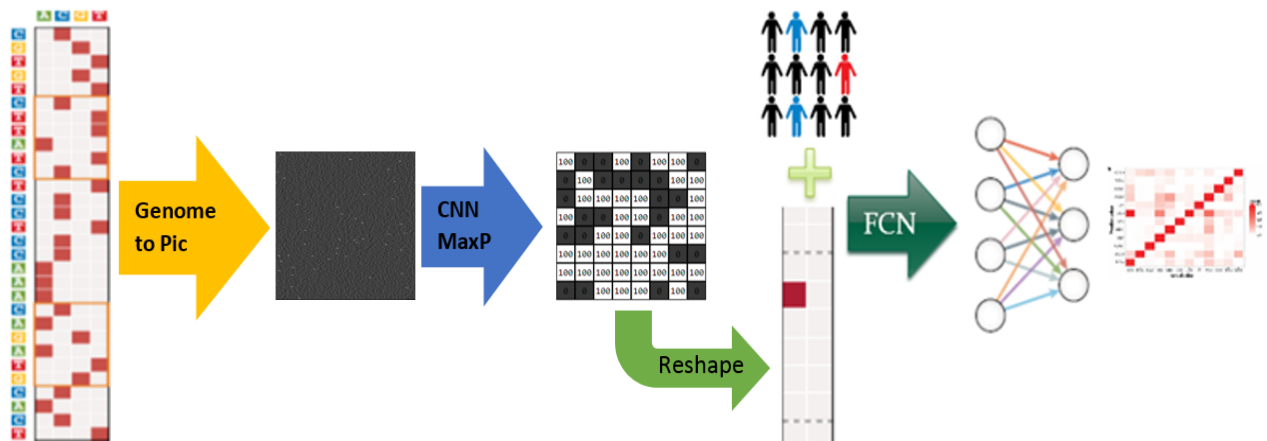


صورت جداگانه محاسبه می شود و در نهایت در مدل سازی پیش بینی کننده بیماری استفاده می گردد.

شبکه عصبی پیچشی

شبکه عصبی پیچشی برای پردازش داده هایی با دو ویژگی مناسب است: (۱) ویژگی ها از نظر فضایی ارتباط دارند. به عنوان مثال، تصویری که در آن پیکسل ها در یک آرایه مرتب شده اند و یا ویدئو، که تصاویر پشت سر هم می باشد. (۲) ویژگی ها در تمام صفحه همگن هستند. ساختار اصلی CNN از لایه های ورودی به شکل صفحه تصویر تشکیل شده است.

لایه های دیگر شامل لایه پیچش، لایه ادغام^{۲۹} و لایه کاملاً متصل^{۳۰} (FCN) و لایه خروجی می باشد (Luo et al. 2018). لایه پیچش ویژگی های محلی را بر اساس وزن نرون ها با در هم پیچیدن تمام بخش های تصویر استخراج می کند. لایه ادغام عملیات مستقلی را بر روی هر نقشه ویژگی انجام می دهد، مانند ادغام متوسط^{۳۱} یا حداکثر ادغام^{۳۲}، که می تواند به طور موثر و ضوح ویژگی را کاهش دهد و تعداد پارامترهای شبکه مورد نیاز برای بهینه سازی کاهش دهد. این CNN ساختارهای اساسی را روی هم می گذارد و خروجی لایه اول را به عنوان ورودی لایه دوم قرار می دهد که این خاصیت توانایی یادگیری عمیق دارد. مراحل اجرای CNN در تصویر ۲ آمده است.



شکل (۲): مراحل اجرای CNN

در اولین مرحله داده ژنومی که در به صورت فایل ped می باشد به فایل تصویری تبدیل می شود. برای اینکار در هر ۴ پیکسل پشت سر هم تنها یک خانه بر اساس نوع ژنوم مشکمی می شود و بقیه خانه ها به رنگ سفید می مانند. برای ساخت بهترین ورودی، طول و عرض تصویر برابر مجذور تعداد ژنوم که برابر با ۶۰۰ هزار SNP است در نظر گرفته می شود. این روش باعث کاهش حجم داده ژنومی با ضریب فشرده سازی ۱۳ درصد می شود. در مرحله بعدی شبکه عمیق مبتنی بر لایه پیچش با ۲۴ فیلتر تشکیل می گردد. فیلترها به صورت 4×4 است که در هر فیلتر در هر ردیف تنها یک ستون

²⁹ Pooling layer

³⁰ Fully Connected Network

³¹ Average pool

³² Maximum pool



افرادی که در هر دوره دارای SBP بالاتر از ۱۴۰، یا DBP بیش از ۹۰ باشد یا شرکت کننده در حال م صرف دارای کنترل فشار خون باشد، برچسب فشار خون بالا داده شد و در غیر این صورت، نرمال در نظر گرفته می شوند. نمونه ژنوم از گلبول‌های سفید با روش استاندارد پروتیناز K استخراج شده با تراشه HumanOmniExpress-24-v1-0 توسط شرکت deCODE مطابق با دستگاہ Illumina ژنوتیپ شدند (Kolifarhood et al. 2021). داده های ژنومی برای هر فرد دارای ۶۰۰ هزار SNP می باشند که در فایل ped ذخیره سازی شده است. بر مبنای (Mattoo 2019) جهت بهبود نتایج و کنترل کیفی، دسته‌بندی بر اساس سن (بالای ۱۸ سال و پایین تر از ۱۸ سال) و بخش‌بندی SNP از منظر نادر و رایج بودن انجام شده است. بر همین اساس، SNP‌هایی که دارای مقادیر گمشده^{۳۵} کمتر از ۵ درصد و میانگین فراوانی آللی^{۳۶} (MAF) بالاتر از ۵ درصد هستند، در نظر گرفته شده است.

اعتبارسنجی و ارزیابی نتایج

معیارهای مختلفی در ادبیات موضوع جهت اعتبارسنجی و ارزیابی نتایج مدل‌سازی بیان شده است که در جدول ۱ آمده است. در این جدول TP, TF, FP, FN به ترتیب مقدار مثبت درست، منفی درست، مثبت غلط و منفی غلط می باشد. این معیارها کارایی مدل ساخته شده را در داده های تحقیق بیان می کنند.

معیار	نحوه محاسبه
Accuracy	$\frac{TP + TN}{TP + TN + FP + FN}$
Recall/sensitivity/true positive rate	$\frac{TP}{TP + FN}$
Specificity	$\frac{TN}{TN + FP}$
Precision	$\frac{TP}{TP + FP}$
False Positive Rate	$\frac{FP}{FP + TN}$
F1-Score	$\frac{2 \cdot (\text{Precision}) \cdot (\text{Recall})}{\text{Precision} + \text{Recall}}$

جدول (۱): معیارهای ارزیابی مدل سازی

علاوه بر این معیارها، یکی از معیارهای پر کاربرد در مقایسه بین مدل ها، معیار سطح زیر منحنی مشخصه عملکرد گیرنده^{۳۷} (AUC) می باشد. یک طبقه بند کامل دارای AUC برابر با ۱ خواهد و یک طبقه بندی ضعیف دارای AUC ۰.۵ است. یک طبقه بندی کننده خوب دارای AUC بیشتر از ۰.۵ خواهد بود و بار سیدن به ۱ عالی در نظر گرفته می شود (Martinez-Ríos et al. 2021).

برای اعتبارسنجی مدل سازی نیز از روش 10-fold استفاده شده است که با تقسیم داده ها به ۱۰ قسمت مساوی بدون تکرار و ایجاد مدل های متفاوت، نتیجه نهایی را بر اساس متوسط نتایج بدست آمده اعلام می کند.

³⁵ missed value

³⁶ Mean Allele Frequency

³⁷ Area Under the receiver operating characteristic Curve

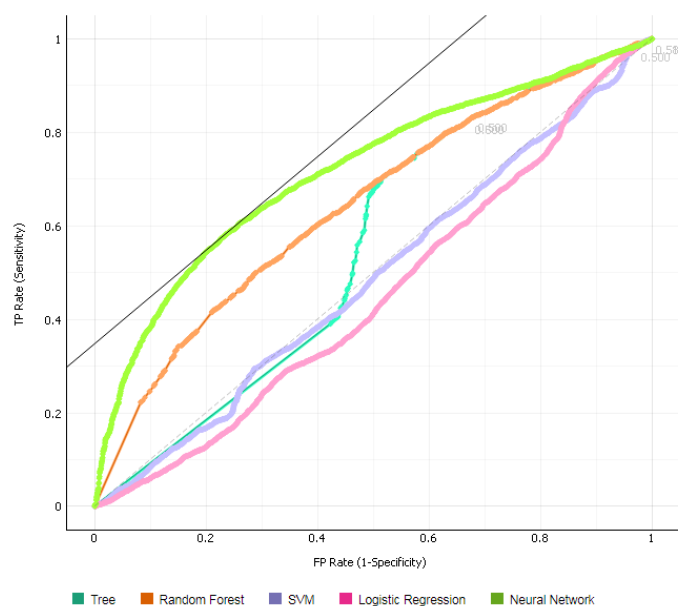


نتایج

پس از محاسبه PRS برای به دست آوردن نتایج طبقه بندی، از مدل های شبکه عصبی ساده^{۳۸}، جنگل تصادفی^{۳۹}، درخت تصمیم^{۴۰}، رگرسیون لجستیک^{۴۱} و ماشین بردار پشتیبان^{۴۲} (SVM) استفاده گردید. نتایج مدل سازی PRS در جدول ۲ و نمودار ROC در شکل ۴ آمده است.

Method	AUC	CA	F1	Precision	Recall
Neural Network	0.714	0.742	0.648	0.652	0.742
Random Forest	0.641	0.702	0.692	0.685	0.702
Tree	0.550	0.704	0.693	0.685	0.704
SVM	0.464	0.675	0.647	0.630	0.675
Logistic Regression	0.425	0.747	0.639	0.558	0.747

جدول (۲): نتایج مدل سازی PRS



شکل (۴): نمودار ROC مدل سازی PRS

در مدل سازی CNN برای لایه FCN می توان روش های مختلف طبقه بندی را در نظر گرفت که برای اینکه بتوان با روش قبلی مقایسه انجام داد، تمامی روش های گفته شده، استفاده گردید. نتایج در جدول ۳ و نمودار ROC در شکل ۵ آمده است.

³⁸ Neural Network

³⁹ Random Forest

⁴⁰ Decision Tree

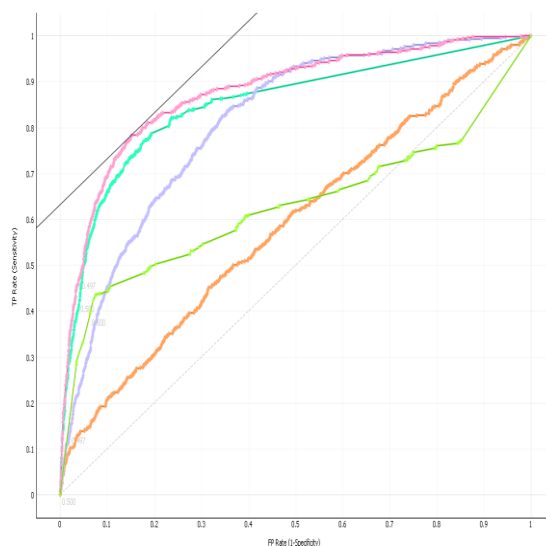
⁴¹ Logistic Regression

⁴² Support Vector Machine



Method	AUC	CA	F1	Precision	Recall
Neural Network	0.877	0.901	0.896	0.893	0.901
Random Forest	0.849	0.896	0.889	0.885	0.896
Logistic Regression	0.811	0.881	0.851	0.845	0.881
Tree	0.621	0.879	0.870	0.864	0.879
SVM	0.608	0.857	0.838	0.825	0.857

جدول (۳): نتایج مدل سازی CNN



شکل (۵): نمودار ROC مدل سازی CNN

نتایج نهایی مقایسه ای PRS و مدل سازی CNN در جدول ۴ آمده است.

	AUC	ACC	F1	PRECISION	RECALL
CNN	۰.۸۷۷	۰.۹۰۱	۰.۸۹۶	۰.۸۹۳	۰.۹۰۱
PRS	۰.۷۱	۰.۷۴	۰.۶۵	۰.۶۵	۰.۷۴

جدول (۴): مقایسه بهترین نتایج مدل سازی CNN و PRS



بحث و نتیجه گیری

در مطالعه حاضر، بیش از ۴۵۰۰۰ نمونه فشار خون در شرکت کنندگان TCGS توسط دو مدل PRS و CNN ارزیابی و پیش‌بینی گردید و بنابر نتایج مدل‌سازی CNN مقدار AUC ۰.۸۷۷ را بدون در نظر گرفتن هیچ عامل خطر اضافی نشان می‌دهد. با توجه به مقادیر AUC که در ادبیات موضوع در مدل‌سازی‌های مختلف مبتنی بر عوامل ژنومی برابر ۰.۸۰۴ (Vaura et al. 2021) و مبتنی بر عوامل ژنومی و بالینی برابر ۰.۸۷۱ (Niu et al. 2021) بدست آمده است، کارایی روش CNN کاملاً نشان داده می‌شود و یافته‌های ما عملکرد بهتری را نشان می‌دهد. ضمناً در عصر ژنومی، مهمترین دستاورد عملی، ایجاد ابزاری برای تشخیص زودهنگام خطر فشار خون بر اساس نشانگرهای ژنومی برای بهبود پیشگیری و درمان در افراد پرخطر است که رویکرد ما ابزار ارتقاء یافته‌ای را برای پیش‌بینی خطر فشار خون بالا در سنین جوانتر ارائه می‌دهد.

تحقیق ما به دلیل استفاده از اعتبارسنجی داده‌های داخلی دارای محدودیت‌هایی است. در این مطالعه، نتایج خارج از داده‌های TCGS برای اعتبارسنجی خارجی استفاده نشده است که این داده‌های نیز دارای ترکیب قومی کمی است. هدف اصلی این مطالعه ارتقای دقت پیش‌بینی خطر پرفشاری خون با ارائه راهکاری مبتنی بر شبکه عمیق CNN بود که در این مقاله با موفقیت این رویکرد در مقایسه با PRS ارائه گردید. برای کار آینده، می‌توان با اضافه کردن متغیرهای بالینی موثر، دقت به طور قابل توجهی افزایش داد.



مراجع

- Abraham, Gad, and Michael Inouye. 2015. "Genomic Risk Prediction of Complex Human Disease and Its Clinical Application." *Current Opinion in Genetics and Development* 33(Cvd): 10–16. <http://dx.doi.org/10.1016/j.gde.2015.06.005>.
- Aulchenko, Yurii S., Dirk Jan De Koning, and Chris Haley. 2007. "Genomewide Rapid Association Using Mixed Model and Regression: A Fast and Simple Method for Genomewide Pedigree-Based Quantitative Trait Loci Association Analysis." *Genetics* 177(1): 577–85.
- Azizi, Fereidoun et al. 2009. "Prevention of Non-Communicable Disease in a Population in Nutrition Transition: Tehran Lipid and Glucose Study Phase II." *Trials* 2009 10:1 10(1): 1–15. <https://trialsjournal.biomedcentral.com/articles/10.1186/1745-6215-10-5> (August 7, 2021).
- Cebamanos, L., A. Gray, I. Stewart, and A. Tenesa. 2014. "Regional Heritability Advanced Complex Trait Analysis for GPU and Traditional Parallel Architectures." *Bioinformatics* 30(8): 1177–79.
- Choi, Shing Wan, Timothy Shin-Heng Mak, and Paul F. O'Reilly. 2020. "Tutorial: A Guide to Performing Polygenic Risk Score Analyses." *Nature Protocols* 15(9): 2759–72. <http://dx.doi.org/10.1038/s41596-020-0353-1>.
- Daneshpour, Maryam S et al. 2017. "Rationale and Design of a Genetic Study on Cardiometabolic Risk Factors: Protocol for the Tehran Cardiometabolic Genetic Study (TCGS)." *JMIR Research Protocols* 6(2): e28. <https://doi.org/10.2196%2Fresprot.6050>.
- Evangelou, Evangelos, and John P.A. Ioannidis. 2013. "Meta-Analysis Methods for Genome-Wide Association Studies and Beyond." *Nature Reviews Genetics* 14(6): 379–89. <http://dx.doi.org/10.1038/nrg3472>.
- Filshtein, Teresa Jenica et al. 2019. "Reserve and Alzheimer's Disease Genetic Risk: Effects on Hospitalization and Mortality." *Alzheimer's and Dementia* 15(7): 907–16.
- Gray, A., I. Stewart, and A. Tenesa. 2012. "Advanced Complex Trait Analysis." *Bioinformatics* 28(23): 3134–36.
- Kolifarhood, Goodarz et al. 2021. "Genome-Wide Association Study on Blood Pressure Traits in the Iranian Population Suggests ZBED9 as a New Locus for Hypertension." *Scientific Reports* 11(1): 1–14. <https://doi.org/10.1038/s41598-021-90925-w>.
- Luo, Yue et al. 2018. "The Prediction of Hypertension Based on Convolution Neural Network." *2018 IEEE 4th International Conference on Computer and Communications, ICC3 2018*: 2122–27.
- Mahajan, Shiwani et al. 2019. "Prevalence, Awareness, and Treatment of Isolated Diastolic Hypertension: Insights From the China PEACE Million Persons Project." *Journal of the American Heart Association* 8(19): 1–17. <https://www.ahajournals.org/doi/10.1161/JAHA.119.012954>.
- Martinez-Ríos, Erick, Luis Montesinos, Mariel Alfaro-Ponce, and Leandro Pecchia. 2021. "A Review of Machine Learning in Hypertension Detection and Blood Pressure Estimation Based on Clinical and Physiological Data." *Biomedical Signal Processing and Control* 68(March): 102813. <https://doi.org/10.1016/j.bspc.2021.102813>.
- Mattoo, Tej K. 2019. "Definition and Diagnosis of Hypertension in Children and Adolescents - UpToDate." *UpToDate (Cv)*: 1–34. https://www.uptodate.com/contents/definition-and-diagnosis-of-hypertension-in-children-and-adolescents?search=tension%20arterial&source=search_result&selectedTitle=1~150&usage_type=default&display_rank=1#H12
- Mills, Katherine T. et al. 2016. "Global Disparities of Hypertension Prevalence and Control." *Circulation* 134(6): 441–50. <https://www.ahajournals.org/doi/10.1161/CIRCULATIONAHA.115.018912>.
- Mills, Katherine T., Andrei Stefanescu, and Jiang He. 2020. "The Global Epidemiology of Hypertension." *Nature Reviews Nephrology* 16(4): 223–37. <http://dx.doi.org/10.1038/s41581-019-0244-2>.
- Mills, Melinda C., and Charles Rahal. 2019. "A Scientometric Review of Genome-Wide Association



- Studies.” *Communications Biology* 2(1): 9. <http://dx.doi.org/10.1038/s42003-018-0261-x>.
- Mosley, Jonathan D. et al. 2020. “Predictive Accuracy of a Polygenic Risk Score Compared with a Clinical Risk Score for Incident Coronary Heart Disease.” *JAMA - Journal of the American Medical Association* 323(7): 627–35.
- Niu, Miaomiao et al. 2021. “Identifying the Predictive Effectiveness of a Genetic Risk Score for Incident Hypertension Using Machine Learning Methods among Populations in Rural China.” *Hypertension Research*. <http://dx.doi.org/10.1038/s41440-021-00738-7>.
- Padmanabhan, Sandosh, and Anna F. Dominiczak. 2020. “Genomics of Hypertension: The Road to Precision Medicine.” *Nature Reviews Cardiology*. <http://dx.doi.org/10.1038/s41569-020-00466-4>.
- Szymczak, Silke et al. 2009. “Machine Learning in Genome-Wide Association Studies.” *Genetic Epidemiology* 33(SUPPL. 1): 51–57.
- Tohidi, M., M. Hatami, F. Hadaegh, and F. Azizi. 2012. “Triglycerides and Triglycerides to High-Density Lipoprotein Cholesterol Ratio Are Strong Predictors of Incident Hypertension in Middle Eastern Women.” *Journal of Human Hypertension* 26(9): 525–32.
- Torkamani, Ali, Nathan E. Wineinger, and Eric J. Topol. 2018. “The Personal and Clinical Utility of Polygenic Risk Scores.” *Nature Reviews Genetics* 19(9): 581–90. <http://dx.doi.org/10.1038/s41576-018-0018-x>.
- Vaura, Felix et al. 2021. “Polygenic Risk Scores Predict Hypertension Onset and Cardiovascular Risk.” *Hypertension* 77(4): 1119–27. <https://www.ahajournals.org/doi/10.1161/HYPERTENSIONAHA.120.16471>.
- Visser, Peter M., Matthew A. Brown, Mark I. McCarthy, and Jian Yang. 2012. “Five Years of GWAS Discovery.” *American Journal of Human Genetics* 90(1): 7–24. <http://dx.doi.org/10.1016/j.ajhg.2011.11.029>.
- Willer, Cristen J., Yun Li, and Gonçalo R. Abecasis. 2010. “METAL: Fast and Efficient Meta-Analysis of Genomewide Association Scans.” *Bioinformatics* 26(17): 2190–91.
- Wu, Xueyi et al. 2020. “Value of a Machine Learning Approach for Predicting Clinical Outcomes in Young Patients With Hypertension.” *Hypertension* 75(5): 1271–78. <https://www.ahajournals.org/doi/10.1161/HYPERTENSIONAHA.119.13404>.
- Yang, Jian et al. 2017. “Concepts, Estimation and Interpretation of SNP-Based Heritability.” *Nature Genetics* 49(9): 1304–10.
- Zhou, Bin et al. 2017. “Worldwide Trends in Blood Pressure from 1975 to 2015: A Pooled Analysis of 1479 Population-Based Measurement Studies with 19·1 Million Participants.” *The Lancet* 389(10064): 37–55.