

شبکه‌های عصبی پیچشی حساس به هزینه برای طبقه‌بندی زیرگروه‌های سرطان

راضیه هاشمی عالم، محبوبه شمسی*، مجید آقایی
برق و کامپیوتر، صنعتی قم، قم، ایران.

چکیده

طبقه‌بندی زیرگروه‌های سرطان وظیفه بسیار مهمی برای تشخیص و پیش‌آگهی سرطان است. در سال‌های اخیر، روش‌های یادگیری عمیق به همین دلیل محبوبیت قابل توجهی به دست آورده‌اند. با این حال، تعیین ساختار شبکه عصبی دشوار است زیرا عملکرد شبکه عمیق تا حد زیادی به ساختار آن بستگی دارد. علاوه بر این، تعداد بالای ژن‌ها در پایگاه داده بیان ژن و عدم تعادل داده‌ها بین طبقات مختلف تأثیر مستقیمی بر پیچیدگی و عملکرد مدل‌های طبقه‌بندی زیرگروه سرطان دارد. برای پرداختن به مشکل داده‌های نامتعادل، یک مدل شبکه عصبی کانولوشن (CNN) با استفاده از یک استراتژی حساس به هزینه برای افزایش دقت مدل در شناسایی کلاس‌های اقلیت پیشنهاد شده است. از سوی دیگر، از تکنیک ضریب فیشر برای کاهش ژن‌ها در مرحله پیش‌پردازش استفاده می‌شود. در روش حساس به هزینه، ماتریس هزینه بر اساس توزیع کلاس‌ها ایجاد می‌شود و سپس از این ماتریس در مرحله تابع هزینه شبکه CNN برای محاسبه میزان خطا استفاده می‌شود. دو مجموعه از مجموعه داده‌های سرطان برای ارزیابی روش پیشنهادی استفاده می‌شود. نتایج با استفاده از سه معیار دقت، فراخوانی و دقت مقایسه می‌شوند. نتایج نشان می‌دهد که انتخاب ژن‌های مناسب برای طبقه‌بندی به همراه استفاده از یادگیری حساس به هزینه برای این منظور می‌تواند عملکرد روش پیشنهادی نسبت به مدل CNN بدون انتخاب ویژگی و یادگیری حساس به هزینه حدود ۱۱٪، ۱۰٪ و ۱۸٪ به ترتیب برای دقت، فراخوانی و صحت افزایش دهد.

کلمات کلیدی: دسته‌بندی، داده‌های نامتوازن، زیرگروه‌های سرطان، داده‌های بیان ژن، یادگیری عمیق CNN.

A Cost-Sensitive Convolution Neural Network for Cancer Subgroups Classification

Razieh hashemi Alam, Mahbobeh Shamsi*, Majid Aghae

Electricity and computers, Qom Industrial, Qom, Iran.

Abstract

Classification of cancer subtypes is very important task for the diagnosis and prognosis of cancer. In recent years, deep learning methods have gained considerable popularity for this reason; however, it is difficult to determine the structure of the neural network because the function of the deep network depends largely on its structure. In addition, the high number of genes in the gene expression database and the imbalanced data between different classes have a direct effect on the complexity and performance of cancer subgroup classification models. To address the problem of unbalanced data, a convolution neural network (CNN) model using a cost-sensitive strategy is proposed to increase the model's accuracy in identifying minority classes. On the other hand, the fisher ratio technique is used to reduce genes in the preprocessing stage. In techniques the cost-sensitive method, a cost matrix is created based on the distribution of classes, and then this matrix is used in the CNN network cost function step to calculate the amount of error. Two sets of cancer datasets are used to evaluate the proposed method. The results show that selecting the appropriate genes for classification along with the use of cost-sensitive learning can increase the performance of the proposed method compared to the CNN model without selecting the feature and cost-sensitive learning about 11%, 10% and 18% in terms of three criteria of accuracy, recall and precision, respectively.

Keywords: Classification, Imbalanced Data, Cancer subgroups, Gene expression, Convolution Neural Networks, Cost-sensitive learning

تاریخچه مقاله:

تاریخ ارسال: ۱۴۰۰/۱۲/۱۵

تاریخ اصلاحات: ۱۴۰۱/۰۱/۰۵

تاریخ پذیرش: ۱۴۰۱/۰۳/۲۵

تاریخ انتشار: ۱۴۰۱/۰۵/۱۱

Keywords:

Classification,
Imbalanced Data,
Cancer subgroups,
Gene expression,
Convolution Neural
Networks,
Cost-sensitive learning

*ایمیل نویسنده مسئول:

shamsi@qut.ac.ir

۱- مقدمه

نشان دهنده این است که آیا یک فرد مبتلا به سرطان است یا نه، طبقه بندی اشتباه یک بیمار به عنوان یک فرد سالم منجر به هزینه بسیار بیشتری در مقایسه با طبقه بندی یک فرد سالم به عنوان یک بیمار خواهد شد. به این دلیل که تشخیص اشتباه ممکن است باعث تأخیر در درمان یا مرگ بیمار شود. یادگیری حساس به هزینه یک استراتژی برای به حداقل رساندن هزینه کلی یادگیری است که باعث می شود یک مدل یادگیری به گونه ای باشد که روند آموزش نسبت به کلاس هایی که هزینه کمتری دارند حساس تر باشد.

هنگام استفاده از یادگیری حساس به هزینه در مدل های یادگیری ژرف، فرایند آموزش نسبت به کلاس هایی که هزینه بالاتری دارند حساس تر است. برخی از تلاش های تحقیقاتی هزینه های خاص کلاس را در طبقه بندی کننده های یادگیری ژرف بررسی کرده اند. اولین بار یادگیری ژرف حساس به هزینه توسط چونگ^۱ و همکاران معرفی شد [۸]، که هزینه ها را در تابع خطا^۲ مرحله قبل از آموزش DNN و CNN ادغام شد. وانگ^۳ و همکاران [۹] با در نظر گرفتن میانگین خطا در هر کلاس، عملکرد از دست دادن میانگین خطای مربع (MSE^4) برای DNN را بهبود بخشید. خان و همکاران [۱۰] یک روش ابتکاری را برای اختصاص دادن هزینه به طور خودکار به هر کلاس با توجه به توزیع داده های کلاس، فرموله کرد. برخی از توابع ضرر که به طور گسترده مورد استفاده قرار می گیرند، مانند MSE، کراس آنترپوی و SVM Hinge، با استفاده از یادگیری حساس به هزینه بهبود یافتند. تلیکانی و همکاران [۱۱] یک SAE حساس به هزینه ایجاد کرد که در آن هزینه ها در عملکرد از بین رفتن آنترپوی متقابل قرار گرفتند. برتری این روش نسبت به سایر رویکردهای یادگیری ژرف حساس به هزینه این است که نیازی به استفاده از ماتریس هزینه دست ساز نیست زیرا هزینه ها از طریق آمار داده ها تعیین می شود.

این مقاله هم استراتژی های یادگیری حساس به هزینه و هم استراتژی کاهش ویژگی را برای رسیدگی به مشکل عدم تعادل طبقاتی و مشکلات ابعادی در طبقه بندی زیرگروه سرطان ادغام می کند. تکنیک ضریب فیشر برای حذف ژن های نامربوط و غیرمفید در مرحله پیش پردازش استفاده می شود. علاوه بر این، عملکرد تلفات متقابل آنترپوی با استفاده از استراتژی یادگیری حساس به هزینه با ادغام هزینه های مربوط به کلاس هنگام محاسبه ارزش تلفات در طول آموزش CNN، بهبود می یابد. در این استراتژی هزینه ها بر اساس آمار داده های دسته های سرطان تعریف می شود. این رویکرد باعث می شود که مدل CNN نسبت به کلاس های سرطان با فرکانس پایین حساس باشد و عملکرد مدل بیان ژن را در این نوع سرطان ها افزایش می دهد. آزمایش های مختلفی روی مجموعه داده GBM انجام می شود و نتایج

طبقه بندی زیرگروه های سرطان بر اساس داده های بیان ژن ابزار بسیار قدرتمندی جهت تحلیل رفتار هم زمان هزاران ژن است. وجود ژن های بسیار زیاد در طبقه بندی داده های حاصل از بیان ژن، باعث به وجود آمدن مشکلات زیادی در تحلیل این داده ها شده است. مشکلاتی از قبیل وجود نویز و اطلاعات غیرمفید در برخی از ژن ها که نه تنها در تحلیل اطلاعات مفید نیستند، بلکه باعث طبقه بندی نادرست اطلاعات نیز می شوند. همچنین، اطلاعات مشابه در برخی از ژن ها باعث عدم اعتبار بعضی از روش های تحلیل داده ها می شود. علاوه بر این، مدل سازی شبکه های عصبی با داده های با ابعاد بالا باعث افزایش هزینه محاسباتی و پیچیدگی روش های طبقه بندی می شود. با این حال، تنها مجموعه کوچکی از ژن ها عامل بیماری هستند، بنابراین حضور ژن های بسیار زیاد باعث کم رنگ شدن اثر ژن های عامل بیماری خواهند شد.

در سال های اخیر، روش های یادگیری ماشین و یادگیری عمیق برای یافتن ژن های مفید و همچنین طبقه بندی داده های بیان ژن مورد استفاده قرار گرفته اند. با این حال، به دلیل اندازه بزرگ داده های بیان ژن، بسیاری از روش های قبلی از روش های آماری برای فیلتر کردن ژن ها استفاده می کردند. در واقع قبل از استفاده از روش های یادگیری ماشینی و یادگیری عمیق، ژن های نامربوط با تکنیک هایی مانند آزمون t حذف می شوند و تنها ژن هایی که طبقه بندی خوبی ارائه می دهند برای ساخت مدل استفاده می شوند. مسئله عدم توازن بر روی عملکرد دسته بندی تأثیر بسیاری دارد. در این میان، الگوریتم هایی که مسئله عدم توازن کلاس را در نظر نمی گیرند، تمایل دارند که توسط کلاس اکثریت تحت پوشش قرار داده شوند و در مقابل توسط کلاس اقلیت نادیده گرفته شوند. در مسائلی که سطح عدم توازن در آن ها زیاد است، برای طراحی یک دسته بندی خوب می بایست سطح عدم توازن با دقت مدیریت شود. تکنیک هایی که مشکلات مربوط به مجموعه داده های نامتعادل را در برمی گیرند می توانند به دو گروه تقسیم شوند. دسته اول روش های سطح داده است که بر روی مجموعه آموزشی کار می کنند و توزیع کلاس ها را تغییر می دهند. دسته دیگر روش های سطح دسته بندی (الگوریتمی) را در برمی گیرد. این روش ها مجموعه داده های آموزشی را بدون تغییر نگه می دارند و الگوریتم های آموزشی طبق با توزیع داده ها تنظیم می شوند.

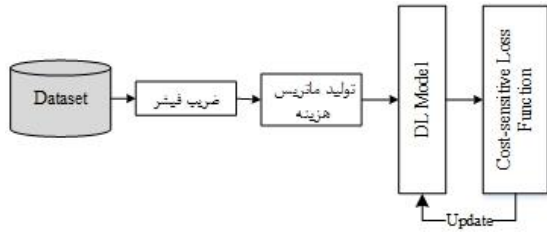
استراتژی یادگیری حساسیت به هزینه یکی از تکنیک های مبتنی بر سطح دسته بندی است که هزینه های مرتبط با دسته بندی نادرست نمونه ها را در نظر می گیرد. در تشخیص سرطان، تفاوت هزینه بین خطاهای طبقه بندی اشتباه بسیار زیاد است. در یک سیستم تشخیص سرطان که در آن هر کلاس

¹Chung

²Loss function

³Wang

⁴Mean Squared Error (MSE)



(شکل-۱): روش پیشنهادی

۳-۱- انتخاب ژن: داده‌های بیان ژن به‌طور کلی از هزاران ژن تشکیل شده است، درحالی‌که تعداد نمونه‌های موجود اغلب اندک است. در میان هزاران ویژگی در داده‌های بیان ژن، تنها چند ژن در واقع با زیرگروه‌های سرطان همراه هستند درحالی‌که بقیه ممکن است به‌عنوان ویژگی‌های زائد یا عامل اغتشاش در نظر گرفته شوند. بنابراین، انتخاب ژن می‌تواند به‌عنوان یک مشکل کاهش ابعاد در نظر گرفته شود که سعی می‌کند ضمن حفظ دقت طبقه‌بندی ژن‌های اصلی، ژن‌های مهم را نیز انتخاب کند.

نسبت فیشر نسبت فواصل کلاس با فواصل طبقه‌بندی شده است. اگر دو طبق بندی در یک مجموعه داده وجود داشته باشد، هر نمونه می‌تواند به‌عنوان نشان داده شود چون $\gamma \in \{+1, -1\}$ و داده‌های بیان ژن می‌توانند باشند به‌عنوان $x_i = \{x_1^i, \dots, x_n^i\}$ برای هر ژن میزان انحراف استاندارد (σ_i^+, σ_i^-) (resp.) و انحراف میانگین (μ_i^+, μ_i^-) (resp.) محاسبه شده و میزان ضریب فیشر بر اساس فرمول زیر محاسبه می‌شود:

$$F_i = \frac{(\mu_i^+ - \mu_i^-)^2}{(\sigma_i^+)^2 + (\sigma_i^-)^2}$$

ژن با بالاترین مقدار F_i آموزنده‌ترین است و میزان بیشترین تفاوت داده‌های بیان ژن را در دو کلاس به‌طور متوسط نشان می‌دهد درحالی‌که کلاس‌های مربوطه به ژن‌هایی که دارای انحراف کوچک هستند، بیشترین اختلاف را نشان می‌دهند. سپس ژن‌هایی با مقادیر F_i بالا به‌عنوان ویژگی‌های برتر انتخاب می‌شوند.

۳-۲- تولید ماتریس هزینه: برای آموزش مدل CNN با استفاده از هزینه‌های مربوط به دسته‌بندی‌های مختلف، ایجاد یک ماتریس هزینه ضروری است. این ماتریس در تابع هزینه برای محاسبه مقدار خطای دسته‌بندی استفاده می‌شود. برخلاف بسیاری از روش‌های قبلی برای تولید ماتریس هزینه که به‌صورت دستی و از طریق کاربرمتمخصص وزن مربوط به هر دسته‌بندی تعیین می‌شود، روش پیشنهادی از یک مکاشفه استفاده می‌کند تا هزینه‌ها به‌صورت خودکار و بدون دخالت کاربر مشخص شوند. این هزینه‌ها با در نظر گرفتن توزیع کلاس‌ها تعیین می‌شوند. شکل ۲ فرآیند تولید ماتریس هزینه γ را نشان می‌دهد.

از نظر دقت، فراخوانی و صحت ارزیابی می‌شوند. نتایج نشان می‌دهد که چارچوب پیشنهادی ما می‌تواند عملکرد مدل CNN را برای بیان ژن، به‌ویژه برای سرطان‌های با فرکانس پایین، بهبود بخشد.

باقیمانده مقاله به شرح زیر سازماندهی شده است: بخش ۲ به‌طور خلاصه کارهای مرتبط در انتخاب ژن و رویکردهای طبقه‌بندی مبتنی بر یادگیری عمیق برای بیان ژن را خلاصه می‌کند. در بخش ۳، یک معماری CNN با استفاده از یادگیری حساس به هزینه معرفی شده است. بخش ۴ نتایج ارزیابی کارایی را ارائه می‌دهد و سپس نتیجه‌گیری در بخش ۵ استنتاج می‌شود.

۲- کارهای مرتبط

سه دسته برای انتخاب ژن وجود دارد که شامل فیلتر، بسته‌بندی و جاسازی شده است. در [۲] از آزمون t برای غلبه بر مشکل پراکندگی برای انتخاب ژن استفاده شد. لیائو و همکاران [۳] از آزمون مجموع رتبه ویلکاکسون به همراه ماشین بردار پشتیبان برای ارزیابی اهمیت ژن‌ها استفاده کرد. رویکردهای Wrapper از یک طبقه‌بندی کننده برای ارزیابی عملکرد یک زیرمجموعه ویژگی استفاده می‌کنند. k-نزدیکترین همسایه [۴]، شبکه عصبی [۵] و ماشین بردار پشتیبان [۶] طبقه‌بندی کننده‌های پرکاربرد برای روش لفاف هستند. هو و همکاران [۷] مدل مجموعه خشن همسایگی را برای پردازش مجموعه داده‌های بیان ژن گسسته و پیوسته پیشنهاد کرد. یک مدل جنگل عصبی در [۱]، مجموعه‌ای از مدل درخت عصبی برای طبقه‌بندی زیرگروه‌های سرطان پیشنهاد شد. مدل جنگل پیشنهادی یک مسئله چند طبقه‌بندی را به بسیاری از مسائل طبقه‌بندی دودویی برای هر جنگل تبدیل می‌کند. یک رویکرد انتخاب ژن با ترکیب نسبت فیشر و مجموعه خشن همسایگی ایجاد شد.

روش DeepGene [۷] یک شبکه عصبی عمیق بهبودیافته برای بیان ژن است. ابتدا از تکنیک Clustered Gene Filtering برای حذف ژن‌های نامربوط استفاده شد. سپس، طبقه‌بندی کننده DNN برای استخراج ویژگی‌های سطح بالا برای طبقه‌بندی استفاده شد.

۳- متدولوژی پیشنهادی

در این بخش، یک مدل CNN مبتنی بر یادگیری حساس به هزینه معرفی می‌شود که شامل چهار مرحله است: انتخاب ژن، تولید ماتریس هزینه، مدل یادگیری ژرف و تابع زیان حساس به هزینه. شکل ۱ مراحل روش پیشنهادی را نشان می‌دهد و در ادامه هر یکی از مراحل به‌طور مفصل بحث می‌شود.

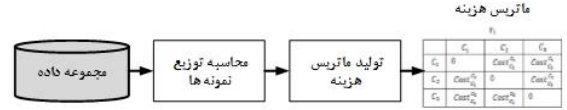
تعیین می شود، در روش پیشنهادی هزینه های مرتبط با هر کلاس به صورت خودکار با استفاده از توزیع داده ها در طول فرآیند یادگیری تنظیم می شود.

هدف ما مجازات کردن انواع خطاهای طبقه بندی بر اساس برخی هزینه های تعیین شده است. این مقدار جریمه برای زمانی که نمونه اقلیت به عنوان کلاس اکثریت طبقه بندی می شود بیشتر از زمانی است که نمونه اکثریت به اشتباه به عنوان کلاس اقلیت طبقه بندی می شود. همان طور که در بخش قبلی اشاره کردیم، کلاس های اقلیت و اکثریت تعیین می شوند و فقط باید هزینه مربوطه را از ماتریس هزینه پیدا کنیم. برتری الگوریتم ما این است که تعیین نوع کلاس ها از نظر اقلیت یا اکثریت لازم نیست. در واقع، هزینه ها فقط بر اساس توزیع کلاس ها اختصاص می یابد. این ویژگی کمک می کند تا الگوریتم در هر مجموعه داده استفاده شود.

این رویکرد قصد دارد با در نظر گرفتن مقادیر هزینه مربوط به هر نوع طبقه بندی نادرست، تابع هزینه cross-entropy اصلاح کند. این روش باعث حساسیت بیشتر مدل CNN نسبت به طبقه بندی نادرست کلاس های اقلیت می شود. در واقع، خروجی لایه Softmax که به شکل احتمالات است، به عنوان ورودی تابع هزینه در نظر گرفته می شود تا مقدار زیان حساس به هزینه محاسبه شود. دلیل انتخاب تابع cross-entropy این است که می تواند در بیشتر موارد نسبت به توابع هزینه دیگر عملکرد بهتری داشته باشد. علاوه بر این، cross-entropy می تواند از کاهش سرعت یادگیری که یکی از مشکلات تابع میانگین خطای مربع (MSE) در یادگیری است جلوگیری کند.

قبل از تشریح استراتژی تابع زیان حساس به هزینه پیشنهادی، نحوه عملکرد لایه Softmax توضیح داده می شود. فرض کنید لایه خروجی به صورت

$\{X, Y\} = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_c)\}$ باشد، که $x_i \in \mathbb{R}^{d \times 1}$ و $y_i \in \mathbb{R}^{c \times 1}$ هستند. اصطلاح d اندازه لایه خروجی و C تعداد کلاس ها هستند. تابع Softmax احتمال اینکه نمونه i (X_i) متعلق به یک کلاس باشد را محاسبه می کند:



(شکل-۲): فرآیند تولید ماتریس هزینه

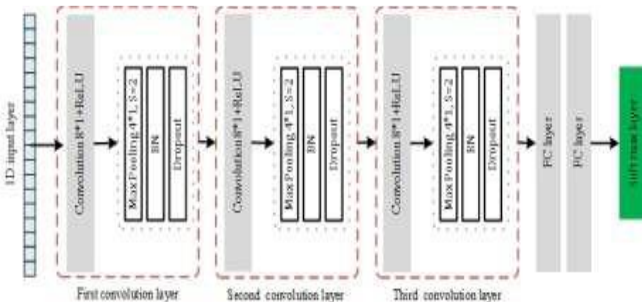
در مرحله اول، توزیع هر کلاس در مجموعه داده محاسبه می شود تا برای تولید ماتریس هزینه استفاده شود. برای تولید ماتریس هزینه، یک فرموله سازی مبتنی بر توزیع داده ها انجام می شود. هزینه بالاتر طبقه بندی نادرست برای کلاس های اقلیت در نظر گرفته می شود در حالی که هزینه طبقه بندی پایین تری برای کلاس های اکثریت تعیین می شود. هزینه طبقه بندی نادرست کلاس i در کلاس j با استفاده از رابطه ۳-۱ محاسبه می شود.

$$\begin{cases} \gamma_{i,j} = \frac{\alpha_i}{\alpha_i + \alpha_j} & i, j = 1, 2, \dots, C \\ \text{subject to } i \neq j \end{cases} \quad (1)$$

در یک ماتریس هزینه، سطر مورب ماتریس به عنوان بردار سودمندی شناخته می شود. این بردار نشان دهنده طبقه بندی های صحیح را نشان می دهد و به صفر تنظیم می شود. همچنین، تمام هزینه ها غیر منفی هستند، یعنی $\gamma_{i,j} > 0$. در این رابطه، α_i و α_j به ترتیب تعداد نمونه های کلاس های i و j هستند.

۳-۳- معماری شبکه عصبی پیچشی: این بخش معماری CNN برای روش پیشنهادی را تشریح می کند (شکل ۳) که دارای یک لایه ورودی یک بعدی و سه لایه کانولوشن است که هر یک از آن ها دارای یک کانولوشن و به دنبال آن لایه های تابع فعال ساز Relu و ادغام حداکثری است. اندازه فیلتر برای لایه کانولوشن 8×8 و $\text{stride} = 1$ است و هر لایه ادغام حداکثری یک ورودی 4×4 را با $\text{stride} = 2$ پردازش می کند. بعد از هر لایه ReLU، از نرمال سازی دسته ای و Dropout با نسبت 0.05 استفاده می شود. بعد از لایه های کانولوشن، دولا به اتصال کامل برای طبقه بندی ژن استفاده شد.

۳-۴- تابع هزینه حساس به هزینه: در این بخش یک تابع هزینه حساس به هزینه پیشنهاد می شود که نسبت به طبقه بندی نادرست کلاس های اقلیت حساس تر است. در طول آموزش، روش یادگیری پیشنهادی به طور مشترک هزینه های وابسته به کلاس و پارامترهای شبکه عصبی را بهینه می کند. در مقایسه با رویکردهای سطح داده (نمونه برداری مجدد)، روش پیشنهادی توزیع داده های اصلی را تغییر نمی دهد، که در نتیجه هزینه های محاسباتی پایین تر در طول فرآیند آموزش به دست می دهد. علاوه بر این، برخلاف روش های حساس به هزینه که از یک ماتریس هزینه دستی که بر اساس نظر یک متخصص



(شکل ۳): حساس به هزینه CNN معماری

می‌شود. در تمام آزمایش‌ها، ۸۰٪ داده‌ها به‌عنوان مجموعه آموزشی، ۱۰٪ به‌عنوان مجموعه اعتبارسنجی و ۱۰٪ به‌عنوان مجموعه آزمایشی استفاده شد.

۴-۱- معیارهای ارزیابی: برای ارزیابی روش پیشنهادی، از چهار معیار دقت (رابطه ۶)، فراخوانی (رابطه ۷) و صحت (رابطه ۸) استفاده می‌شود.

$$\text{Accuracy} = \frac{TP+TN}{FP+FN+TP+TN} \quad (۶)$$

$$\text{Recall} = \frac{TP}{TP+FN} \quad (۷)$$

$$\text{Precision} = \frac{TP}{TP+FP} \quad (۸)$$

۴-۲- مجموعه داده‌ها: در این پژوهش، مجموعه داده زیرگروه سرطان بیان ژن RNA-Seq برای GBM (گلیوبلاستوما چندحالتی) برای آزمایش‌ها استفاده شد. مجموعه داده GBM دارای ۲۱۲ نمونه سرطانی با حدود ۱۲ هزار ژن است که در چهار گروه ژنی دسته‌بندی می‌شوند. در مجموعه داده GBM دو کلاس TCGA02 و TCGA06 دارای توزیع یکسانی هستند (حدود ۸۰ نمونه) و دو کلاس دیگر دارای توزیع اقلیتی هستند.

۴-۳- نتایج و بحث: در این بخش عملکرد روش پیشنهادی بر روی مجموعه داده GBM و همچنین نسخه متوازن شده مجموعه داده GBM بررسی می‌شود. جدول ۲ عملکرد روش پیشنهادی بر روی مجموعه داده GBM بر اساس دو معیار فراخوانی و صحت را نشان می‌دهد. روش فیشر با ژن‌های اصلی مقایسه شده‌اند. همان‌گونه که در جدول مشاهده می‌شود، روش پیشنهادی بر روی ژن‌های انتخاب‌شده توسط تکنیک ناهنجاری بهترین عملکرد را داشته است. همچنین شکل ۴ ماتریس‌های آشفتگی برای روش پیشنهادی بر روی داده اصلی را نشان می‌دهد.

(جدول ۲): عملکرد روش پیشنهادی برای مجموعه داده GBM

کلاس	انتخاب ژن	فراخوانی	صحت
TCGA02	مجموعه اصلی	۰.۶۷	۰.۷۵
	فیشر	۰.۸۸	۰.۷۲
TCGA06	مجموعه اصلی	۰.۴	۰.۶۷
	فیشر	۰.۶	۱.۰
TCGA08	مجموعه اصلی	۰.۵	۰.۴
	فیشر	۰.۷۵	۰.۷۵
TCGA12	مجموعه اصلی	۱.۰	۰.۶۷
	فیشر	۰.۷۵	۰.۷۵

$$f_{\theta}(x) = \frac{1}{\sum_{j=1}^C e^{y_j}} = \begin{bmatrix} p(y_i = 1|x_i) \\ p(y_i = 2|x_i) \\ \dots \\ p(y_i = C|x_i) \end{bmatrix} \quad (۲)$$

متغیر θ پارامتر نگاشت برای کلاس j است $(b_j + W_j x)$. رویکرد پیشنهادی در این پژوهش مجازات کردن طبقه‌بندی نادرست در تابع هزینه cross-entropy بر اساس هزینه‌های تعیین‌شده در ماتریس هزینه (Y) برای به حداکثر رساندن نزدیکی پیش‌بینی به کلاس واقعی را مدنظر قرار می‌دهد. مقدار کل هزینه هر دسته با N نمونه با استفاده از معادله ۳ محاسبه می‌شود:

$$\mathcal{L}(O, y) = -\frac{1}{N} \sum_{j=1}^N \mathcal{L}(O_j, y_j) \quad (۳)$$

که در آن مقدار cross-entropy میانگین مقادیر زیان برای کل N طبقه‌بندی است. مقدار زیان برای هر پیش‌بینی توسط معادله ۴ محاسبه می‌شود:

$$\mathcal{L}(O_i, y_i) = -\sum_{c=1}^C (y_{i,c} \log p(y_i = c|x_i; \theta_i)) \quad (۴)$$

در این رابطه، $y_{i,c}$ یک شاخص دودویی (۰ یا ۱) که به پیش‌بینی صحیح مشاهده برای نمونه O اشاره دارد. مقدار $y_{i,c}$ برای کلاس اشتباه پیش‌بینی شده ۱ و برای کلاس واقعی ۰ است. احتمال طبقه‌بندی اشتباه با در نظر گرفتن هزینه مربوط به کلاس تغییر می‌یابد (معادله ۵):

$$p(y_i = 1|x_i) = \frac{y_{i,j} \cdot \exp(O_i)}{\sum_{i=1}^C \exp(O_i)} \quad (۵)$$

بر اساس معادله ۳-۶، ضرب هزینه مربوط به کلاس‌های اقلیت، مقدار احتمال جدید را به‌شدت کاهش می‌دهد و بنابراین، منجر به افزایش مقدار زیان طبقه‌بندی در رابطه ۳-۵ می‌شود. به‌این‌ترتیب، کلاس‌های اقلیت بیشتر از کلاس‌های اکثریت بر روی تابع هزینه تأثیر می‌گذارند.

۴- نتایج ارزیابی

در این فصل، عملکرد مدل CNN حساس به هزینه پیشنهادی با مدل‌های دیگر مقایسه می‌شود. این مدل‌ها شامل نسخه‌های غیرحساس برای CNN است. کتابخانه Keras و Tensorflow به‌عنوان Backend برای اجرای مدل‌های DL مورد استفاده قرار گرفتند. تمام مدل‌ها با ۱۰۰ دوره آموزش دیده بودند. استراتژی توقف اولیه برای جلوگیری از مشکل بیش‌برازش مورد استفاده قرار گرفت که در آن زمانی که مقدار خطا بر روی داده‌های اعتبارسنجی برای چندین دوره تغییر نکرده باشد، فرایند آموزش متوقف می‌شود. از تابع Adam به‌عنوان بهینه‌ساز برای شبکه‌های عصبی استفاده

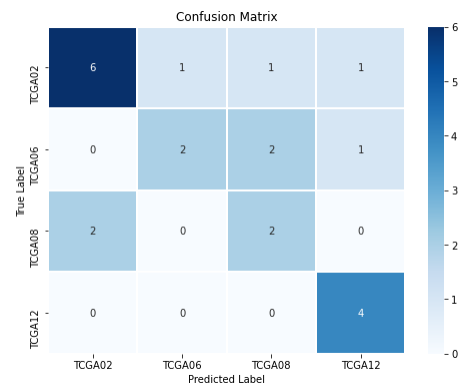
⁵Glioblastoma Multiforme

۵- نتیجه گیری

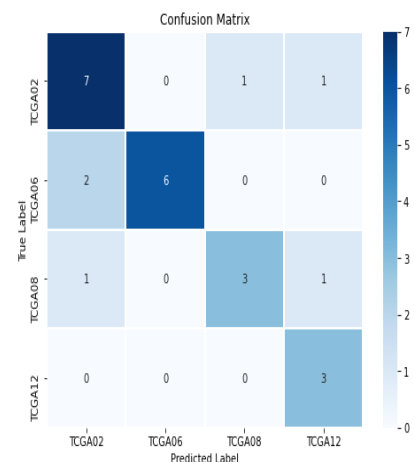
یادگیری عمیق (DL) یک تکنیک پرکاربرد در ناحیه بیان ژن است. با این حال، این الگوریتم‌ها در مورد داده‌های با ابعاد بالا و همچنین داده‌های نامتعادل با چالش‌های زیادی روبرو هستند. تعداد زیاد ویژگی‌ها، پیچیدگی مدل‌های یادگیری عمیق را افزایش می‌دهد و همچنین عدم تعادل بین کلاس‌های سرطان، عملکرد مدل طبقه‌بندی را کاهش می‌دهد. برای رسیدگی به این چالش‌ها، یک رویکرد CNN پیشنهاد شد که با تکنیک انتخاب ژن ادغام شد. یک استراتژی حساس به هزینه نیز برای مقابله با داده‌های نامتعادل استفاده شد. در این استراتژی، طبقه‌بندی‌های اشتباه مختلف دارای هزینه‌های مشخصی هستند که در هنگام محاسبه میزان خطا اعمال می‌شود و در روش پیشنهادی، ماتریس هزینه با استفاده از یک تابع فرمول‌بندی شده تعیین می‌شود. برای ارزیابی روش پیشنهادی از مجموعه داده استفاده شد. معیارهای دقت، یادآوری، دقت و F1-Score برای مقایسه روش پیشنهادی با مدل‌های مشابه استفاده شد. اجرای حساس به هزینه CNN با نسخه‌های غیر حساس مقایسه شد. نتایج نشان داد که روش پیشنهادی توانایی تشخیص بالاتری برای طبقات اقلیت دارد. به طور متوسط، مدل پیشنهادی عملکرد تشخیص سرطان زیرگروه را حدود ۳٪ افزایش داده است.

۵- منابع

- [1] Xu, J., Wu, P., Chen, Y., Meng, Q., Dawood, H. and Khan, M.M., 2019. A novel deep flexible neural forest model for classification of cancer subtypes based on gene expression data. *IEEE Access*, 7, pp.22086-22095.
- [2] Zhu, S., Wang, D., Yu, K., Li, T. and Gong, Y., 2008. Feature selection for gene expression using model-based entropy. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 7(1), pp.25-36.
- [3] Liao, C., Li, S. and Luo, Z., 2006, November. Gene selection using wilcoxon rank sum test and support vector machine for cancer classification. *In International Conference on Computational and Information Science* (pp. 57-66). Springer, Berlin, Heidelberg.
- [4] L. Li, C. Weinberg, T. Darden, and L. G. Pedersen, "Gene selection for sample classification based on gene expression data: Study of sensitivity to choice of parameters of the GA/KNN method," *Bioinformatics*, vol. 17, no. 12, pp. 1131-1142, Dec. 2001.
- [5] J. Khan et al., "Classification and diagnostic prediction of cancers using gene expression



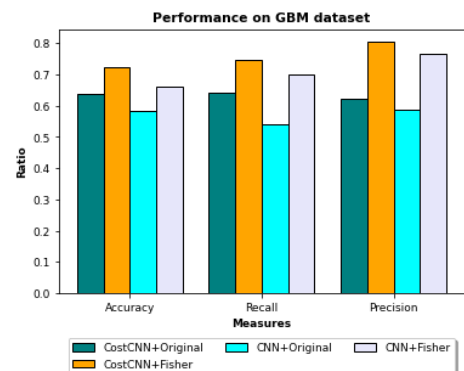
(الف) مجموعه داده اصلی



(ب) ضریب فیشر

(شکل-۴): ماتریس آشفتگی برای مدل با مجموعه ژن‌های انتخاب شده متفاوت برای مجموعه داده GBM اصلی

شکل ۵ میانگین عملکرد مدل‌های بیان ژن را نشان می‌دهد. بدترین عملکرد برای ژن‌های منتخب با تکنیک نسبت فیشر بود. در مقابل، بالاترین عملکرد مربوط به کاربرد روش پیشنهادی بر روی مجموعه داده انتخابی با استفاده از تکنیک ترکیبی بود.



(شکل-۵): عملکرد مدل‌های یادگیری ژرف برای دسته‌بندی



دکتر مجید آقایی، نرم افزار، دانشکده مهندسی کامپیوتر نرم افزار هیئت علمی دانشگاه صنعتی قم
aghaee@qut.ac.ir

روش ارجاع به مقاله : ر. هاشمی عالم، م. شمس، م. آقایی. شبکه های عصبی پیچشی حساس به هزینه برای طبقه بندی زیرگروه های سرطان. دوفصلنامه محاسبات و سامانه های توزیع شده سال پنجم، شماره اول، شماره پیاپی ۹، صفحه ۱۵ تا ۲۱، سال ۱۴۰۱.

How to cite: Razieh hashemi alam, mahbobeh shamsi, majid aghaei. A cost-sensitive convolution neural network for cancer subgroups classification. Journal of Distributed Computing and Systems(JDACS), Vol 5, Issue 1, Page 15-21, 2022.

- pro_ling and arti_cial neural networks," Nature Med., vol. 7, no. 6, pp. 673_679, 2001.*
- [6] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Mach. Learn., vol. 46, nos. 1_3, pp. 389_422, 2002.*
- [7] Yuan, Y., Shi, Y., Li, C., Kim, J., Cai, W., Han, Z. and Feng, D.D., 2016. DeepGene: an advanced cancer type classifier based on deep learning and somatic point mutations. *BMC bioinformatics, 17(17), pp.243-256.*
- [8] Y.-A. Chung, H.-T. Lin, and S.-W. Yang, "Cost-aware pretraining for multiclass cost-sensitive deep learning," *arXiv preprint arXiv:1511.09337, 2015.*
- [9] S. Wang, W. Liu, J. Wu, L. Cao, Q. Meng, and P. J. Kennedy, "Training deep neural networks on imbalanced data sets," in *International Joint Conference on Neural Networks. IEEE, 2016, pp. 4368–4374.*
- [10] S. H. Khan, M. Hayat, M. Bennamoun, F. A. Sohel, and R. Togneri, "Cost-sensitive learning of deep feature representations from imbalanced data," *IEEE Transactions on Neural Networks and Learning Systems, vol. 29, no. 8, pp. 3573–3587, 2017.*
- [11] A. Telikani and A. H. Gandomi, "Cost-sensitive stacked auto-encoders for intrusion detection in the internet of things," *Internet of Things, p. 100122, 2019.*



راضیه هاشمی عالم، کارشناسی ارشد نرم افزار، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی قم
hashemialam.r@qut.ac.ir



دکتر محبوبه شمس، دکتری مهندسی کامپیوتر نرم افزار (پردازش تصویر)، استادیار، هیئت علمی دانشگاه صنعتی قم
shamsi@qut.ac.ir