

## دسته بندی داده های چند برچسبی با استفاده از وابستگی بین کلاس ها و تخمین تعداد برچسب ها

سجاد کمالی<sup>۱</sup>، حمید بیگی<sup>۲</sup>

<sup>۱</sup>دانشجوی ارشد دانشگاه شریف، skamali@ce.sharif.edu

<sup>۲</sup>استاد دانشگاه شریف، beigy@sharif.edu

چکیده - با افزایش حجم داده ها نیاز به دسته بندی و تحلیل داده ها به صورت خودکار، از جایگاه ویژه ای برخوردار شده است. در این مقاله به دسته بندی داده های چند برچسبی خواهیم پرداخت. داده های چند برچسبی داده هایی هستند که در آن ها نمونه ها می توانند بیش از یک برچسب داشته باشند. به عبارت دیگر هر نمونه توسط یک بردار از برچسب ها نمایش داده می شود.

در این مقاله یک روش مبتنی بر نزدیک ترین همسایگی برای دسته بندی داده های چند برچسبی پیشنهاد شده است. در روش ارائه شده نخست تعداد برچسب ها برای هر نمونه جدید تخمین زده می شود و سپس برای هر نمونه جدید، با استفاده از ارتباط و وابستگی بین برچسب ها، برچسب های آن را پیش بینی می گردد. در این مقاله از دسته بندی نزدیک ترین همسایگی به عنوان دسته بندی پایه استفاده شده است اما با این تفاوت که برای محاسبه فاصله بین نمونه ها از معیار فاصله اقلیدسی استفاده نمی شود. با ارزیابی روش پیشنهادی بر روی مجموعه داده های موجود در حوزه داده های چند برچسبی و مقایسه آن با دیگر روش های مطرح در این زمینه، نشان داده می شود که روش پیشنهادی در مقایسه با روش های دیگر در پیش بینی برچسب های نمونه های جدید، از لحاظ پیچیدگی زمانی و دقت از کارایی بهتری برخوردار است.

کلید واژه- داده های چند برچسبی، دسته بندی داده ها، الگوریتم های دسته بندی

### ۱-مقدمه

تعداد حالت های ممکن برای مجموعه برچسب های متفاوت به صورت نمایی افزایش می یابد و نمی توان از الگوریتم های موجود در داده های تک برچسبی برای حل این مسائل استفاده کرد. زیرا استفاده از دسته بندی های داده های تک برچسبی باعث پیچیدگی زمانی زیادی می شود. روش های یادگیری چند برچسبی سعی در کاهش این هزینه ها دارند. از طرفی برخی از کاربردها همانند توابع تاثیر ژن ها، ذاتا چند برچسبی هستند و حتی با هزینه های بالا هم نمی توان آنها را به مسئله تک برچسبی تبدیل کرد. بنابراین به یادگیری چند برچسبی برای این مسائل ضروری است. در یادگیری چند برچسبی هدف یافتن تابع نگاشتی است که فضای ویژگی های را به فضای مجموعه برچسب ها نگاشت کند [1]. یک رویکرد مناسب برای حل این مشکل استفاده از ارتباط و وابستگی بین برچسب می باشد.

در دسته بندی تک برچسبی هر نمونه با یک برچسب در ارتباط می باشد که این برچسب ویژگی های آن نمونه را مشخص می کند. در دسته بندی چند برچسبی هر نمونه ممکن است با چندین برچسب در ارتباط باشد و اجتماع این برچسب ها ویژگی های نمونه ها را مشخص می کنند به عبارت دیگر در دسته بندی چند برچسبی هر نمونه توسط برداری از برچسب ها مشخص می شود. در بسیاری از کاربردهای واقعی نیز داده ها چند برچسبی می باشند. به عنوان مثال یک فیلم سینمایی ممکن است مربوط به چندین ژانر مختلف باشد. و یا یک خبر مربوط به چندین حوزه مختلف باشد. معمولا برای حل این مسائل از راهکارهای تک برچسبی استفاده می شود که اغلب، این راهکارها، باعث افزایش هزینه های نگه داری و بازیابی می شود و دارای پیچیدگی زمانی زیادی نیز هستند. با افزایش تعداد برچسب ها یادگیری داده های چند برچسبی دچار چالش می شود زیرا با این افزایش،

دسته بندی ها یکی از دو برچسب را به عنوان برچسب آن ها پیش بینی می کند در پایان برای ترکیب این دسته بندی ها از روش رای گیری و رتبه بندی استفاده می شود. یک روش جدید موثر برای رای گیری بر اساس روش وزن دهی Qweighted به نام QWML برای روش فاصله جفتی ارائه شده است [۴].

## ۲-۲- روش های تطبیق الگوریتم

در این دسته از روش ها تلاش می شود که الگوریتم های موجود در دسته بندی داده های تک برچسبی به گونه ای تغییر داده شوند که توانایی دسته بندی داده های چندبرچسبی را داشته باشند.

روش درخت تصمیم: در این روش سعی شده الگوریتم C4.5 به گونه تغییر داده شود تا این الگوریتم برای دسته بندی داده های چندبرچسبی هم مناسب باشد. در الگوریتم C4.5ML- رابطه موجود در الگوریتم C4.5 برای محاسبه این تروپی تغییر یافته و به صورت رابطه (۱) محاسبه می شود [۵].

$$entropy(E) = - \sum_{i=1}^N (P((c_i) \log(c_i) + q(c_i) \log q(c_i)) \quad (1)$$

منظور از E مجموعه نمونه های آموزشی و منظور از نسبت تکرارهای برچسب i ام به مجموع تکرارهای کل برچسب ها می باشد و است.

روش نزدیک ترین همسایگی: در این روش ابتدا با استفاده از فاصله اقلیدسی نزدیکترین همسایه های نمونه جدید پیدا می شود سپس در بین این همسایه ها احتمال این که نمونه جدید هر کدام از برچسب ها را بگیرد محاسبه می شود. سپس این احتمالات را به عنوان احتمال اولیه اینکه یک نمونه جدید با چه کلاس هایی در ارتباط است به یک دسته بند بیز داده می شود و برچسب های خروجی دسته بند بیزین که احتمال بیش ۰.۵ دارند به عنوان برچسب های نمونه جدید در نظر گرفته می شوند [۶].

روش HOMER: این روش یک سلسله مراتب از دسته بندی ها می سازد؛ که هر کدام از این دسته بندی ها با یک مجموعه کوچک از برچسب ها (نسبت به کل برچسب ها) در ارتباط هستند. ایده اصلی این است که الگوریتم دسته بندی داده های چند برچسبی، در مجموعه داده های با برچسب های زیاد، به شکل یک درخت سلسله مراتبی تبدیل شوند. در این درخت هر گره شامل

در ادامه این مقاله در بخش ۲ مروری بر پژوهش های پیشین در حوزه یادگیری چند برچسبی خواهیم داشت در بخش ۳ روش پیشنهادی ارائه می شود در بخش ۴ به ارزیابی و مقایسه نتایج حاصل از پیاده سازی روش پیشنهادی می پردازیم و در بخش ۵ نتایج و جمع بندی ارائه می شود.

## ۲-مروری بر پژوهش های پیشین

انواع پژوهش های انجام شده در حوزه یادگیری داده های چندبرچسبی به سه دسته کلی روش های تطبیق الگوریتم، روش های تغییر مسئله روش های جمعی تقسیم بندی شده است [2].

### ۲-۱- روش های تغییر مسئله

این روشها تلاش می کنند که داده های چندبرچسبی را به گونه ای تغییر دهند که این داده ها تبدیل به داده های تک برچسبی شوند. سپس با استفاده از روش ها و الگوریتم های موجود در حوزه یادگیری تک برچسبی کار یادگیری و پیش بینی برچسب ها را برای این داده ها انجام می شود. که برخی از مهم ترین روش ها در این حوزه عبارتند از:

روش ارتباط دودویی (BR): در این روش به ازای تعداد برچسب ها، دسته بند ساخته می شود. و یادگیری و پیش بینی برای هر کدام از این دسته بندی ها به صورت جداگانه انجام می شود. یعنی به ازای هر برچسب یک دسته بند جداگانه خواهیم داشت.

روش مجموعه توانی (LP): در این روش به ازای هر زیر مجموعه ای از مجموعه توانی برچسب های موجود، یک برچسب جداگانه در نظر گرفته می شود. برای نمونه اگر L برچسب داشته باشیم  $2^L$  مجموعه برچسب خواهیم داشت که هر کدام به عنوان یک برچسب مجزا در نظر گرفته می شود و به این ترتیب داده های چندبرچسبی به داده های تک برچسبی تبدیل می شوند [۳].

روش فاصله جفتی: در این روش به ازای هر جفت برچسب یک دسته بند ساخته می شود. یعنی اگر L برچسب داشته باشیم  $L*(L-1)/2$  دسته بند ساخته می شود. هر دسته بند برای مقایسه دو برچسب به کار می رود. برای داده های آزمایشی هر کدام از این

۲. تعیین کردن تعداد برچسب های هر نمونه

۳. هزینه محاسباتی پایین

روش های تغییر مسئله ( همانند روش ارتباط دودویی، روش مجموعه توانی و روش فاصله جفتی ) نسبت به روش های دیگر در محاسبه احتمال حضور هر برچسب برای نمونه های جدید موفق تر هستند اما کاستی بزرگ این روش ها هزینه زمانی بسیار بالای آنهاست و زمانی که حجم مجموعه داده ها یا تعداد برچسب ها زیاد شود، استفاده از این روش ها در عمل غیر ممکن می شود. از طرفی این روش ها برای مشخص کردن تعداد برچسب های نمونه های جدید هیچ راه حلی ارائه نمی دهند و معمولا تعداد برچسب های نمونه های جدید را به اشتباه تشخیص می دهند.

روش های تطبیق الگوریتم ها ( همانند روش درخت تصمیم، روش نزدیکترین همسایگی، روش HOMER و...) نسبت به دیگر روش ها اغلب دارای پیچیدگی زمانی پایین تری هستند و با بزرگ شدن حجم مجموعه داده ها یا تعداد برچسب ها دچار چالش جدی نمی شوند. اما این روش ها نسبت به دیگر روش ها دارای دقت پایین تری در پیش بینی برچسب ها هستند این روش ها نیز همانند روش های تغییر مسئله هیچ راهکاری برای مشخص کردن تعداد برچسب های نمونه های جدید ارائه نمی دهند. این روش ها معمولاً در تشخیص تعداد برچسب های نمونه های جدید دچار مشکل می شوند.

روش های جمعی رفتارهای متفاوتی دارند، برخی روش های جمعی (همانند روش RAKEL) با وجود اینکه دارای پیچیدگی زمانی پایینی هستند اما در تشخیص احتمال برچسب های هر نمونه، دقت کافی ندارند و بسیار حساس به پارامترهایشان هستند. برخی دیگر از روش های جمعی (همانند روش دسته بندی زنجیره ای) مانند روش های تغییر مسئله دارای پیچیدگی زمانی زیادی هستند و زمانی که حجم مجموعه داده ها یا تعداد برچسب ها زیاد شود استفاده از این روش ها در عمل غیر ممکن می شود. در راستای حل مسئله دسته بندی داده های چند برچسبی، ساختار شکل (۱) برای حل مساله بیان می شود که در ادامه هر کدام از زیر بخش ها با جزئیات بیشتر بیان می شود.

مجموعه ای از برچسب ها می باشد. این درخت  $|L|$  برگ دارد که هر کدام از این برگ ها شامل یک مجموعه مولفه با برچسب جدا می باشد  $\{l_j\}$  و هر گره داخلی شامل مجموعه برچسب هایی است که از اجتماع برچسب های فرزندانش بدست می آید و ریشه شامل همه برچسب ها می باشد. هر گره داخلی از سلسله مراتب نیز شامل یک دسته بند چند برچسبی ( $h_n$ ) می باشد. کار این دسته بند این است که یک یا چند ابر برچسب برای فرزندانش پیش بینی کند [۷].

## ۲-۳- روش های جمعی

در دسته بندی داده ها به کمک روش های جمعی، ابتدا نمونه ها را توسط چندین دسته بند جداگانه، دسته بندی می کنند. سپس برای مشخص کردن برچسب نهایی نمونه ها، برای هر نمونه، بین دسته بندها رای گیری انجام می دهند.

روش K مجموعه برچسب تصادفی (RAKEL): در این روش ابتدا به صورت تصادفی مجموعه برچسب های اولیه به K مجموعه برچسب مساوی افزای می شود. سپس برای هر کدام از این مجموعه برچسب های ایجاد شده، به صورت جداگانه دسته بندی انجام می شود. یعنی برای هر کدام از این مجموعه برچسب ها، تمام مجموعه توانی آن ها را ساخته و کار دسته بندی انجام می شود. در پایان برای مشخص کردن خروجی نهایی یک نمونه جدید، بین خروجی های این K دسته بند رای گیری انجام می شود [۸].

مدل دسته بندی زنجیره ای (CC): در مدل دسته بندی زنجیره ای L دسته بند دودویی وجود دارد. که هر دسته بند دودویی برای یک برچسب است و این دسته بندها با هم یک زنجیر را می سازند. در این زنجیر هر دسته بند مربوط به یکی از برچسب های موجود در مجموعه برچسب ها است [۹].

## ۳- روش پیشنهادی

هدف نهایی در دسته بندی داده های چند برچسبی ارائه الگوریتم هایی است که توانایی پیش بینی برچسب های نمونه های جدید را داشته باشند. بنابراین لازم است الگوریتم های ارائه شده در این راستا سه مسئله مهم زیر را در نظر بگیرند:

۱. مشخص کردن احتمال حضور هر برچسب برای هر نمونه

### ۳-۲- استفاده از احتمال شرطی برچسبها

برای اینکه در پیش بینی برچسب های نمونه های جدید کارایی بهتری حاصل شود می توان از ارتباط و وابستگی بین برچسب ها استفاده کرد. اگر بدانیم که یک نمونه حتما یک برچسب خاص را می گیرد می توان از وابستگی بین این برچسب و دیگر برچسب ها برای دسته بندی داده ها استفاده کرد. به عنوان نمونه، فرض کنید برچسب  $L_1$  جزء برچسب های واقعی یک نمونه است. اکنون اگر بدانیم که برچسب  $L_2$  در هیچ یک از مثال های آموزشی به همراه برچسب  $L_1$  نیامده است و برچسب  $L_3$  همواره با برچسب  $L_1$  رخ داده است می توان گفت، اگر نمونه ای برچسب  $L_1$  را داشته باشد به احتمال بالایی برچسب  $L_3$  را دارد و برچسب  $L_2$  را ندارد. رابطه (۳) احتمال شرطی مربوط به برچسب های  $L_1$  و  $L_2$  را نشان می دهد.

$$P(L_1|L_2) = 0 \quad (3)$$

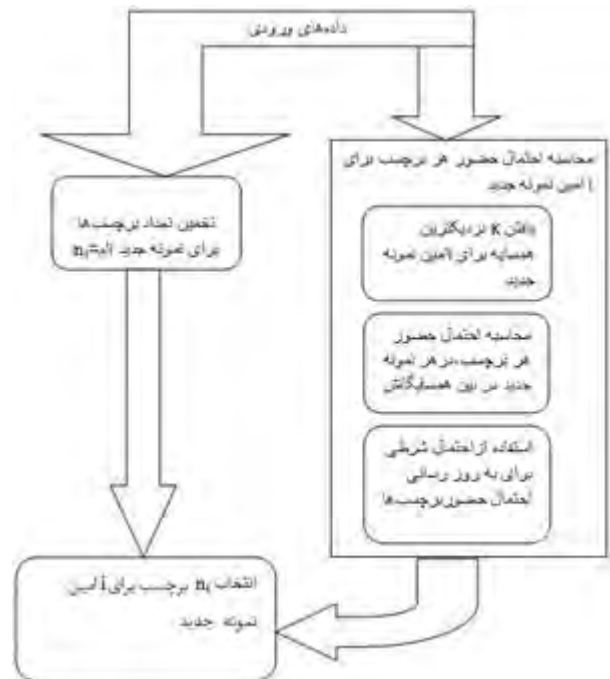
رابطه (۴) احتمال شرطی مربوط به برچسب های  $L_1$  و

$$L_3 \text{ را نشان می دهد.}$$

$$P(L_1|L_3) = 1 \quad (4)$$

نکته اصلی برای استفاده از وابستگی بین برچسب ها در دسته بندی داده های چندبرچسبی این است که ابتدا برای یک نمونه جدید برچسبی یافته شود که با احتمال زیادی مربوط به آن نمونه باشد. سپس برای بقیه برچسب ها که احتمال کمتری دارند از احتمال شرطی استفاده شود و در پایان کار دسته بندی انجام شود. برای یافتن احتمال هر برچسب برای هر نمونه جدید از روش نزدیک ترین همسایگی استفاده خواهد شد.

برای هر نمونه جدید ابتدا با استفاده از روش نزدیکترین همسایگی،  $K$  نزدیکترین همسایه آن پیدا شده سپس در بین برچسب های این  $K$  همسایه، احتمال وقوع هر برچسب محاسبه می شود. یعنی برای هر برچسب تعداد رخداد آن برچسب در بین  $K$  همسایه، بر مجموع رخداد های کل برچسب ها در  $K$  همسایه تقسیم می شود. سپس برای هر نمونه، برچسبی که بیشترین تکرار را داشته باشد به عنوان برچسب اصلی این نمونه در نظر گرفته شده و احتمالات بدست آمده برای برچسب های دیگر از روش نزدیکترین همسایگی با استفاده از احتمال شرطی این برچسب و دیگر برچسب ها به روز رسانی می شود.



شکل ۱: ساختار روش پیشنهادی

### ۳-۱- تخمین تعداد برچسب های هر نمونه

بسیاری از الگوریتم های ارائه شده برای دسته بندی داده های چند برچسبی، اغلب پس از مشخص کردن احتمال حضور هر برچسب در هر نمونه جدید، از یک تابع آستانگی برای تعیین برچسب های نهایی استفاده می نمایند. این تابع آستانگی معمولاً به صورت رابطه (۲) تعریف می شود.

که منظور از  $PL_i$  احتمال حضور برچسب  $i$ ام برای نمونه جدید  $i$ ام و منظور از  $L_i$  برچسب  $i$ ام برای نمونه  $i$ ام است. میزان دقت و کارایی الگوریتم هایی که از رابطه (۳-۱) برای تعیین تعداد برچسب های نهایی نمونه های جدید استفاده می کنند، به میزان زیادی به نحوه انتخاب متغیر آستانگی  $a$  وابسته است. به عنوان نمونه در الگوریتم نزدیکترین همسایگی  $a=0.5$  در نظر گرفته می شود.

به منظور مدیریت چالش مذکور برای تخمین تعداد برچسب های نمونه های جدید، تخمین تعداد برچسب های یک نمونه جدید قبل از دسته بندی به عنوان یک فاز اولیه در نظر گرفته می شود و در آن تخمین تعداد برچسب ها انجام می پذیرد.

### ۳-۳- تغییر معیار فاصله اقلیدسی

هم این است که در یک داده چندبرچسبی ویژگی های موجود در بردار ویژگی های هر نمونه مربوط به چندین برچسب هستند بنابراین یک نمونه با چند برچسب دارای اجتماعی از ویژگی هایی است که هر برچسب دارا است. به عبارت دیگر وجود هر برچسب در مجموعه برچسب های یک نمونه موجب می شود آن نمونه دارای تعدادی از ویژگی ها شود. بنابراین در جدول (۳-۱) می توان گفت (به صورت نادقیق) نمونه های اول و دوم تمام برچسب های نمونه آزمایشی را دارند به اضافه چند برچسب اضافی (که باعث شده است ویژگی های بیشتری را داشته باشند) و نمونه سوم چند مورد از برچسب های موجود در مجموعه برچسب های نمونه آزمایشی را ندارد. پس می توان گفت نمونه های اول و دوم نسبت به نمونه سوم به نمونه آزمایشی نزدیک تر هستند و فاصله ها یکسان نیست.

برای حل این مشکل، برای محاسبه فاصله نمونه آزمایشی با هریک از نمونه های آموزشی، به جای استفاده از فاصله اقلیدسی برای محاسبه فاصله، بردار ویژگی های نمونه ها با هم مقایسه می شوند. به عنوان نمونه فرض کنید یک نمونه جدید، ویژگی f1 را دارد و ویژگی f2 را ندارد. برای پیدا کردن نزدیکترین همسایه های این نمونه جدید در بین داده های آموزشی، به این صورت عمل می شود.

برای هریک از نمونه های آموزشی یک متغیر امتیاز تعریف می شود. که مقدار این امتیاز برابر است با مجموع امتیازهایی که هر یک از ویژگی های نمونه آموزشی به دست آورده اند.

اگر یک نمونه از داده های آموزشی ویژگی f1 را دارا باشد به متغیر امتیاز آن مقدار  $\alpha$  اضافه می شود. اگر یک نمونه از داده های آموزشی یک ویژگی ای که نمونه جدید دارد، را دارا باشد به این معنی است که در این ویژگی، این نمونه آموزشی با نمونه جدید همسو است. بنابراین یک امتیاز مثبت می گیرد که نشان دهنده این است که در این ویژگی این نمونه آزمایشی به نمونه جدید نزدیک است.

اگر یک نمونه از داده های آموزشی ویژگی f2 را نداشته باشد به متغیر امتیاز آن مقدار  $\beta$  ( $\beta < \alpha$ ) اضافه می شود. اگر یک نمونه از داده های آموزشی یک ویژگی ای که نمونه جدید ندارد را نداشته باشد به این معنی است که این دو نمونه به یکدیگر نزدیک هستند ولی نه به اندازه حالت قبل بنابراین امتیازی پایین تر از

در روش دسته بندی داده های چندبرچسبی به روش K نزدیکترین همسایگی، ابتدا نزدیکترین همسایه ها در بین مثال های آموزشی برای نمونه جدید پیدا شده و سپس احتمال وقوع هر برچسب بین k همسایه برای نمونه جدید به عنوان احتمال اولیه بدست می آید. منظور از فاصله بین دو نمونه، این است که ویژگی های این دو نمونه به چه میزان به یکدیگر شبیه هستند. فاصله اقلیدسی از رابطه (۵) برای محاسبه فاصله میان دو نمونه X, Y که هریک سه ویژگی با مقادیر  $x_1, x_2, x_3, y_1, y_2, y_3$  استفاده می کند.

$$\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2} \quad (5)$$

اگر از فاصله اقلیدسی برای یافتن نزدیکترین همسایه برای نمونه جدید موجود در جدول (۱) استفاده شود مشاهده می شود که معیار فاصله اقلیدسی فاصله هر سه نمونه را با نمونه آزمایشی برابر بدست می آورد.

جدول ۱: نمونه های چند برچسبی

نمونه	ویژگی				برچسب			
	F1	F2	F3	F4	L1	L2	L3	L4
نمونه اول	۱	۰	۱	۱	۱	۰	۱	۱
نمونه دوم	۰	۱	۱	۱	۱	۱	۰	۱
نمونه سوم	۰	۰	۱	۰	۱	۰	۰	۰
نمونه آزمایشی	۰	۰	۱	۱	?	?	?	?

دلیل این موضوع این است که فاصله اقلیدسی دو نمونه آموزشی اول، که نسبت به نمونه جدید (نمونه آزمایشی) یک ویژگی اضافی دارند و نمونه آموزشی سوم که نسبت به نمونه جدید یک ویژگی کم دارد برابر است.

نمونه های اول و دوم تمام ویژگی هایی که نمونه آزمایشی دارا می باشد را دارند و نیز یک ویژگی بیشتر دارند ولی نمونه سوم یکی از ویژگی هایی که نمونه آزمایشی دارد، را ندارد. در دسته بندی داده های تک برچسبی این سه نمونه فاصله برابر با نمونه آزمایشی دارند، چون قرار است در نهایت یک برچسب برای نمونه جدید پیش بینی شود، مشکلی بوجود نمی آید. اما در داده های چندبرچسبی تا حدودی شرایط متفاوت است. دلیل آن



میانگین دقت که در روش های مبتنی بر رتبه بندی استفاده می شوند [۲۹] برای ارزیابی استفاده شده اند.

معیار پوشش: این معیار مشخص می کند که به طور میانگین در بردار احتمال برچسب ها (که به صورت نزولی مرتب هستند) چه تعداد برچسب (چه تعداد گام) باید بررسی شود تا برچسب های صحیح نمونه جدید مشخص گردد. هرچه مقدار این معیار کوچک تر باشد به معنای کارایی بهتر در الگوریتم پیش بینی برچسب ها است. رابطه ریاضی این معیار در رابطه (۶) آورده شده است.

$$coverage(f) = \frac{1}{N} \sum_{i=1}^N \max_{l_i \in Y} rank_f(x_i, l) - 1 \quad (6)$$

در جدول (۲) نتایج مربوط به اجرای روش پیشنهادی و دیگر روش ها برای معیار پوشش آورده شده است.

جدول ۳: مقایسه نتایج اجرای روش پیشنهادی و با دیگر روش ها بر اساس معیار coverage.

Data Set	BR	CC	HOMER	MLC4.5	MLKNN	RAKEL	Proposed method
Medical	1.61	1.47	5.32	3.05	2.84	8.52	2.80
Enron	12.5	12.4	24.2	17.0	13.1	.5	13.9
Bibtex	2.9	21.1	65.6	58.0	56.2	DNF	43.6

معیار میانگین دقت: این معیار مشخص می کند که به طور میانگین چه نسبتی از برچسب هایی که رتبه بالاتر از یک برچسب واقعی  $l \in Y_i$  دارند عضو  $Y_i$  هستند. بنابراین هرچه مقدار این معیار بزرگ تر باشد به معنای کارایی بهتر در الگوریتم پیش بینی برچسب ها است. رابطه این معیار در رابطه (۷) آورده شده است.

$$average\ precision(f) = \frac{1}{N} \sum_{l \in Y_i} \frac{L_i}{rank_f(x_i, l)} \quad (7)$$

در جدول (۴) نتایج مربوط به اجرای روش پیشنهادی و دیگر روش ها برای معیار میانگین دقت آورده شده است.

حالت قبل می گیرد. اگر یک نمونه از داده های آموزشی ویژگی  $f_2$  را داشته باشد به متغیر امتیاز آن مقدار  $\gamma$  ( $\gamma < \beta$ ) اضافه می شود. اگر نمونه جدید یک ویژگی را نداشته باشد ولی نمونه های آموزشی داشته باشند این بدان معنی است که این نمونه آموزشی با نمونه جدید همسو نیست. بنابراین امتیاز کمتری خواهد گرفت (امتیاز منفی) تا مشخص کننده این عدم تطابق باشد. در واقع می توان فرض کرد که نمونه آموزشی، برچسبی را دارا است که جز برچسب های اصلی نمونه جدید نیست و به خاطر این برچسب، نمونه آموزشی ویژگی که این نمونه دارا نیست را دارد.

اگر یک نمونه از داده های آموزشی ویژگی  $f_1$  را نداشته باشد به متغیر امتیاز آن مقدار  $\theta$  ( $\theta < \gamma$ ) اضافه می شود. بیشترین فاصله بین یک نمونه جدید و نمونه های آموزشی در یک ویژگی، زمانی حاصل می شود که یک نمونه آموزشی ویژگی ای که نمونه جدید دارد را نداشته باشد. این حالت عکس حالت قبل است، که می توان فرض کرد در این حالت، نمونه جدید یک برچسبی را دارد که نمونه آموزشی ندارد و به دلیل وجود این برچسب، نمونه جدید این ویژگی را دارد و نمونه آموزشی این ویژگی را ندارد. در پایان  $K$  تا از نمونه های آموزشی که بیشترین امتیاز را دارند به عنوان همسایه های نمونه جدید مشخص می شوند.

#### ۴- ارزیابی و مقایسه نتایج

برای ارزیابی روش پیشنهادی از مجموعه داده های موجود در جدول (۲) که شامل سه مجموعه داده enron, bibtex, medical می باشند، استفاده می کنیم.

جدول ۲: مجموعه داده های مورد ارزیابی

نمونه ها	برچسب ها	ویژگی ها	مجموعه داده ها
978	45	1449b	Medical
1702	53	1001b	Enron
7393	159	1836b	Bibtex

از آنجا که روش پیشنهادی مبتنی بر رتبه بندی برچسب ها می باشد. بنابراین از بین معیارهای موجود در حوزه دسته بندی داده های چند برچسبی معیارهای پوشش، خطای همینگ و

موجود مقایسه شده است. در جدول (۶) زمان اجرای الگوریتم های مختلف بر حسب دقیقه آورده شده است.

جدول ۶: زمان مورد نیاز برای اجرای الگوریتم های مختلف بر حسب دقیقه

Data Set	BR	CC	HOMER	MLC4.5	MLKNN	RAKEL	Proposed method
Medical	7.3	11	5.8	1.03	0.4	35.3	2.5
Enron	122	164	60	5.06	3	215	11
Bibtex	3889	4365	1017	190	62	DNF	150

## ۵- نتیجه گیری

در این مقاله یک روش با استفاده از احتمال بین برچسب ها مبتنی بر روش نزدیکترین همسایگی برای دسته بندی داده های چندبرچسبی ارائه شد. ایده اصلی ارائه شده در این مقاله استفاده از احتمال شرطی بین برچسب ها برای دسته بندی داده های چند برچسبی می باشد برای اینکه بتوان از احتمال شرطی بین برچسب ها استفاده کرد باید در گام اول برای هر نمونه یک برچسب را به درست پیش بینی کرد (تا جای ممکن معیار تک خطا را کاهش داد). در این مقاله برای تشخیص این برچسب از روش نزدیکترین همسایگی بر اساس فاصله غیر اقلیدسی استفاده شده است. نتایج آزمایش ها بر روی مجموعه داده ها بیانگر این نکته است که روش پیشنهادی با توجه به پیچیدگی زمانی آن، نسبت به روش های دیگر، در پیش بینی برچسب های نمونه های جدید بهبود داشته است.

همچنین نتایج آزمایشات بیانگر این موضوع می باشد که با انتخاب مقادیر  $\alpha=1$ ,  $\beta=0$ ,  $\gamma=-1$  و  $\theta=-4$  بهترین نتایج برای الگوریتم پیشنهادی حاصل می شود به همین دلیل برای محاسبه معیارهای ارزیابی از این مقادیر استفاده شده است. ویژگی های روش پیشنهادی را می توان به صورت زیر خلاصه کرد.

زمان اجرای روش پیشنهادی برخلاف روش های BR, RAKEL, CC, ECC به تعداد برچسب ها مجموعه داده وابستگی شدیدی ندارد و با افزایش تعداد برچسب ها هزینه زمانی بالایی را تحمیل نمی کند.

روش پیشنهادی نسبت به روش هایی مانند MLKNN, MLC4.5, HOMER که هزینه زمانی قابل قبولی

جدول ۴: مقایسه نتایج اجرای روش پیشنهادی و با دیگر روش ها بر اساس معیار average precision

Data Set	BR	CC	HOMER	MLC4.5	MLKNN	RAKEL	Proposed method
Medical	.896	.901	.786	.823	.784	.676	.801
Enron	.693	.695	.604	.629	.635	.522	.661
Bibtex	.597	.599	.407	.392	.349	DNF	.507

معیار دیگر برای ارزیابی صحت در یادگیری چندبرچسبی، معیار خطای همینگ است. این معیار جز روش مبتنی بر برچسب است که به صورت رابطه (۳-۴) تعریف می شود. هرچه مقدار این معیار کمتر باشد به معنای کارایی بهتری می باشد.

$$Hamming - Loss(D) = \frac{1}{NL} \sum_{i=1}^N |\hat{y}_i \Delta y_i| \quad (8)$$

در جدول (۵) نتایج مربوط به اجرای روش پیشنهادی و دیگر روش ها برای معیار فقدان همینگ آورده شده است. منظور از DNF در این جدول ها این است که الگوریتم مربوطه نیاز به زمان زیادی داشته است.

جدول ۵: مقایسه نتایج اجرای روش پیشنهادی و با دیگر روش ها بر اساس معیار hamming loss

Data Set	BR	CC	HOMER	MLC4.5	MLKNN	RAKEL	Proposed method
Medical	.077	.077	.012	.015	.018	.012	.014
Enron	.045	.064	.051	.053	.051	.045	.048
Bibtex	.012	.012	.014	.016	.014	DNF	.014

اما نکته دیگری که باید در ارزیابی روش پیشنهادی با روش های پیش برسی شود، میزان پیچیدگی زمانی الگوریتم ها است. برخی از الگوریتم ها مانند الگوریتم ارتباط دودویی، اغلب از لحاظ دقت پیش بینی برچسب های نمونه های جدید از دیگر الگوریتم ها کارا تر است اما از لحاظ پیچیدگی زمانی، نیاز به زمان اجرای بالایی دارد. به همین دلیل در ارزیابی روش پیشنهادی با دیگر روش های موجود، علاوه بر مقایسه معیارهای ذکر شده، روش پیشنهادی از نظر میزان زمان مورد نیاز با دیگر الگوریتم های

- [3] G. Tsoumakas and I. Katakis, "Multi label classification An overview," International Journal of Data Mining and Warehousing, Vol. 3, No. 3, pp. 1-13, 2007.
- [4] E. L. Mencia, S. H. Park and J. Furnkranz, "Efficient voting prediction for pairwise multilabel classification," in Proceedings of the 18th European Conference on Machine Learning, pp. 658-668, 2010.
- [5] A. Clare and R. D. King, "Knowledge discovery in multi-label phenotype data," in Proceedings of the 5th European Conference on PKDD, pp. 42-53, 2001.
- [6] M. L. Zhang and Z. H. Zhou, "ML-kNN: a lazy learning approach to multi-label learning," Pattern Recognition, Vol. 40, No. 7, pp. 2038-2048, 2007.
- [7] G. Tsoumakas, I. Katakis and P. Vlahavas, "Effective and efficient multilabel classification in domains with large number of labels," In ECML/PKDD Workshop on Mining Multidimensional Data, 2008.
- [8] G. Tsoumakas and I. P. Vlahavas, "Random k-labelsets: An ensemble method for multilabel classification," In Proceedings of the 18th European Conference on Machine Learning, PP. 406-417, Springer, 2007.
- [9] J. Read, B. Pfahringer and G. Holmes, "Multi-label classification using ensembles of pruned sets," in Proceedings of the 8th IEEE International Conference on Data Mining, pp. 995-1000, 2008.

دارند تقریباً دارای دقت بالاتری در پیش‌بینی برچسب‌های نمونه‌های جدید است.

در روش پیشنهادی سعی شده است از ارتباط بین برچسب‌ها برای تشخیص برچسب‌های نمونه‌های جدید استفاده شود. که در روش‌های قبلی از این ویژگی استفاده نشده است. روش پیشنهادی دارای دو فاز مجزا است. در فاز نخست تعداد برچسب‌ها برای هر نمونه جدید تخمین زده می‌شود و در فاز دوم احتمال اینکه هر نمونه چه برچسب‌هایی را خواهد گرفت محاسبه می‌شود.

#### مراجع

- [1] J. Read, "Scalable Multi-label Classification," Ph.D. These, Department of Computer Science, University of Waikato, 2010.
- [2] G. Madjarov, D. Kocev and D. Gjorgjevikj, "an extensive experimental comparison of methods for multi-label learning," pattern Recognition," Vol. 45, No. 9, pp. 3084-3104, 2012.