

## مروری بر الگوریتم‌های یادگیری دسته جمعی جهت تشخیص هرزنامه

\*علی حلافی<sup>۱</sup>، علی هارون آبادی<sup>۲</sup>، علی چوبین<sup>۳</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد، دانشگاه آزاد اسلامی واحد اهواز، Hallafi\_a@yahoo.com

<sup>۲</sup> استادیار، گروه نرم افزار، دانشگاه آزاد اسلامی واحد تهران مرکزی، تهران، a.haronabadi@gmail.com

<sup>۳</sup> دانشجوی کارشناسی ارشد، دانشگاه آزاد اسلامی واحد اهواز، alichooabin@gmail.com

### چکیده

دردنیای امروز اینترنت یکی از مهمترین روش‌های ارتباطی است، ونقش مهمی در زندگی امروز ایفا می‌کند. پست الکترونیک نیز یکی از مهمترین وگسترده ترین خدماتی است که اینترنت در اختیار کاربران قرار می‌دهد. این خدمت بنا به دلایلی از جمله سادگی، هزینه کم بسیار مورد توجه قرار گرفته است. همچنین عده‌ای عوامل سودجو از محیط مجازی و این خدمت برای تبلیغ کالا و یا مسائلی دیگر به فکر سوء استفاده هستند، بطوری که در راستای این عمل یکسری نامه‌های نامعتبر به روزانه به طور ناخواسته وارد پست الکترونیکی ما می‌شوند که از آنها به نام هرزنامه یا نامه‌های تجاری ناخواسته یاد می‌شود که در حال افزایش هستند. راه‌های متعدد و فراوانی جهت مقابله با این پدیده ارائه شده است که هر کدام یکسری مزایا و معایب دارند، شناخت صحیح عملکرد هر کدام از این الگوریتم‌ها و میزان دقت در تشخیص و فیلتر کردن نامه‌های الکترونیکی نامعتبر کاری ضروری می‌باشد. ما در این جستار به الگوریتم‌های یادگیری جمعی (به جای الگوریتم‌های منفرد) و معرفی آنها پرداخته‌ایم.

**کلمات کلیدی:** یادگیری، الگوریتم‌های دسته جمعی، تکنیک‌های فیلترینگ و هرزنامه

### مقدمه

هرزنامه<sup>۱</sup> یا نامه‌های ناخواسته<sup>۲</sup> به نامه‌های الکترونیکی نامعتبری اطلاق می‌شود که به شکل چشمگیری در حال افزایش بوده و به جرأت می‌توان گفت که بیش از نیمی از فضای پست‌های الکترونیکی ما را به خود جای داده است که مسائل و مشکلات زیادی را برای کاربران محیط مجازی به وجود آورده که عبارتند از: اتلاف فضای ذخیره سازی، هدر رفت پهنای باند، ضایع شدن حداقل زمان ممکن کاربران جهت پاکسازی پست الکترونیکی که معمولاً به صورت مرتب این عمل انجام می‌شود، از جنبه‌های دیگر نیاز مشکلات ناخواسته‌ای را منجر شده که به صورت اشاره‌ای عبارتند از کلاهبرداری اینترنتی، لو رفتن اطلاعات حساب بانکی و از بُعد دیگر داشتن ضمیمه‌های همچون ویروس، تروجان و موارد این چنین را شامل می‌شود که می‌بایست نکات و هشدارهای امنیتی در این حوزه را شناخت و سعی در بهره‌گیری از امکانات موجود (نرم‌افزاری و سخت‌افزاری) جهت داشتن پست الکترونیکی عاری از هرزنامه یا با حداقل آن را داشت. در این جستار ما می‌خواهیم مروری بر روش‌های یادگیری دسته‌جمعی که یکی از کاربردهای آن جهت تشخیص و فیلتر کردن هرزنامه به صورت خودکار را مورد بررسی و بازبینی قرار دهیم.

### تاریخچه

نامه‌های الکترونیکی ناخواسته تا کنون ۳ دوره تاریخی را پشت سر گذاشته‌اند:

دوره نخست: این دوره که از سال ۱۹۸۷ تا اواسط ۱۹۹۰ به طول انجامید سال‌های اولیه حضور نامه‌های نامعتبر در دنیای اینترنت بوده است. در این دوره هرزنامه‌ها به صورت دستی فرستاده می‌شدند از این رو ارسال توده‌ای عظیمی از نامه‌های نامعتبر کاری دشوار و نیاز به نیروی انسانی نسبتاً زیاد بود.

دوره دوم: این دوره مربوط به شروع دوران حضور ماشین‌ها و نرم‌افزارها در فرستادن و پالایش هرزنامه‌ها است. این دوره با ایجاد یک ماشین هرزنامه‌نویس

در سال ۱۹۹۴ آغاز شد، به دنبال آن در سال ۱۹۹۷ اولین نرم‌افزار پالایش هرزنامه‌ها طراحی شد [۱].

دوره سوم: این دوره مربوط به جنگ نرم‌افزارها برای فرستادن و پالایش هرزنامه‌ها می‌باشد. این دوره تقریباً از سال ۲۰۰۲ و با انتشار مقاله‌ی "تدبیری

برای هرزنامه" توسط پل گراهام آغاز شد، در این مقاله استفاده از یادگیری ماشین و تفکیک‌کننده‌های آماری برای پالایش هرزنامه‌ها پیشنهاد شد. این دوره تاکنون نیز ادامه دارد که در طول آن نرم‌افزارهای ارسال هرزنامه‌ها قدرتمند شدند، و هرزنامه‌ها را با سرعت و گوناگونی فراوان تولید می‌کنند، و در مقابل نرم-

<sup>1</sup> Spam

<sup>2</sup> Unwanted Email

افزارهای پالایش هرزنامه‌ها به رویکردهای یادگیری ماشین استفاده از دسته‌بندی آماری و ترکیبی روی آورده است. این دوره اهمیت بیشتری نسبت به بقیه دوره دارد [۳ و ۲].

## انواع هرزنامه

به جهت وسعت و گستردگی که هرزنامه‌ها دارند، طبقه بندی خاصی ندارند و می‌توان را برحسب معیارهای متفاوتی دسته بندی کرد. برای مثال فرستادن هرزنامه دلایل متعددی می‌تواند داشته باشد از جمله دلایل عقیدتی، کلاهبرداری اطلاعاتی و تبلیغاتی. بر همین اساس هرزنامه‌ها به ۳ دسته زیر تقسیم می‌شوند [۴].

- هرزنامه تبلیغاتی<sup>۳</sup>: این نوع هرزنامه‌ها برای تبلیغ و معرفی کالا و خدمات به کار می‌روند و یکی از مهمترین علل گستردگی آن هزینه پایین تبلیغات نسبت به روش‌های دیگر است.
- هرزنامه‌های مالی<sup>۴</sup>: این نوع هرزنامه‌ها با سناریوهای مختلف سعی در فریب کاربران و دریافت پول از آن‌ها را دارد، بعنوان مثال به کاربر گفته می‌شود که در فلان قرعه‌کشی برنده شده‌اید و می‌بایست برای دریافت جایزه هزینه‌ی را بپردازید.
- هرزنامه‌های کلاهبرداری<sup>۵</sup>: هدف این هرزنامه‌ها دسترسی به اطلاعات شخصی و محرمانه کاربران مانند نام کاربری و رمز عبور کارت اعتباری آن‌ها است.

## نیاز به فیلترینگ هرزنامه

گرچه اینترنت روش‌های ارتباطی جدیدی را عرضه کرده است ولی سیل هرزنامه‌ها و هرزنامه‌نویس روز به روز در حال افزایش و تغییر ماهیت هستند. که نیازمند راه‌های فیلترینگ هوشمند و خودکار همچنان ضروری بنظر می‌رسد. یادگیری ماشین: یادگیری ماشین یکی از حوزه‌های مهم هوش مصنوعی است. الگوریتم‌هایی در این حوزه قرار دارند که قابلیت یادگیری را داشته باشند، که باعث افزایش کارایی خود در واحد زمان را داشته باشند [۵].

انواع فیلترینگ‌های هرزنامه: در حال کلی فیلترینگ هرزنامه را می‌توان به دو دسته تقسیم کرد:

۱. مبتنی بر روش‌های یادگیری ماشین: در پالایش (فیلترینگ) هرزنامه‌ها از الگوریتم‌های یادگیری ماشین نیز می‌توان بهره برد. بدین صورت که هدف این الگوریتم‌ها ایجاد تمایز بین هرزنامه‌ها و نامه‌های معتبر می‌باشد. این روش‌ها توانایی استخراج دانش از یک مجموعه داده یا مجموعه مستندات را دارا می‌باشند و در حقیقت دانشی را که از یک مجموعه داده استخراج می‌کنند برای شناسایی هرزنامه‌های جدید مورد استفاده قرار می‌دهند.

این الگوریتم‌ها دسته‌بندی مختلفی دارند که به اختصار توضیح می‌دهیم:

- ۱.۱ یادگیری با مربی: منبع اصلی این دسته‌بندی داده‌های آموزشی برچسب دار می‌باشد. که دارای دو فاز است فاز آموزش<sup>۶</sup> و فاز آزمون<sup>۷</sup>.
- ۱.۲ یادگیری بدون مربی: تفاوت اصلی این الگوریتم‌ها با الگوریتم‌های یادگیری با مربی در فاز آموزش است، در این الگوریتم‌ها داده‌ای برچسب‌داری برای آموزش وجود ندارد تا با استفاده از برچسب آن‌ها یک مدل دسته بندی بسازد، بلکه این الگوریتم‌ها با توجه به شباهت داده‌ای جدید آن‌ها را به گروه‌های تقسیم می‌کند. مسائل خوشه بندی<sup>۸</sup> از جمله زیر مجموعه‌های یادگیری بدون مربی محسوب می‌شود. البته علاوه بر دو دسته فوق، دسته‌ی دیگری وجود دارد که بعنوان یادگیری نیمه نظارتی<sup>۹</sup> شناخته می‌شود.

۲. مبتنی بر روش‌های عدم یادگیری ماشین: این فیلترینگ‌ها همانند فیلترینگ‌های مبتنی بر قانون، از الگوها و قوانین ثابتی جهت تشخیص هرزنامه‌ها استفاده می‌کنند. پیاده‌سازی اینگونه فیلترینگ‌ها اصولاً کاری ساده است و توانایی تشخیص هرزنامه‌ها پایین است. روش‌هایی مانند لیست سیاه و سفید از دسته فیلترینگ‌ها می‌باشند [۶].

<sup>3</sup> Advertisement

<sup>4</sup> Financial

<sup>5</sup> Phishing

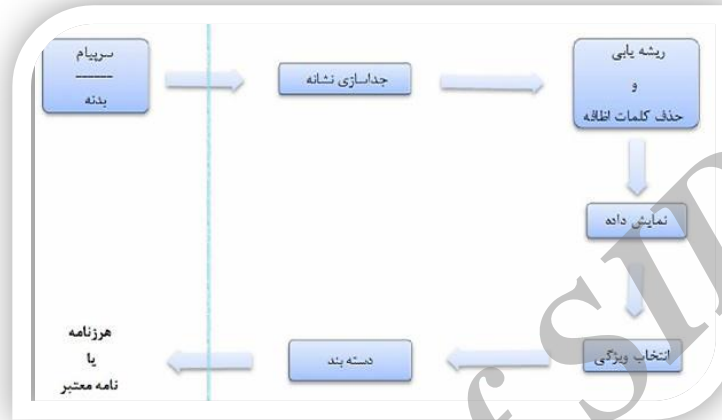
<sup>6</sup> Train

<sup>7</sup> Test

<sup>8</sup> Clustering

<sup>9</sup> Semi Supervised

مراحل اصلی فیلترینگ هرزنامه: مراحل اصلی فیلترینگ هرزنامه که مبتنی بر محتوا است در شکل (۱) آمده است.



شکل (۱) مراحل اصلی فیلترینگ

۱. جداسازی نشانه<sup>۱۰</sup>: استخراج کلمات از بدنه‌ی نامه
۲. حذف حروف اضافه: حذف حروفی که اغلب در بسیاری از متون و به طور مکرر دیده می‌شوند از این رو دارای ارزش محتوایی نیستند.
۳. ریشه‌یابی: باز گرداندن کلمات به ریشه‌ی اصلی خود.
۴. نمایش: تبدیل کردن مجموعه کلمات به فرم مشخص و مورد نیاز برای الگوریتم مورد نظر، مانند ساخت بردار ویژگی با استفاده از روش کیسه‌ی از کلمات.
۵. انتخاب ویژگی: انتخاب زیر مجموعه‌ی از کلمات نشان داده شده، که حاوی اطلاعات مفیدتری هستند. به عبارت دیگر حذف ویژگی‌های نامناسب از بردار ویژگی نشان داده شده در مرحله قبل.

#### الگوریتم‌های یادگیری ماشین در فیلترینگ هرزنامه

از آنجا که فیلترینگ هرزنامه را می‌توان گونه‌ای از دسته‌بندی متون به حساب آورد، بسیاری از الگوریتم‌های یادگیری ماشین که در دسته‌بندی متون کاربرد دارند، در شناسایی هرزنامه نیز می‌توانند مفید باشند. برخی از این الگوریتم‌ها عبارتند از درخت تصمیم، شبکه عصبی، بیزین، ماشین بردار پشتیبان، یادگیری جمعی و غیره می‌باشند. که در ادامه به معرفی الگوریتم‌های یادگیری جمعی می‌پردازیم.

#### الگوریتم‌های مبتنی بر یادگیری جمعی

الگوریتم‌های یادگیری جمعی<sup>۱۱</sup> که به نام‌های مختلفی چون ماشین‌های رایزن<sup>۱۲</sup>، ترکیبی از چند دسته‌بند، مخلوطی از کارشناسان نیز شناخته می‌شوند. از اوایل دهه‌ی ۷۰ وارد پژوهش‌های پیرامون دسته‌بندی شدند ای الگوریتم‌ها جزء دسته‌بندی "یادگیری با مربی" هستند. تفاوت اصلی این الگوریتم‌ها با روش‌های دیگر این است که به جای استفاده از یک دسته‌بند<sup>۱۳</sup>، از چند دسته‌بند<sup>۱۴</sup> که به آن‌ها دسته‌بندهای پایه گفته می‌شود، برای دسته‌بندی استفاده می‌کنند. که در نهایت این دسته‌بندها با یکدیگر به روشی خاص مانند رأی‌گیری، میانگین‌گیری، انتخاب بهترین فرد ترکیب می‌شوند [۷].

#### ساختار یادگیری جمعی

<sup>10</sup>

<sup>11</sup> Ensemble Learning

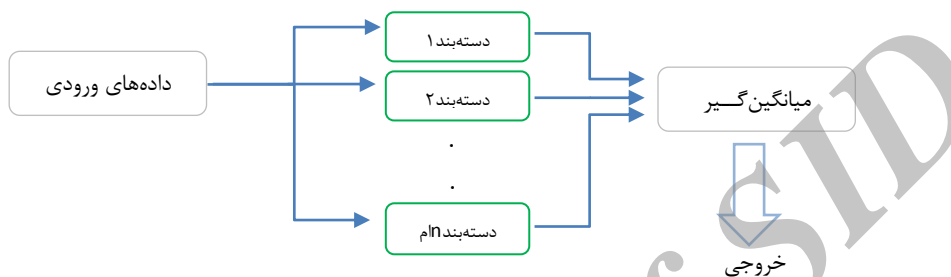
<sup>12</sup> Committee Machine

<sup>13</sup> Single Classifier

<sup>14</sup> Multiple Classifier

الگوریتم‌های یادگیری جمعی ساختارهای گوناگونی دارند و می‌توان آن‌ها را براساس معیارهای گوناگون به دسته‌های مختلفی تقسیم کرد، این الگوریتم‌ها از لحاظ تنوع الگوریتم‌های یادگیری پایه به دو دسته همگن و ناهمگن تقسیم می‌شوند، در یادگیری جمعی همگن تمامی دسته‌بندی‌های پایه انواع الگوریتم‌های مبتنی بر یادگیری یکسان، مانند بیزین ساده، استفاده می‌کنند. در مقابل، یادگیری جمعی ناهمگن شامل تعدادی دسته‌بندی پایه است که از الگوریتم‌های یادگیری متفاوتی بهره می‌برند. به عنوان مثال یکی از دسته‌بندی‌های پایه از ماشین بردار پشتیبان، یکی از K-نزدیکترین همسایه و دیگری از بیزین ساده استفاده می‌کنند [۸].

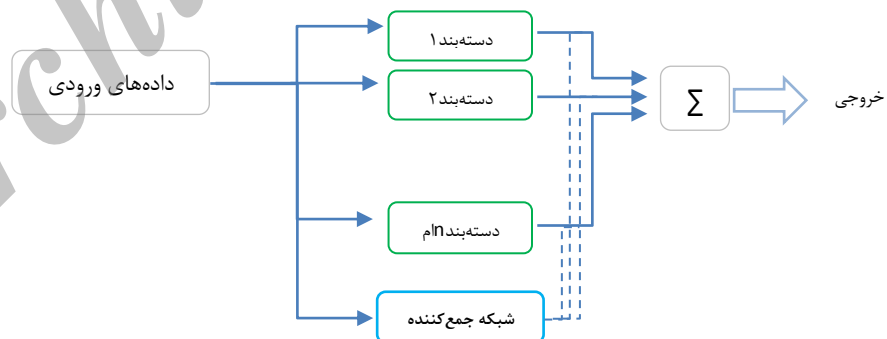
الگوریتم‌های یادگیری جمعی همچنین بر اساس چگونگی ترکیب دسته‌بندی‌های پایه به دو دسته ایستا و پویا تقسیم می‌شوند. در الگوریتم‌های ایستا هیچ‌گونه سیگنال ورودی در ترکیب دسته‌بندی‌ها استفاده نمی‌شود. از جمله معمول‌ترین روش‌ها روش میانگین‌گیری<sup>۱۵</sup> اشاره کرد. در این روش خروجی نهایی دسته‌بندی به صورت ترکیب خطی از خروجی عوامل یادگیر به دست می‌آید. شکل (۳) ساختار میانگین‌گیری دسته‌جمعی را نشان می‌دهد.



شکل (۳) ساختار میانگین‌گیری دسته‌جمعی [۸]

سازوکارهای پویا شامل الگوریتم‌های هستند که یک سیگنال ورودی در فرآیند ترکیب خروجی دسته‌بندی‌های پایه تأثیر دارد. در این الگوریتم‌ها علاوه بر عامل یادگیر قسمتی وجود دارد که در مورد چگونگی ترکیب خروجی عوامل یادگیر می‌باشد، به عبارت دیگر این بخش از الگوریتم با استفاده از ورودی و نیز اطلاعات مربوط به هر دسته‌بندی پایه، وزن هر دسته‌بندی را در جواب نهایی مشخص می‌کند.

روش "مخلوطی از کارشناسان" یا همان "پیش‌بینی با نظر خبرگان" از جمله معروف‌ترین روش‌های یادگیری جمعی پویا محسوب می‌شود. در این روش‌ها علاوه بر عوامل یادگیر پایه، یک واحد به نام شبکه‌ی جمع‌کننده<sup>۱۶</sup> نیز وجود دارد که به عنوان عامل تنظیم‌کننده‌ی تأثیر هر عامل یادگیر در جواب نهایی عمل می‌کند. به عبارتی در این روش‌ها دو مدل وجود دارد: یک یادگیری با مربی که عوامل یادگیر پایه از آن بهره می‌برند، و دیگری یادگیری بدون مربی که در بخش شبکه جمع‌کننده انجام می‌شود. در شکل (۴) این ساختار نشان داده شده است.



شکل (۴) ساختار ترکیبی از کارشناسان [۸]

## مراحل ایجاد یادگیری جمعی

- همانطور که پیش از این نیز گفته شد، یادگیری جمعی شامل مجموعه‌ای از دسته‌بندی‌های پایه و یک روش برای ترکیب نتایج دسته‌بندی‌ها می‌باشد.
- بنابراین برای ایجاد یک یادگیر جمعی دو مرحله وجود دارد.
- طراحی دسته‌بندی‌های پایه

<sup>15</sup> Averaging

<sup>16</sup> Gating Network

- انتخاب روشی برای ترکیب دسته‌بندهای پایه

## ۱. الگوریتم‌های یادگیری جمعی

در ادامه استفاده از داده‌های آموزشی در یادگیری با ساختار همگن اشاره خواهد شد که بعنوان نمونه الگوریتم‌های بگینگ<sup>۱۷</sup> و بوستینگ<sup>۱۸</sup> مورد بررسی قرار می‌گیرند که خود پایه‌ی بسیاری از الگوریتم‌ها هستند.

### ۱.۱. بگینگ

این الگوریتم یکی از ساده‌ترین الگوریتم‌های موجود در حوزه‌ی یادگیری جمعی محسوب می‌شود، که به آسانی پیاده‌سازی می‌شود و با این حال در آزمایش‌ها کارایی مناسبی از خود نشان داده‌است. بحث "گوناگونی" که در ایجاد دسته‌بندها به عنوان عامل اساسی در ساخت یک الگوریتم جمعی شناخته می‌شود، در این الگوریتم با استفاده از تولید زیرمجموعه‌های تصادفی همراه با جایگزینی<sup>۱۹</sup> از کل مجموعه داده آموزش، برای هر دسته‌بند ایجاد می‌شود. پس هر دسته‌بند بر روی یک مجموعه‌ی خاص آموزش می‌بیند، سپس این دسته‌بندها با استفاده از روش میانگین‌گیری و یا رأی اکثریت با یکدیگر ترکیب می‌شوند تا جواب نهایی را تولید کنند. این روش بیشتر زمانی مناسب است که حجم داده‌های مورد استفاده برای آموزش اندک باشد [۲]. نکته‌ی دیگری که در این الگوریتم حائز اهمیت است و در واقع یک نکته‌ی مثبت تلقی می‌شود این است که این الگوریتم توانایی ساخت دسته‌بندها به صورت با یکدیگر را دارد. چرا که دسته‌بندها به یکدیگر وابسته نیستند و هر یک بر روی زیر مجموعه‌ای تصادفی از کل فضای نمونه آموزش می‌بیند. از جمله الگوریتم‌هایی که براساس این ساختار طراحی و ارائه شده‌اند می‌توان به الگوریتم‌های Random Forest و Rotation Forest اشاره کرد.

### ۱.۲. بوستینگ

این الگوریتم جزء الگوریتم‌های قوی در زمینه‌ی یادگیری جمعی به حساب می‌آید. در این روش همانند روش بگینگ، تعدادی دسته بند وجود دارند که هر کدام با استفاده از یک مجموعه داده که به صورت تصادفی از مجموعه داده‌ای اصلی انتخاب می‌شوند آموزش دیده‌اند و با رأی‌گیری اکثریت با یکدیگر ترکیب شده‌اند. اما نکته‌ی اصلی در این الگوریتم این است که در این روش دسته‌بندها نمی‌توانند به صورت موازی و مستقل از هم آموزش ببینند بلکه به ترتیب ویکی پس از دیگری آموزش می‌بینند. در واقع در این روش‌ها اصولاً توزیع تصادفی که برای یک دسته‌بند استفاده می‌شود به گونه‌ای است که در آن احتمال انتخاب داده‌هایی که توسط دسته‌بندهای قبلی نادرست برچسب‌گذاری شده‌اند، بیشتر خواهد بود. بدین ترتیب می‌توان گفت که دسته‌بندهای انتهایی متخصص داده‌های سخت‌تر خواهند بود. الگوریتم AdaBoost از جمله معروف‌ترین الگوریتم‌ها با ساختار بوستینگ می‌باشد.

#### ۱.۲.۱. Ada-Boosting مبتنی بر بوستینگ

در سال ۱۹۹۶ توسط Freund, Schapire خلق شد. دارای دو رویکرد زیر می‌باشد:

- ✓ انتخاب نمونه‌ها براساس خطای دسته‌بند قبلی (دارای رواج بیشتر)
- ✓ خطای وزن‌دار از مواردی که به اشتباه دسته‌بندی شده‌اند بالا است اما نه برای کلیه‌ی الگوریتم‌ها (تمام موارد تعبیه شده‌اند، اما وزن‌ها متفاوت است)

#### مراحل Ada-Boosting

- تعریف  $\epsilon_k$  بعنوان مجموع احتمالات برای نمونه‌های که به اشتباه دسته‌بندی شده‌اند برای دسته‌بند  $C_k$
- ضرب احتمالات مواردی که بصورت نادرست دسته‌بندی شده‌اند بوسیله:

$$\beta_k = (1 - \epsilon_k) / \epsilon_k$$

- نرمال کردن مجدد احتمالات
- ترکیب دسته‌بندها  $C_1, \dots, C_k$  با استفاده از رأی‌گیری وزن‌دار مکانی که  $C_k$  بصورت  $\log(\beta_k)$  وزن دهی شده است.

#### ۲.۲.۱. Arcing مبتنی بر بوستینگ

در سال ۱۹۹۶ توسط Breiman ابداع شد.

مراحل کار:

- برای آمین نمونه در مجموعه آموزش،  $m_1$  به تعداد دفعاتی که به اشتباه دسته‌بندی بوسیله  $K$  دسته‌بند قبلی ارجاع می‌کند.

<sup>17</sup> Bagging

<sup>18</sup> Boosting

<sup>19</sup> Subsampling with Replacement

• احتمال  $P_i$  انتخاب نمونه  $i$  در دسته بند بعدی برابر:

$$P_i = \frac{1+m_i^4}{\sum_{j=1}^N 1+m_j^4}$$

اندازه گیری تجربی

مقایسه تجربی:

- ۲۳ مجموعه داده از مخازن UCI
- ۱۰ بار اعتبار سنجی متقاطع
- پس انتشار شبکه عصبی
- درخت های طبقه بندی
- منفرد و ساده (شبکه های عصبی چندتایی با وزن های مقداردهی اولیه متفاوت)، بگینگ، Ada-boost و Arcing
- # از دسته بند بصورت جمعی
- "دقت" به معیار کارایی

جدول ۱: مقایسه بر اساس دو دسته بند

شبکه عصبی				دسته بند مبتنی بر درخت			
Ada	Arcing	بگینگ	ساده	Ada	Arcing	بگینگ	
.85	.87	.88	1	.37	.38	-.10	شبکه عصبی ساده
.78	.78	1	.88	.35	.35	-.11	شبکه عصبی بگینگ
.99	1	.78	.87	.60	.61	.14	شبکه عصبی Arcing
1	.99	.78	.85	.63	.62	.17	شبکه عصبی Ada
.17	.14	-.11	-.1	.69	.68	1	بگینگ مبتنی بر درخت دسته بند
.62	.61	.35	.38	.96	1	.68	Arcing مبتنی بر درخت دسته بند
.63	.60	.35	.37	1	.96	.69	Ad مبتنی بر درخت دسته بند

نتایج :

- دسته بند جمعی عموماً بهتر از دسته بند منفرد است.
- دسته بند دسته جمعی با شبکه عصبی ها و درخت قویاً وابسته هستند.
- Ada بوستینگ و Arcing قویاً وابسته هستند حتی در مقابل الگوریتم های متفاوت (بوستینگ بیشتر مبتنی بر مجموعه داده است تا نوع الگوریتم دسته بند)
- عموماً شبکه های عصبی بهتر از درخت هستند.

نتیجه گیری:

- بوستینگ دارای عملکرد ضعیف با نويز بیشتر است.

منابع مورد استفاده:

- [۱] Y. Yang and J. O. Pedersen, "A comparative Study on Feature selection in Text Categorization", In Proceedings of 14th International Conference on Machine Learning, pp. ۴۲۰-۴۱۲, ۲۰۰۸.
- [۲] L. S. Larkey and W. B. Croft, "Combining Classifiers in Text Categorization", In Proceedings of 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieve, pp. 289-297, 1996.
- [۳] D. Opitz, "Feature Selection for Ensembles", In Proceedings of 16th National Conference on Artificial Intelligence, pp. 379-384, 1999.
- [۴] G. P. C. Fung, J. X. Yu, H. Wang, D. W. Cheung and H. Liu, "A Balanced Ensemble Approach to Weighting Classifiers for Text Classification", In Proceedings of the 6th International Conference on Data Mining (ICDM '06), pp. 869-873, 2006.
- [۵] K. Wood, W. P. Kegelmeyer and K. Bowyer, "Combination of Multiple Classifiers using Local Accuracy Estimates", IEEE Transaction on Pattern Analysis and Machine, vol. 19x, pp. ۴۱۰-۴۰۵, ۱۹۹۹.
- [۶] E. M. D. Santos, R. Sabourin and P. Maupin, "A Dynamic Overproduce-and-Choose Strategy for the Selection of Classifier Ensembles", Pattern Recognition, vol. 41, pp. 2993-3009, 2008.
- [۷] C. Schaffer, "Selection a Classification Method by Cross-Validation", Machine Learning, vol. 13, pp. 135-143, 1993.
- [۸] G. Brown, "Ensemble Learning", In: Encyclopedia of Machine Learning, Springer Press, Heidelberg, 2010.
- [۹] N. Littlestone and M. K. Warmuth, "The Weighted Majority Algorithm. Information and Computation", 108(2):212- 261, 1994.

Archive of SID