

فهرست مطالب

۲	مقدمه	۱
۴	معرفی تئوری اطلاعات	۲
۴	تئوری اطلاعات	۱-۲
۹	فرآیندهای تصادفی	۲-۲
۱۳	فرآیندهای مارکف	۳-۲
۱۴	رابطه بین دسته‌بندی و Entropy	۴-۲
۱۵	دسته‌بندی داده‌ها	۳
۱۵	تعریف ریاضی مساله دسته‌بندی داده‌ها	۱-۳
۱۸	درخت‌های تصمیم	۲-۳
۱۹	شبکه‌های عصبی	۳-۳
۲۱	Naive Bayes	۴-۳
۲۶	بررسی تخمین‌های پارامترهای Naive Bayes	۴
۲۶	پایداری	۱-۴
۲۷	جابه‌جایی	۲-۴
۲۷	انحراف از معیار	۳-۴
۲۸	گسترش Naive Bayes	۴-۴
۲۹	روش دسته‌بندی داده‌ها بر اساس تئوری اطلاعات	۵
۲۹	علت نیاز به روش دیگر هنگامی که روش Bayes وجود دارد	۱-۵
۳۲	مدل احتمالات	۲-۵
۳۹	چگونگی محاسبه Entropy برای چند فرآیند خاص	۳-۵
۴۱	رابطه بین دسته‌بندی و Entropy	۴-۵
۴۱	ایده کلی روش	۵-۵
۴۴	اثبات هم‌ارزی روش مبتنی بر تئوری اطلاعات با Naive Bayes	۶-۵
۵۳	کارآیی	۷-۵
۵۵	یک مثال	۶
۵۵	معرفی مساله	۱-۶
۵۶	تبدیل مساله به مدل مارکف مربوطه	۲-۶
۶۰	بحث در مورد تخمین‌ها	۳-۶
۶۲	نتایج حاصل	۴-۶
۶۳	یک مثال	۷
۶۳	معرفی مساله	۱-۷
۶۳	تبدیل مساله به مدل مارکف مربوطه	۲-۷
۶۴	نتایج حاصل	۳-۷
۶۷	کارهای آتی	۸

امروزه دسته‌بندی داده‌ها^۱ به صورت وسیعی در علوم کامپیوتر^۲ و هوش مصنوعی^۳ کاربرد دارند. مساله دسته‌بندی داده‌ها معمولاً به صورت تشخیص یک دسته مناسب از بین چند دسته موجود برای یک داده ورودی تعریف می‌شود که در آن دسته‌های موجود بر اساس تعدادی از اعضای موجود در آن دسته (مجموعه آموزش) توصیف شده‌اند. در حالت کلی نمی‌توان یک روش برای دسته‌بندی ارائه داد که بتواند داده‌ها را همیشه به صورت درست دسته‌بندی نماید، اما می‌توان روش‌هایی بهینه و یا نیمه‌بهینه برای این کار ارائه داد. روش‌های مختلفی برای دسته‌بندی داده‌ها مطرح شده است که بسیاری از این روش‌ها، ابتکاری^۴ هستند. روش‌های دسته‌بندی داده‌ها را به دو دسته کلی روش‌های هوش مصنوعی و روش‌های آماری تقسیم می‌کنند که روش‌های هوش مصنوعی بیشتر به دنبال استفاده از تکنیک‌های هوش مصنوعی برای حل این مساله هستند و روش‌های آماری بیشتر تلاش در جهت مدل‌سازی مساله توسط یک فرآیند تصادفی دارند و سپس از ابزارها و روش‌های موجود در آمار و احتمال برای حل مساله استفاده می‌کنند. روش Naive Bayes یکی از ساده‌ترین روش‌ها از نظر پیاده‌سازی است که در موارد خاص دارای بهترین کارایی است و می‌توان از آن برای مسائل پیچیده‌تر (مسائلی که در شرط‌های لازم برای کارایی Naive Bayes نمی‌گنجد) نیز استفاده کرد. Naive Bayes را به طور کلی به عنوان یک روش آماری می‌شناسند.

تئوری اطلاعات، شاخه‌ای از ریاضیات و علوم ارتباطی است که به بررسی مساله‌هایی از قبیل مقدار آشفتگی در پیام‌های ارسالی از یک منبع، خطاهای ممکن حین انتقال پیام از مبدا به مقصد و همچنین فشرده‌سازی داده‌ها می‌پردازد. این شاخه از علم با انتشار مقاله‌ای از Shannon در سال ۱۹۴۸ رسماً متولد شد.

در مساله‌های دسته‌بندی می‌توان مساله را به نوعی منبع ارسال‌کننده پیام تشبیه کرد که تعدادی از پیام‌های آن منبع را در اختیار داریم و می‌خواهیم در مورد امکان ارسال شدن یک پیام توسط آن منبع نظر بدهیم، یعنی پیش‌بینی کنیم که آیا منبع اطلاعاتی می‌تواند یک پیام مشخص را ارسال کند یا خیر. روش‌های مختلفی برای ایجاد یک منبع اطلاعاتی با استفاده از تعدادی از پیام‌های آن منبع وجود دارد.

در این پایان‌نامه تلاش شده است تا یک روش برای دسته‌بندی داده‌ها بر اساس تئوری

¹Data Classification

²Computer Science

³Artificial Intelligence

⁴Heuristic

اطلاعات (ITDC) ارائه شود و به صورت ریاضی ثابت شود که این روش (ITDC) در صورت برقرار بودن شرایطی مانند Naive Bayes عمل می‌کند و در حقیقت یک گسترش برای Naive Bayes محسوب می‌شود و سپس با دو مثال عملی به بررسی کارایی ITDC پرداخته شده است. شرط‌های مورد نیاز برای عملکرد مشابه Naive Bayes و ITDC بیشتر معطوف به مجموعه آموزش است که شامل به اندازه کافی بزرگ بودن و همچنین گستردگی کافی آن مجموعه است. شرایط مورد نظر در حقیقت شرایطی هستند که در آن‌ها Naive Bayes و سایر الگوریتم‌های مبتنی بر Bayes به درستی (با دقت خوبی) کار می‌کنند.

ساختار پایان‌نامه به این صورت می‌باشد که در فصل دوم به معرفی تئوری اطلاعات و مرور کلی بر مفاهیم اولیه آن پرداخته شده است. فصل سوم، شامل معرفی دقیق مساله دسته‌بندی و همچنین بررسی اجمالی چند روش موجود با تاکید بر Naive Bayes خواهد بود و در ادامه فصل چهارم به بررسی دقیق‌تر Naive Bayes می‌پردازد. در فصل پنجم، به معرفی ITDC پرداخته شده است. فصل‌های ششم و هفتم نیز شامل چند مساله معروف در زمینه دسته‌بندی داده‌ها می‌باشد.

۲ معرفی تئوری اطلاعات

در این بخش به معرفی مفاهیم اساسی تئوری اطلاعات^۵ می پردازیم. تئوری اطلاعات در مقاله‌ای توسط Claude E. Shannon به نام A mathematical theory of computation [۲۶]، که شامل نتایج پایه‌ای در رابطه با منابع بدون حافظه و کانال‌ها و معرفی مدل‌های جامع‌تر ارتباطی، از جمله منابع و کانال‌ها با تعداد حالات متناهی بود وجود آمد [۱۲].

۱-۲ تئوری اطلاعات

Shannon به مطالعه ارتباطات در مجردترین حالت پرداخت. برای این منظور او نیاز داشت تا پیام‌ها (که انتقال آن‌ها بین فرستنده و گیرنده یک ارتباط را تشکیل می‌دهد) را بررسی کند. برای این منظور Shannon به جای در نظر گرفتن پیام‌ها، یک سامانه^۶ را به عنوان فرستنده پیام در نظر گرفت. هر یک از این سامانه‌ها می‌توانند تعداد متناهی و یا نامتناهی پیام را با احتمالات متفاوت ایجاد نمایند.

Entropy

Shannon نیاز به یک معیار برای سنجش مقدار اطلاعاتی که توسط سامانه مبدا هنگامی که یک پیام توسط آن برای ارسال انتخاب می‌شد (به عبارت دیگر مقدار عدم اطمینانی که در پیام خروجی مبدا به وجود می‌آید) داشت. Hartly پیش از Shannon، با مثال‌های زیادی نشان داده بود که تعداد خروجی‌های ممکن یک منبع (پیام‌ها) در مدت زمان T به صورت نمایی با T رشد می‌کند [۱] و بنابراین مشاهده، پیشنهاد ارائه تعریف میزان اطلاعات یک منبع بر اساس لگاریتم این رشد را داده بود [۲]. Shannon این ایده را با وارد کردن خصوصیت‌های آماری منبع تولید کننده پیام گسترش داد.

تعریف پیشنهادی Shannon در ساده‌ترین حالت ممکن که در آن خروجی یک منبع به وسیله یک متغیر تصادفی^۷ X که مقادیر خود را از یک مجموعه مانند $\Sigma = \{1, 2, \dots, N\}$ با احتمال $p_i, i \in \Sigma$ انتخاب می‌کند (با فرض مخالف صفر بودن همه p_i ‌ها) به صورت زیر

⁵Information Theory

⁶System

⁷Random Variable

است:

$$H(X) = \sum_i p_i \log \frac{1}{p_i} \quad (2-1)$$

$H(X)$ را آشفته‌گی منبع^۸ می‌نامند. اگر لگاریتم در عبارت بالا در پایه دو محاسبه گردد، آشفته‌گی منبع بر اساس $\frac{\text{bits}}{\text{symbols}}$ (تعدادبیت‌ها به ازای هر سمبل) خواهد بود. در مبحث فشرده‌سازی^۹، از مقدار Entropy یک مجموعه داده به عنوان بهترین نرخ فشرده‌سازی استفاده می‌شود [۱۹].

Shannon نشان داد که $H(X)$ در خاصیت‌هایی که از یک معیار اندازه‌گیری مقدار اطلاعات به طور بدیهی انتظار می‌رود صدق می‌کند. این خاصیت‌ها به صورت زیر هستند: [۱۹]

۱. اگر منبع تولید کننده پیام دارای احتمالات مساوی باشد (یعنی $\forall i, j \in \Sigma, p_i = p_j$) میزان آشفته‌گی منبع برابر خواهد بود با:

$$H(X) = \log |\Sigma|$$

و با تعریف Hartley [۸] هم خوان است.

۲. اگر منبعی یک پیام را با احتمال ۱ تولید کند و در نتیجه سایر پیام‌ها دارای احتمالی برابر صفر باشند، میزان آشفته‌گی منبع برابر خواهد بود با:

$$H(X) = 0$$

که با شهود ما از یک چنین منبعی مطابقت دارد. در حقیقت در مورد خروجی این منبع هیچ عدم قطعیتی وجود ندارد و به صورت قطعی خروجی آن را شناخته شده است.

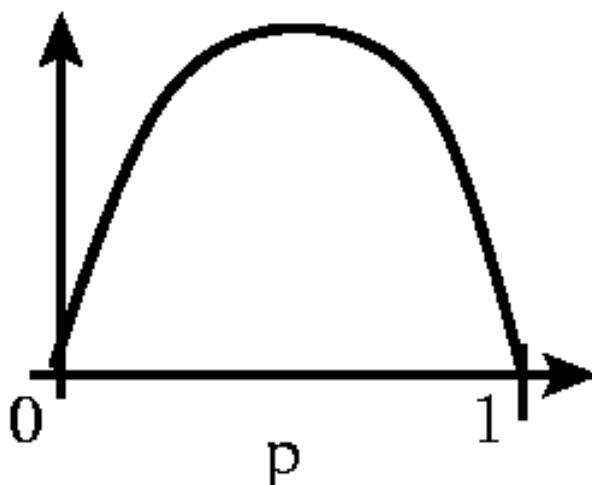
۳. اگر یک مجموعه از پیام‌ها مانند Σ موجود باشد به طوری که $\Sigma' \subset \Sigma$ و $\forall i, j \in \Sigma' : p_i = p_j$ و $\forall i \in \Sigma - \Sigma' : p_i = 0$ آن‌گاه میزان آشفته‌گی منبع تولید کننده Σ برابر خواهد بود با:

$$H(X) = \log |\Sigma'|$$

این موضوع با شهود ما از این‌که مقدار آشفته‌گی در پیام‌های تولید شده توسط یک منبع با خصوصیات توصیف شده فوق معادل است با مقدار آشفته‌گی منبعی که پیام‌های Σ' را تولید می‌کند صدق می‌کند.

⁸Entropy of the Source

⁹Compression



شکل ۱-۲: میزان آشفتگی به ازای مقادیر مختلف احتمالات در حالت دو متغیره

۴. مقدار آشفتگی (عدم قطعیت) هر منبع بزرگتر یا مساوی با صفر است یعنی $H(X) \geq 0$.

$$\forall i : p_i \log \frac{1}{p_i} = -p_i \log p_i$$

$$\forall i : 0 \leq p_i \leq 1$$

اگر $0 < p_i \leq 1$ باشد، داریم $\log p_i \leq 0$ و در نتیجه $-p_i \log p_i \geq 0$ خواهد بود. در حالت خاص $p_i = 0$ از حالت حدی برای تعریف مقدار $p_i \log p_i$ استفاده می‌کنیم:

$$\lim_{t \rightarrow 0} t \log t = 0$$

در نتیجه مجموع چند مولفه نامنفی (مثبت و یا برابر با صفر) همیشه مثبت خواهد بود.

شکل ۱-۲، نشان‌دهنده مقدار آشفتگی منبع تولیدکننده دو پیام با احتمال‌های p و $1-p$ به‌ازای مقادیر مختلف $p \in [0, 1]$ است. تصویر نشان‌گر این موضوع است که بیشینه مقدار آشفتگی (عدم قطعیت) منبع تولیدکننده دو پیام هنگامی است که $p = 0.5$ باشد، یعنی هنگامی که منبع تولیدکننده پیام‌ها یک تولیدکننده تصادفی پیام‌ها با احتمال برابر باشد مقدار آشفتگی آن بیشینه است.

آشفتگی توامان X و Y

آشفتگی بین X و Y ^{۱۰} به صورت زیر تعریف می‌شود:

$$H(X, Y) = \sum_{(x,y) \in X \times Y} -P(x, y) \log P(x, y) \quad (2-2)$$

اگر دو متغیر تصادفی X و Y از هم مستقل باشند یعنی $P(x, y) = P(x)P(y)$ خواهیم داشت:

$$H(X, Y) = H(X) + H(Y)$$

که به آسانی با استفاده از رابطه ۲-۲ با جایگزین کردن $P(x, y)$ با $P(x)P(y)$ بدست می‌آید.

آشفتگی شرطی X با دانستن مقدار Y

آشفتگی شرطی X با دانستن مقدار Y ^{۱۱} برابر است بامیانگین، روی Y ، از مقدار آشفتگی X با دانستن مقدار Y و به صورت زیر تعریف می‌شود:

$$\begin{aligned} H(X|Y) &= - \sum_{y \in Y} P(y) \left(\sum_{x \in X} P(x|y) \log P(x|y) \right) \\ &= - \sum_{(x,y) \in X \times Y} P(x, y) \log(P(x|y)) \end{aligned}$$

این کمیت، میزان تردید (عدم قطعیت) باقی مانده در مورد x هنگامی که مقدار y مشخص شود را بیان می‌کند [۷].

قاعده زنجیر در مورد میزان اطلاعات

قانون Bayes در احتمالات بیان می‌کند:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

¹⁰Joint Entropy

¹¹Conditional Entropy of X given Y

و با توجه به رابطه Bayes داریم:

$$\log(P(x, y)) = \log(P(x)P(y|x)) = \log P(x) + \log P(y|x)$$

و با توجه به تعریف آشفنگی توامان X و Y داریم:

$$\begin{aligned} H(X, Y) &= \sum_{(x, y) \in X \times Y} -P(x, y) \log P(x, y) = - \sum_{(x, y) \in X \times Y} P(x, y) (\log P(x) + \log P(y|x)) \\ &= - \sum_{(x, y) \in X \times Y} P(x, y) \log P(x) - \sum_{(x, y) \in X \times Y} P(x, y) \log P(y|x) \\ &= H(X) + H(Y|X) \end{aligned}$$

رابطه بالا بیان می‌کند که میزان تردید در اطلاعات (توامان) X و Y برابر با جمع تردید در مقدار X و تردید در مقدار Y با دانستن مقدار متغیر X می‌باشد. به راحتی می‌توان دید که $[30] H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$.

اطلاعات مشترک

مقدار اطلاعات مشترک^{۱۲} بین X و Y با رابطه زیر تعریف می‌شود:

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

این کمیت بیان می‌کند که چه میزان از عدم قطعیت درباره X با دانستن مقدار Y کاسته می‌شود و یا به صورت معادل، چه میزان اطلاعات در مورد Y را می‌توان از X بدست آورد [۱۴]. این معیار یک معیار متقارن است و می‌توان این موضوع را با استفاده از قاعده زنجیر اثبات کرد [۷].

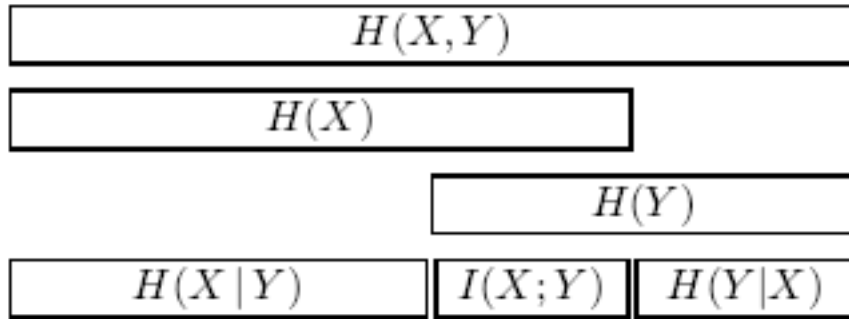
آشفنگی نسبی

مقدار آشفنگی نسبی^{۱۳} یا Kullback-Leibler divergence بین دو توزیع احتمال گسسته $p(x)$ و $q(x)$ به صورت زیر تعریف می‌شود:

$$D_{KL}(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

¹²Mutual Information

¹³Relative Entropy



شکل ۲-۲: رابطه بین آشفته‌گی توامان، آشفته‌گی شرطی و آشفته‌گی هر یک از متغیر

آشفته‌گی نسبی مقداری بزرگ‌تر یا مساوی با صفر دارد و تساوی فقط هنگامی اتفاق می‌افتد که $p = q$ باشد. این نکته که آشفته‌گی نسبی که گاهی از آن با نام "فاصله KL" یاد می‌شود کمیت مقارنی نیست قابل توجه است، یعنی در حالت کلی $D_{KL}(p||q) \neq D_{KL}(q||p)$. آشفته‌گی نسبی در دسته‌بندی داده‌ها و شبکه‌های عصبی نیز علاوه بر تئوری اطلاعات کاربرد دارد [۲۷]

شکل ۲-۲، نشان‌دهنده رابطه بین مفاهیم مطرح شده در بالا است.

۲-۲ فرآیندهای تصادفی

خاصیت مارکف

یک فرآیند تصادفی دارای خاصیت مارکف است اگر احتمال هر یک از حالت‌های بعدی فقط بر اساس حالت کنونی و نه بر اساس حالات گذشته قابل بیان باشد. [۱۱]

اگر بخواهیم تعریف فوق را به صورت ریاضی بیان کنیم خواهیم داشت:

$$\forall h > 0 : p[X(t+h) = y | X(s) = x(s), s \leq t] = p[X(t+h) = y | X(t) = x(t)]$$

یک فرآیند تصادفی را مستقل از زمان می‌نامیم اگر

$$p[X(t+1) = y | X(t) = x(t)] = p[X(1) = y | X(0) = x(0)]$$

در تئوری اطلاعات، به فرایندهای تصادفی^{۱۴} که دارای خاصیت مارکف^{۱۵} باشند فرآیند مارکف گویند^{۱۶} [۳]. زنجیرهای مارکف^{۱۷} نوع خاصی از فرآیندهای مارکف هستند که با زمان گسسته کار می‌کنند.

زنجیرهای مارکف

همان‌طور که اشاره شد زنجیرهای مارکف یک فرآیند تصادفی زمان-گسسته^{۱۸} هستند که دارای خاصیت مارکف می‌باشند.

یک زنجیر مارکف دنباله‌ای از متغیرهای تصادفی X_1, X_2, \dots با خاصیت مارکف است و می‌توان آن را به صورت زیر نمایش داد:

$$p(X_{n+1} = x | X_n = x_n, \dots, X_1 = x_1, X_0 = x_0) = p(X_{n+1} = x | X_n = x_n)$$

مقادیر ممکن برای متغیر X_i باید یک مجموعه شمارا باشد که بدان فضای حالت^{۱۹} گفته می‌شود. برای نمایش یک زنجیر مارکف معمولاً از یک گراف جهت‌دار استفاده می‌شود به طوری که یال‌ها با احتمال تغییر حالت برچسب گذاری شده است.

یکی از مثال‌های زنجیرهای مارکف، ماشین‌های دارای حالات متناهی^{۲۰} است. اگر ماشین M در لحظه n در حالت x باشد آن‌گاه ماشین با احتمالی که فقط وابسته به حالت‌های x و y است و زمان در آن تاثیری ندارد به حالت y منتقل می‌شود. یک زنجیر مارکف ایستا^{۲۱} نامیده می‌شود اگر در رابطه زیر صدق کند:

$$\forall n \in N; p(X_{n+1} = x | X_n = y) = p(X_n = x | X_{n-1} = y)$$

¹⁴Stochastic Process

¹⁵Markov Property

¹⁶Markov Process

¹⁷Markov Chain

¹⁸Discrete-Time

¹⁹State Space

²⁰Finite State Machine

²¹Stationary Markov Chain

خصوصیت‌های زنجیره‌های مارکف

اگر p_{ij}^n را برابر با احتمال انتقال از حالت i به حالت j در زمان n در نظر بگیریم داریم:

$$p_{ij}^n = p(X_n = j | X_{n-1} = i)$$

و اگر زنجیر مارکف مستقل از زمان^{۲۲} باشد می‌توانیم p_{ij} را برابر با احتمال انتقال از حالت i به حالت j تعریف کنیم:

$$p_{ij} = p(X_1 = j | X_0 = i)$$

در ادامه تمرکز ما بیشتر روی زنجیره‌های مارکف مستقل از زمان خواهد بود. مجموع احتمال‌های خروجی از هر حالت برابر با یک است:

$$\forall i : \sum_k p_{ik} = 1$$

حالت j از حالت i قابل دسترسی است اگر حالت i با احتمال بزرگ‌تر از صفر به حالت j می‌رود یعنی:

$$p_{ij} \neq 0$$

گوییم حالات i و j با هم در ارتباط هستند اگر هم i از j قابل دسترسی باشد و هم j از i . یک مجموعه از حالات که اعضای آن دو به دو با هم در ارتباط هستند را یک کلاس ارتباطی^{۲۳} گوییم. کلاس ارتباطی C را بسته نامند اگر و فقط اگر احتمال ترک کردن کلاس (انتقال به حالتی که عضو C نیست) صفر باشد.

یک زنجیر مارکف را ساده‌نشده^{۲۴} می‌نامند اگر کل حالات آن در یک کلاس ارتباطی قرار داشته باشند. این بدان معناست که هر حالتی در یک چنین زنجیری از هر حالت دیگری قابل دستیابی است و در نتیجه در صورت حذف هر یک از حالات آن، مجموعه پیام‌های تولیدی توسط آن زنجیر به کلی تغییر می‌کند.

می‌گوییم حالت i دارای دوره تناوب k است اگر هر بازگشت دوباره به حالت i مستلزم طی مضربی از k گام باشد. برای مثال، هنگامی که برای بازگشت به حالت i مجبور باشیم تا تعداد زوجی گام برداریم، آن‌گاه i متناوب با دوره تناوب ۲ است.

$$k = \gcd\{n : p(X_n = i | X_0 = i) > 0\}$$

²²Time-Stationary Markov Chain

²³Communicating Class

²⁴Irreducible

در رابطه بالا \gcd تابع محاسبه کننده بزرگترین مقسوم علیه مشترک است و k دوره تناوب i خواهد بود. اگر $k = 1$ باشد حالت را غیرنوسانی^{۲۵} می نامند و در غیر این صورت آن را تناوبی با دوره تناوب k می نامند.

بازگشت

به حالت i گذرا گویند اگر با شروع از حالت i با احتمالی مخالف صفر هیچ گاه به i باز نگردیم. اگر T_i اولین باری باشد که بعد از خروج از حالت i به i بازمی گردیم:

$$T_i = \min\{n : X_n = i | X_0 = i\}$$

حالت i یک حالت گذرا است اگر T_i با احتمالی مخالف یک متناهی باشد یعنی:

$$p(T_i < \infty) < 1$$

اگر حالتی گذرا نباشد می گوئیم آن حالت بازگشتی^{۲۶} یا پایدار^{۲۷} است. میانگین زمان بازگشت برای یک حالت به صورت $M_i = E[T_i]$ تعریف می شود. اگر M_i متناهی باشد حالت i را بازگشتی مثبت^{۲۸} و در غیر این صورت آن را ناپایدار^{۲۹} گویند. حالات پایدار حالتی هستند که در تولید رشته های نامتناهی می توان بینهایت بار از آن ها عبور کرد. یعنی می توان نشان داد که یک حالت پایدار است اگر و فقط اگر

$$\sum_{n=0}^{\infty} p_{ii}^n = \infty$$

فرآیندهای Ergodic

حالت i را Ergodic گویند اگر غیرنوسانی باشد و امید ریاضی طول بازگشت آن متناهی باشد [۲]. اگر همه حالات یک زنجیر مارکف Ergodic باشند آن زنجیر را Ergodic گویند.

²⁵Aperiodic

²⁶Recurrence

²⁷Persistent

²⁸Positive Recurrent

²⁹Null Persistent

بررسی حالات دائمی و محدود کردن توزیع‌ها

اگر زنجیر مارکف مستقل از زمان باشد می‌توان آن را با استفاده از یک ماتریس p_{ij} که درایه i و j آن احتمال انتقال از حالت i به حالت j را نشان می‌دهد بیان کرد. بردار Π ، بردار توزیع ایستایی^{۳۰} برای یک زنجیر است اگر مجموع عناصر Π برابر با یک باشد و در رابطه صدق کند:

$$\Pi_j = \sum_i \Pi_i p_{i,j}$$

یک زنجیر مارکف ساده نشدنی دارای بردار توزیع ایستایی است اگر و فقط اگر همه حالات آن بازگشتی مثبت باشند:

$$\Pi = \Pi P$$

در حقیقت، بردار توزیع ایستایی Π برابر با بردار ویژه^{۳۱} چپ ماتریس P با مقدار ویژه ۱ است. اگر ماتریس P دارای یک بردار ویژه Π با مقدار ویژه ۱ باشد می‌توان آن را با رابطه زیر بدست آورد:

$$\lim_{k \rightarrow \infty} P^k = 1\Pi$$

که در آن 1 یک بردار ستونی است که مقدار همه درایه‌های آن برابر با ۱ است.

۳-۲ فرآیندهای مارکف

نوع دیگری از فرآیندهای تصادفی، فرآیندهای مارکف گسسته است که به صورت زیر تعریف می‌شوند [۲۶]: مجموعه متناهی از حالات مانند S_1, S_2, \dots, S_n در کنار مجموعه‌ای از احتمالات جابه‌جایی مانند آن‌چه در زنجیرهای مارکف وجود دارد تشکیل می‌شود. برای تبدیل این فرآیند مارکف به یک منبع اطلاعات، فرض می‌شود که هر جابه‌جایی از یک حالت به حالت دیگر (یال در گراف متناظر) یک سمبل تولید کند. در بخش ۵ به صورت دقیق‌تر به بررسی این گونه فرآیندهای مارکف می‌پردازیم.

³⁰Stationary Distribution

³¹Eigenvector

۴-۲ رابطه بین دسته‌بندی و Entropy

هدف نهایی در دسته‌بندی داده‌ها این است که روشی برای تشخیص تعلق و یا عدم تعلق یک داده به هر یک از دسته‌ها ارائه شود. راه‌های مختلفی برای این منظور وجود دارد که به بخشی از آن‌ها در فصل دوم اشاره شد. در دسته‌بندی معمولاً تعدادی از اعضای یک دسته در اختیار است که در این صورت یکی از روش‌های تعیین عضویت می‌تواند استفاده از میزان شباهت بین داده با اعضای هر دسته باشد. به این ترتیب که با استفاده از یک معیار به سنجش میزان شباهت داده با اعضای هر یک از دسته‌ها پردازیم و مشابه‌ترین دسته را به عنوان دسته برنده اعلام کنیم.

همان‌طور که دیدیم، یکی از راه‌های محاسبه میزان شباهت استفاده از مفهوم مشابه میزان تردید (عدم قطعیت) است. در بخش بعدی به توضیح چگونگی استفاده از این مفهوم برای دسته‌بندی خواهیم پرداخت.

۳ دسته‌بندی داده‌ها

مساله دسته‌بندی داده^{۳۲} یکی از مساله‌های اساسی در هوش مصنوعی^{۳۳} است که سال‌ها مورد مطالعه قرار داشته است و راه حل‌های متعددی برای حل این مساله ارائه شده است. در این بخش، ابتدا تلاش خواهیم کرد تا مساله را به صورت ریاضی بیان کنیم و سپس به توضیح چند روش معمول در این زمینه بپردازیم.

۱-۳ تعریف ریاضی مساله دسته‌بندی داده‌ها

مساله دسته‌بندی داده‌ها درباره حدس زدن و یا پیش‌بینی یک دسته برای یک داده ورودی است. داده ورودی معمولاً یک مجموعه از اعداد و یا داده‌های نمادین^{۳۴} است که به صورت یک بردار d -تایی مانند x نمایش داده می‌شود ولی ممکن است در برخی کاربردها داده ورودی به صورت یک تصویر و یا یک نمودار نیز باشد. دسته‌ای که x بدان تعلق دارد را با y نشان می‌دهیم. در ساده‌ترین حالت، y یکی از دو مقدار $\{1, -1\}$ را اختیار می‌کند [۴]. در این پایان‌نامه توجه اصلی به دسته‌بندی داده‌ها بین دو کلاس (دسته‌بندی دوگانه^{۳۵}) خواهد بود و از کلمه دسته‌بندی به معنای دسته‌بندی دوگانه استفاده شده است. به راحتی می‌توان با استفاده از یک ساختار سلسله‌مراتبی، از دسته‌بندی دوگانه برای حل هر مساله دسته‌بندی دیگری استفاده کرد. دلیل این انتخاب ساده‌تر بودن دسته‌بندی دوگانه ولی در عین حال حفظ قابلیت کلی مساله‌های دسته‌بندی در حالت عمومی است. در دسته‌بندی به دنبال ایجاد تابعی به صورت $g: X \rightarrow \{-1, 1\}$ هستیم که دسته پیشنهادی برای داده (بردار) ورودی x ، به صورت $y = g(x)$ محاسبه شود. به چنین تابعی، تابع دسته‌بندی‌کننده^{۳۶} می‌گوییم. در برخی از تابع‌های دسته‌بندی کننده برد تابع را به $[-1, 1]$ (به جای $\{-1, 1\}$) تغییر می‌دهند. خطای یک تابع دسته‌بندی کننده روی یک مجموعه از داده‌ها $\{(x_i, y_i)\}$ به صورت زیر تعریف می‌شود:

$$Err(g, S) = \sum_{(x_i, y_i) \in S} \left| \frac{g(x_i) - y_i}{2} \right|$$

³²Data Classification

³³Artificial Intelligence

³⁴Symbolic Data

³⁵Binary Classification

³⁶Classifier function

برای بیان مساله یادگیری، از یک مدل احتمالاتی استفاده می‌کنیم که در آن (X, Y) یک زوج تصادفی از $X \times \{-1, 1\}$ شامل داده و دسته مربوط به آن است. علت استفاده از مدل تصادفی این است که ممکن است یک $X = x$ به هر دو دسته $Y = \{-1, +1\}$ متناظر شود. یکی از دلایل‌های این حالت می‌تواند کافی نبودن اطلاعات ویژگی‌های انتخابی باشد. به عنوان مثال، حالتی را در نظر بگیرید که دسته یک مجموعه از داده‌ها توسط سه ویژگی به طور کامل مشخص شود ولی ما فقط از دو تا از این ویژگی‌ها در X استفاده کرده‌ایم. دلیل دیگر برای این حالت، می‌تواند طبیعت احتمالاتی مساله باشد. توزیع احتمال زوج (X, Y) را می‌توان براساس توزیع احتمال X (توسط $P(X)$) و $\eta(x) = P\{Y = 1 | X = x\}$ می‌توان خطای کارآیی تابع g را بر اساس تابع توزیع احتمال خطا اندازه‌گیری کرد:

$$L(g) = P\{g(x) \neq Y\}$$

در صورتی که تابع η یا تخمین مناسبی از آن در دسترس باشد می‌توان با استفاده از تصمیم‌گیری حریصانه^{۳۷} یک تابع دسته‌بندی کننده با کمترین احتمال خطا ارائه داد:

$$g^*(x) = \begin{cases} 1 & \text{if } \eta(x) > \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

به راحتی می‌توان اثبات کرد:

$$\forall g : X \mapsto \{-1, 1\} : L(g^*) \leq L(g)$$

کمترین مقدار خطا $L^* = L(g^*)$ را خطای Bayes^{۳۸} می‌نامیم [۹]. به صورت دقیق‌تر، می‌توان ثابت کرد [۶] که

$$L(g) - L^* = E[\mathbb{1}_{g(x) \neq g^*(x)} | 2\eta(X) - 1 | \geq 0]$$

معمولا دسته‌بندی بر اساس تابع g^* را، دسته‌بندی بر اساس روش Bayes^{۳۹} می‌نامند. در روش‌های آماری^{۴۰} عموماً یک مجموعه از داده‌ها $D_n = \{(X_i, Y_i) : 1 \leq i \leq n\}$ را در اختیار داریم. فرض می‌کنیم که داده‌های مجموعه D_n توسط یک توزیع تصادفی مستقل^{۴۱} از بین زوج‌های ممکن از $X \times \{-1, 1\}$ استخراج می‌شود و در نتیجه $(X_1, Y_1), \dots, (X_n, Y_n)$ با دقت خوبی دارای توزیعی مانند (X, Y) می‌باشد.

³⁷Greedy

³⁸Bayes Error

³⁹Bayes Classifier

⁴⁰Statistical Model

⁴¹Independent Identically Distribution

تابع دسته‌بندی کننده‌ای که بر پایه $D_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ بدست می‌آید را g_n می‌نامیم. کارایی تابع g_n را می‌توان بر پایه احتمال خطای آن محاسبه کرد:

$$L(g_n) = P\{g_n(X) \neq Y | D_n\}$$

توجه اصلی در تئوری و کاربردهای عملی دسته‌بندی، ایجاد تابع g_n به طوری است که احتمال خطای آن نزدیک‌ترین مقدار ممکن L^* است.

یک راه ساده برای مساله دسته‌بندی داده‌ها این است که یک کلاس مانند C از تابع‌های دسته‌بندی کننده مانند $\{-1, 1\}$ $g: X \mapsto \{-1, 1\}$ تشکیل دهیم و سپس با استفاده از D_n یک تخمین برای احتمال خطای هر $g \in C$ یعنی $L(g)$ حساب کرده و بر اساس آن یک تابع را انتخاب کنیم. شاید بدیهی‌ترین راه برای تخمین احتمال خطای $L(g) = P\{g(X) \neq Y\}$ استفاده از

$$L_n(g) = \frac{1}{n} \sum \{1_{g(X_i) \neq Y_i}\}$$

$L_n(g)$ را خطای تجربی^{۴۲} می‌نامیم.

هدف در الگوریتم‌های دسته‌بندی مبتنی بر یادگیری ماشین^{۴۳} تولید توابعی است که می‌توانند هر نمونه از داده‌های ورودی را به یک دسته مناسب نظیر کنند. روش‌های گوناگون یادگیری ماشین برای کاربردهای گوناگون بر اساس مدل‌های احتمالاتی و یا آماری و الگوریتم‌های گوناگون ارائه شده است. در بین این روش‌ها، شبکه‌های عصبی^{۴۴} در کاربردهایی که در آن عناصر بردارهای اعداد حقیقی هستند، و درخت‌های تصمیم^{۴۵}، Naive Bayes در کاربردهایی که عناصر بردار ورودی آن نمادین هستند از معروف‌ترین روش‌ها هستند [۱۰]. در ادامه به توضیح مختصر دو روش اول می‌پردازیم و روش آخر (Naive Bayes) را به علت ساختار ریاضی آن و رابطه نزدیکی که با روش مورد بحث در این پایان‌نامه دارد را به طور دقیق‌تر بررسی خواهیم کرد.

تقریباً همه الگوریتم‌های دسته‌بندی دارای دو مرحله هستند:

۱. مرحله آموزش: در این مرحله یک مدل (بسته به الگوریتم مورد استفاده) انتخاب می‌شود و بر اساس مجموعه‌های آموزش پارامترهای مدل مربوطه تعیین می‌شود. این مرحله، معمولاً فقط یک‌بار اجرا می‌شود و بارها از نتیجه آن استفاده می‌شود بنابراین پیچیدگی محاسبه در این مرحله، زیاد مورد توجه قرار نمی‌گیرد.

⁴²Empirical Error

⁴³Machine Learning

⁴⁴Neural Network

⁴⁵Decision Tree

۲. مرحله استفاده: در این مرحله، از مدل و پارامترهای بدست آمده در مرحله آموزش استفاده شده و با به کارگیری آن‌ها، در مورد بردار ورودی اظهار نظر می‌شود. برخلاف مرحله قبل، این مرحله بارها مورد استفاده قرار می‌گیرد و بنابراین پیچیدگی محاسبه آن بسیار مهم است.

۲-۳ درخت‌های تصمیم

درخت‌های تصمیم دارای یک ساختار درختی هستند که راس‌های هر سطح آن متناظر با یک ویژگی بوده و دارای دو و یا بیشتر فرزند است که بسته به مقدار آن ویژگی در بردار ورودی یک و یا چند تا از فرزندان انتخاب می‌شوند و در نهایت برگ‌ها به یکی از دسته‌ها مربوط می‌شوند.

در مرحله آموزش این روش، ابتدا ساختار درخت (ترتیب قرارگرفتن ویژگی‌ها در درخت) تعیین می‌شود و سپس با استفاده از مجموعه آموزشی دسته‌هایی که برگ‌ها بدان‌ها باید اشاره کنند تعیین می‌شود. پیچیدگی درخت تصمیم تولید شده (تعداد راس‌های درخت) بسیار وابسته به ساختار درخت است. روش‌های ابتدایی^{۴۶} برای تعیین ترتیب قرارگرفتن ویژگی‌ها در درخت تصمیم وجود دارد ولی هیچ روش اثبات شده بهینه‌ای در این رابطه وجود ندارد.

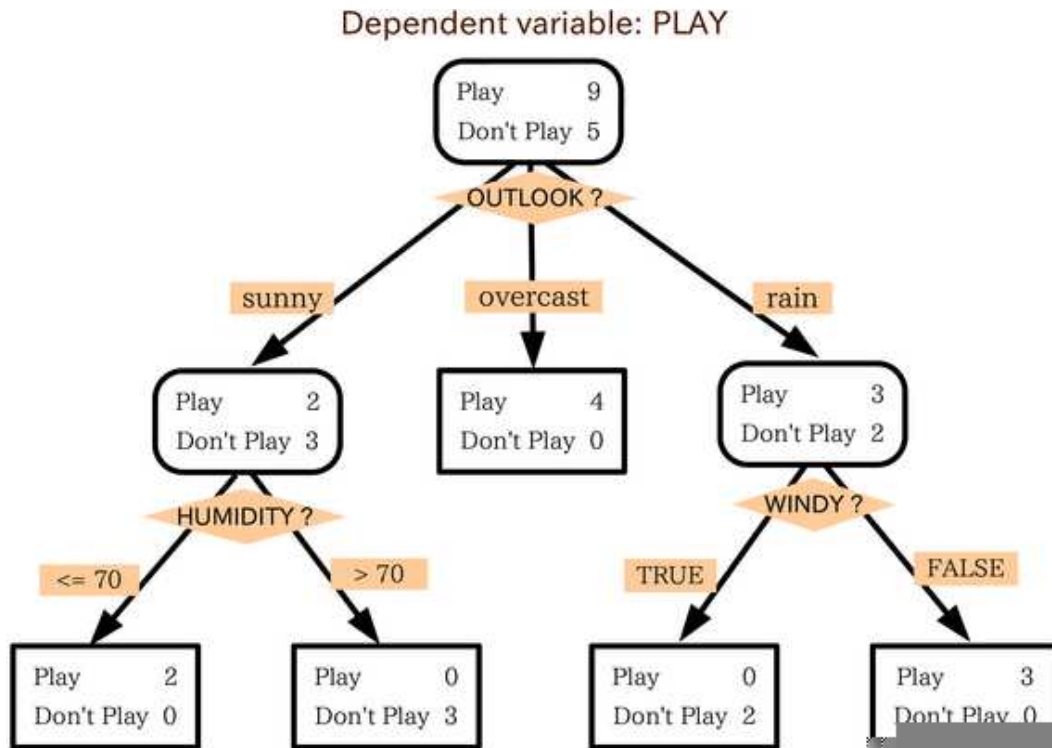
در هنگام استفاده از یک درخت تصمیم، با شروع از ریشه در هر راس، براساس مقدار ویژگی اختصاص داده شده به راس، به یکی از فرزندان آن راس می‌رویم. و این روند را تا رسیدن به یک برگ ادامه می‌دهیم و دسته متناظر با برگ و یا برگ‌های متناظر به عنوان دسته (دسته‌های) مربوطه گزارش می‌شود [۲۴].

مزیت اساسی درخت‌های تصمیم بر سایر الگوریتم‌های دسته‌بندی، در کم بودن هزینه محاسبه برای انتخاب دسته پیشنهادی و قابلیت کار با داده‌های عددی و نمادین است و همچنین در حالتی که داده دارای ویژگی‌های گم‌شده^{۴۷} باشد درخت‌های تصمیم کارآیی خوبی دارند. یکی دیگر از امتیازات درخت‌های تصمیم قابلیت نمایش ساختارهای تصمیم پیچیده در آن‌ها است. شکل ۳-۱، یک نمونه از درخت‌های تصمیم را نشان می‌دهد.

یکی از بهترین الگوریتم‌هایی که برای تولید درخت‌های تصمیم به کار گرفته می‌شود الگوریتم C4.5 [۲۰] است.

⁴⁶Hueristic

⁴⁷Missing Variables



شکل ۳-۱: یک نمونه از درخت‌های تصمیم برای بازی گلف

۳-۳ شبکه‌های عصبی

شبکه‌های عصبی مصنوعی^{۴۸} که معمولاً با نام شبکه عصبی^{۴۹} شناخته می‌شوند از یک گروه به هم متصل از سلول‌های عصبی مصنوعی^{۵۰} (نرون) تشکیل شده و از یک مدل ریاضی یا محاسباتی برای پردازش اطلاعات بر حسب اتصالات بین نرون‌ها استفاده می‌کند. در بیشتر مواقع شبکه عصبی مصنوعی یک سامانه سازگار شونده^{۵۱} است که ساختار آن در طول فرآیند بر اساس اطلاعات داخلی و یا خارجی تغییر می‌کند.

در کاربرد، شبکه عصبی یک ابزار برای مدل‌سازی داده‌های آماری غیر خطی است. از آن می‌توان برای مدل کردن رابطه‌های پیچیده بین ورودی‌ها و خروجی‌ها و یا پیدا کردن الگو استفاده کرد.

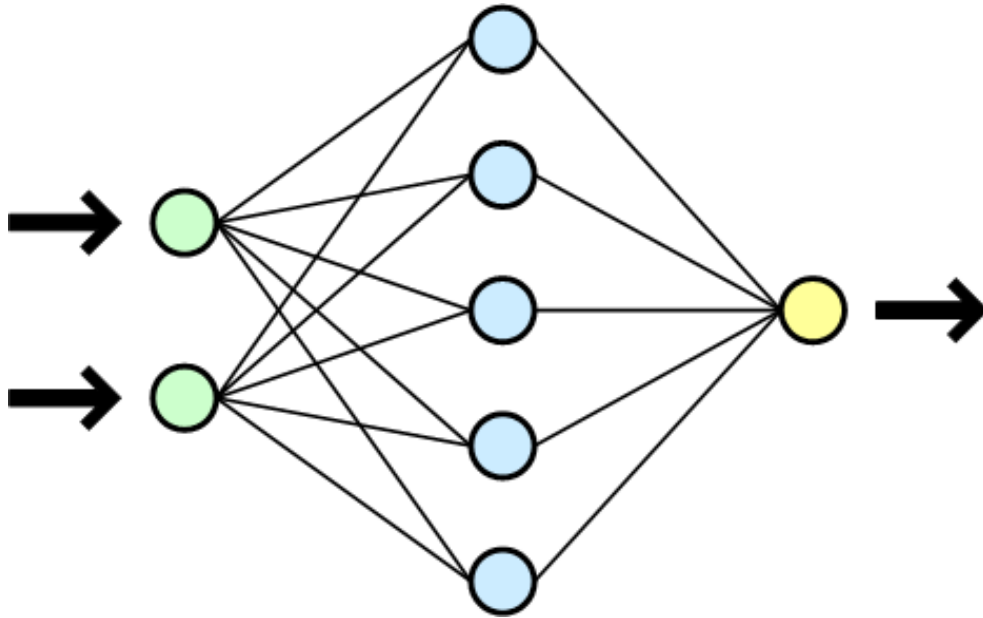
تعریف دقیق و مورد قبول عمومی برای شبکه عصبی وجود ندارد. یکی از تعاریف به صورت زیر است: شبکه عصبی از یک شبکه از عناصر پردازش‌گر ساده (Neurons) تشکیل شده است

⁴⁸Artificial Neural Network

⁴⁹Neural Network

⁵⁰Artificial Neurons

⁵¹Adaptive System



شکل ۳-۲: یک نمونه شبکه عصبی با سه لایه

و می‌تواند ساختارها و رفتارهای پیچیده را با چگونگی ارتباط بین عناصر پردازش‌گر شبکه و داده‌های ورودی و خروجی بیان کرد. ایده اصلی در شبکه‌های عصبی مصنوعی از سیستم عصبی انسان گرفته شده است که در آن فعالیت‌ها به صورت جمعی و به صورت موازی بوسیله هر واحد ساده انجام می‌شود [۱۳].

شکل ۳-۲ یک سامانه شبکه عصبی را نشان می‌دهد. هر نرون دارای یک و یا چند یال ورودی و دقیقاً یک یال خروجی است. مقدار خروجی هر نرون به صورت زیر محاسبه می‌شود:

$$\sum_{i=1}^n g(w_i I_i)$$

که در آن، هر w_i یک عدد حقیقی است که به هر یال ورودی نسبت داده می‌شود و I_i مقدار ورودی i -ام است. در حقیقت w_i ها مقدار اثر هر ورودی را تعیین می‌کنند و از این رو معمولاً w_i را «وزن» نیز می‌نامند. و تابع فعالیت^{۵۲} نرون است. تابع فعالیت نرون‌ها می‌تواند هر تابع تام از اعداد حقیقی به بازه $[-1, +1]$ باشد ولی معمولاً از تابع Sigmoid به شکل $\frac{1}{1+e^{-at}}$ استفاده می‌شود [۱۳].

معمولاً در شبکه‌های عصبی، هر نرون در لایه خروجی نشان‌گر یک دسته است و دسته پیشنهادی توسط یک شبکه عصبی دسته‌ای است که خروجی نرون مربوطه (در حالتی که

⁵²Activation Function

چند نرون به یک کلاس مربوط شده اند بیشینه مقدار خروجی نرون‌ها) بیشترین باشد. شبکه‌های عصبی بر حسب نوع ارتباطاتی که نرون‌ها می‌توانند با هم داشته باشند و همچنین بر حسب تابع فعالیت نرون‌ها به دسته‌های مختلفی تبدیل می‌شوند. شبکه‌های عصبی بازخورد خطا^{۵۳}[۳۴] از معروف‌ترین شبکه‌های عصبی می‌باشند.

۴-۳ Naive Bayes

روش استدلال Bayes^{۵۴} از روش‌های آماری برای استنتاج استفاده می‌کند. در این روش فرض این است که ویژگی‌های^{۵۵} مورد علاقه بر اساس توزیع احتمالی که وجود دارد ولی احتمالاً ما از آن بی‌اطلاع هستیم قابل توصیف هستند و می‌توان بر اساس یک تخمین از توزیع احتمال (که بر اساس مجموعه آموزشی به دست می‌آید) تصمیم بهینه گرفت. روش‌هایی که بر این اساس کار می‌کنند را الگوریتم‌های یادگیری Bayes^{۵۶} می‌نامند. در بین الگوریتم‌های یادگیری Bayes، می‌توان به روش Naive Bayes اشاره کرد که به علت سادگی پیاده‌سازی بسیار مورد توجه است. این الگوریتم در کنار سادگی، از کارایی بسیار خوبی مخصوصاً در زمینه دسته‌بندی متون^{۵۷} برخوردار است [۲۱]. در [۱۷]، به مقایسه کارایی دسته‌بندی با استفاده از Naive Bayes با سایر روش‌ها از جمله درخت‌های تصمیم و شبکه‌های عصبی پرداخته شده است. آمار ارائه شده در این مقاله نشان می‌دهد که Naive Bayes یک رقیب خوب برای سایر الگوریتم‌های دسته‌بندی داده‌ها است و حتی در بعضی از مساله‌ها، بهتر از سایر روش‌ها عمل می‌نماید.

تئوری Bayes در نظریه احتمال یک روش برای محاسبه کردن احتمال رخ دادن h به شرط اتفاق افتادن D یعنی $P(h|D)$ با استفاده از احتمال رخ دادن D به شرط اتفاق افتادن h یعنی $P(D|h)$ و احتمال رخ دادن D و h به صورت زیر ارائه می‌دهد [۳۱]:

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

برای بیان الگوریتم Naive Bayes و سایر الگوریتم‌های وابسته به روش Bayes، در جزئیات مساله دسته‌بندی باید تغییراتی ایجاد کنیم. فرض کنید به ازای هر دسته مانند

⁵³Back Propagation

⁵⁴Bayesian Resoning

⁵⁵Features

⁵⁶Bayesian Learning Methods

⁵⁷Document Classification

یک بردار (نامعلوم) θ^c وجود دارد که داده‌های متعلق به آن دسته را به صورت مستقل تولید می‌کند. داده‌هایی که دسته مربوط به آن‌ها شناخته شده است مجموعه آموزشی (به ازای هر دسته مانند c مجموعه مانند D^c) نامیده می‌شوند و بقیه داده‌های ممکن (داده‌های مجاز) مجموعه آزمون را تشکیل می‌دهند [۲۱].

پیدا کردن (تخمین زدن) مقادیر پارامتر مجهول یکی از مهم‌ترین بخش‌های روش Naive Bayes است. روش‌های متفاوتی برای برآورد کردن این بردار (θ^c) وجود دارد که در ادامه به بیان جزئیات دو روش پرکاربردتر خواهیم پرداخت:

ML Naive Bayes

منطقی‌ترین راه برای تخمین زدن بردار θ^c استفاده از تنها داده شناخته شده یعنی D^c است. در روش بیشترین احتمال^{۵۸}، فرض بر این است که هیچ دانشی در مورد چگونگی توزیع احتمال دسته‌ها در دست نیست و بنابراین توزیع یکنواخت برای آن در نظر گرفته می‌شود. بردار θ^c توسط بردار $\hat{\theta}^c$ چنان تخمین زده می‌شود که با بیشترین احتمال ممکن مجموعه آموزشی را تولید کند یعنی:

$$\hat{\theta}_c = \arg \max_{\theta} p(D^c | \theta)$$

و سپس از بردار $\hat{\theta}^c$ به جای θ^c استفاده می‌شود. این روش دارای نقاط ضعفی است که شاید از مهمترین آن‌ها، وابستگی بسیار زیاد به چگونگی انتخاب مجموعه آموزش است. برای مثال، اگر یک ویژگی در مجموعه آموزش پدیدار نشده باشد در این روش فرض می‌شود که ویژگی مذکور در هیچ یک از داده‌های عضو آن دسته، ظاهر نخواهد شد. برای بدست آوردن احتمال تولید مجموعه آموزشی با فرض در اختیار بودن بردار θ می‌توان از روابط زیر استفاده کرد:

$$P(D^c | \theta) = \frac{N^c!}{\prod_k N_k^c!} \prod_k \theta_k^{N_k^c}$$

که در آن N_k^c نشان‌دهنده تعداد دفعاتی است که ویژگی k (f_k) در مجموعه آموزشی مربوط به دسته c ظاهر شده و $N^c = \sum_k N_k^c$ برابر با کل ویژگی‌هایی است که (با در نظر گرفتن تکرار) در D^c وجود دارد. فرض کنید θ_k ، k -امین مولفه بردار θ باشد و نشان‌دهنده احتمال پدیدار شدن ویژگی k -ام در هر مولفه از داده‌های عضو D^c باشد. در تخمین ML بر اساس D^c (مقداری برای $\hat{\theta}^c$ که احتمال $p(D^c | \hat{\theta}^c)$ را بیشینه می‌کند) به ازای تمام k ‌ها، $\hat{\theta}_k^c = \frac{N_k^c}{N^c}$ خواهد بود.

⁵⁸Maximum Likelihood

همان‌طور که اشاره شد Naive Bayes، از $\hat{\theta}^c$ برای ارزیابی این که آیا داده آزمون، d ، به دسته c تعلق دارد (توسط تولید کننده دسته c تولید شده است) استفاده می‌کند. از آنجایی که تخمین بردار پارامتر و در نتیجه قضاوت در مورد داده‌های آزمون بسیار وابسته به چگونگی انتخاب داده‌های مجموعه آموزشی، D^c ، است این مجموعه (D^c) می‌تواند تاثیر زیادی روی کارایی داشته باشد.

انتخاب جواب بهینه طبق رابطه Bayes چنین خواهد بود:

$$H_{ML}(d) = \arg \max_c p(c|D, d) = \arg \max_c p(d|\hat{\theta}^c)p(c)$$

اگر همه دسته‌ها، دارای شانس تقریباً برابر باشند (یعنی $p(c)$ دارای توزیع یکسان^{۵۹} باشد)، فرآیند انتخاب دسته جواب را می‌توان به صورت زیر ساده کرد:

$$H_{ML}(d) = \arg \max_c p(d|\hat{\theta}^c) = \operatorname{argmax}_c \prod_k \left(\frac{N_k^c}{N^c}\right)^{N_k^d}$$

N_k^d برابر است با تعداد دفعاتی که ویژگی k -ام در داده آزمون (d) ظاهر شده است. همان‌طور که اشاره شد در بسیاری از کاربردها، ممکن است N_k^c برای یک یا چند ویژگی صفر باشد و در این صورت مقدار $P(d|c)$ برابر با صفر خواهد بود و ممکن است به گمراهی H_{ML} منتهی شود. برای جلوگیری از این حالت، به یک سامانه Naive Bayes علاوه بر داده‌های آموزشی، یک دنباله از اعداد a_k به عنوان ورودی داده می‌شود و از این دنباله که a_k با ویژگی k -ام در ارتباط است و ($a = \sum a_k$) به صورت زیر استفاده می‌شود:

$$H_{ML}(d) = \arg \max_c \prod_k \left(\frac{N_k^c + a_k}{N^c + a}\right)$$

به این روش (اضافه کردن یک دنباله از اعداد به Naive Bayes) M-Estimation گفته می‌شود [۲۹]. در کاربرد، تعیین این دنباله به صورت بهینه کار غیرممکنی است و بنابراین تنها از این خاصیت رابطه جدید، که در صورت برابر با صفر بودن N_k^c ، $P(d|c)$ صفر نخواهد بود استفاده می‌شود (با انتخاب a_k ها به صورت پیش فرض). در بیشتر به کارگیری‌های M-Estimation، از دنباله یکنواخت ($\forall i, j : a_i = a_j$) استفاده می‌شود که در آن $a_i = 1$ است [۲۱].

MAP Naive Bayes

MAP Naive Bayes با فرض این که هیچ دانشی از چگونگی توزیع احتمال دسته‌ها در دسترس نیست کار می‌کند ولی اگر چنین دانشی موجود باشد آن‌گاه می‌توان از روش MAP Naive

⁵⁹Uniform Distribution

Bayes⁶⁰ استفاده کرد. MAP estimation یک گسترش از ML است و مانند ML Bayes کار می‌کند یعنی بازهم بهترین $\hat{\theta}^c$ بر حسب داده‌های مجموعه آموزشی انتخاب کرده و سپس از آن و توزیع احتمال دسته‌ها برای دسته‌بندی استفاده می‌کنیم. در MAP estimation، θ^c توسط $\hat{\theta}^c$ (مشابه ML) تخمین زده می‌شود:

$$\hat{\theta}^c = \operatorname{argmax}_{\theta} p(D^c|\theta)$$

فرآیند تصمیم در MAP به صورت زیر است:

$$H_{MAP}(d) = \operatorname{argmax}_c p(c|d) = \operatorname{argmax}_c P(d|\theta^c)p(c)$$

که در آن $p(c)$ احتمال دسته c است.

برای روشن‌تر شدن چگونگی کاربرد تئوری Bayes در دسته‌بندی و همچنین روش‌های ML و MAP در تخمین پارامترها و استدلال، می‌توان به مساله زیر اشاره کرد [۱۸] فرض کنید در مساله تشخیص یک نوع سرطان یکی از دو فرض ”(۱) مراجعه‌کننده به نوع خاصی از سرطان دچار است” و یا ”(۲) مراجعه‌کننده دچار بیماری سرطان نیست” برقرار باشد. داده‌های موجود در این مساله از یک آزمایشگاه به صورت نتیجه یک آزمون دارای دو حالت مثبت (\oplus) و یا منفی (\ominus) است، و همچنین آزمون انجام شده در آزمایشگاه، یک نشان‌گر غیردقیق از وجود بیماری است. از میان ۱۰۰۰ آزمایش‌شونده در حالتی که نتیجه آزمون مثبت است ۹۸۰ نفر واقعا بیمار بوده‌اند و هنگامی که آزمون جواب منفی می‌دهد ۹۷۰ نفر سالم بوده‌اند. و در سایر موارد پاسخ آزمون نادرست است.

$$P(\oplus | \text{cancer}) = 0.98 \quad P(\ominus | \text{cancer}) = 0.02$$

$$P(\oplus | \neg \text{cancer}) = 0.03 \quad P(\ominus | \neg \text{cancer}) = 0.97$$

- ۱- فرض کنید بیماری وجود دارد که نتیجه آزمون او مثبت بوده است. آیا این بیمار سرطان دارد؟
 - ۲- اگر بدانیم تنها ۸٪/۰ از افراد جامعه ممکن است مبتلا به این‌گونه از سرطان باشند، آیا بیماری با آزمون مثبت، سرطان دارد؟
- در این مساله فقط یک ویژگی (نتیجه آزمون که می‌تواند \oplus و یا \ominus باشد) و دو دسته (cancer و $\neg \text{cancer}$) وجود دارد و داده آزمون به صورت $d = \langle \oplus \rangle$ است.
- ۱- در حالت اول، دانشی در مورد توزیع احتمال دسته‌ها وجود ندارد و برای حل مساله از ML استفاده کنیم:

⁶⁰Maximum A Posteriori

$$p(d|cancer) = p(\oplus|cancer) = 0.98, \quad p(d|\neg cancer) = p(\oplus|\neg cancer) = 0.3$$

$$\Rightarrow H_{ML}(d) = cancer$$

۲- در این حالت، دانشی از چگونگی توزیع احتمال دسته‌ها وجود دارد و بنابراین باید (بهتر است) از ML استفاده شود.

$$P(cancer) = 0.008 \quad P(\neg cancer) = 0.992$$

$$P(\oplus|cancer) = 0.98 \quad P(\ominus|cancer) = 0.02$$

$$P(\oplus|\neg cancer) = 0.03 \quad P(\ominus|\neg cancer) = 0.97$$

احتمال مبتلا بودن مراجعه کننده با نتیجه آزمون مثبت ($d = \oplus$) به سرطان طبق MAP از رابطه زیر بدست می آید:

$$H_{cancer}(d) = p(d|\theta^{cancer})p(cancer) = 0.98 * 0.008 = 0.00784$$

و احتمال مبتلا نبودن وی به سرطان نیز از رابطه مشابه زیر بدست می آید:

$$H_{\neg cancer}(d) = p(d|\theta^{\neg cancer})p(\neg cancer) = 0.03 * 0.992 = 0.02976$$

حال اگر بخواهیم با استفاده از Naive Bayes، داده d را به یکی از دو دسته اختصاص دهیم باید اعلام کنیم که مراجعه کننده سالم است.

۴ بررسی تخمین‌های پارامترهای Naive Bayes

Naive Bayes یک الگوریتم حریصانه است که براساس احتمال عضویت داده در کلاس تصمیم می‌گیرد. برای انجام این کار، Naive Bayes نیاز به تخمین زدن پارامترهای هر کلاس (که در بخش قبل بدان‌ها اشاره شد) دارد و کیفیت تخمین پارامترها به طور مستقیم بر کارایی الگوریتم تاثیر دارد. در این بخش نشان خواهیم داد که تخمین‌های به کار رفته در Naive Bayes پایدار^{۶۱} هستند و سپس کارکرد روش‌های تخمین را روی داده‌های آموزشی متنهایی با تحلیل جابه‌جایی^{۶۲} و انحراف از معیار^{۶۳} بررسی خواهیم کرد. نشان خواهیم داد که جابه‌جایی در تخمین با افزایش داده‌های آموزشی، به سمت صفر میل خواهد کرد. تحلیل نشان خواهد داد که در صورت عدم وجود مجموعه آموزشی به اندازه کافی گسترده در یک دسته، می‌تواند روی کارایی کل سامانه تاثیر داشته باشد. انحراف از معیار سامانه برابر با حاصل جمع انحراف هر یک از مولفه‌های آن سامانه است. در نتیجه اگر یک دسته دارای انحراف بالایی باشد، مقدار انحراف کل سامانه نیز بزرگ خواهد بود. [۲۱].

۱-۴ پایداری

پایداری بدین معناست که با افزایش اندازه مجموعه آموزش، مقدار تخمین زده شده به مقدار واقعی آن میل می‌کند. ML، m-estimation، و MAP سعی می‌کنند که با استفاده از $\{a_k\}$ (در ML همه آن‌ها برابر با صفر هستند) تخمین $\hat{\theta}$ از θ را به صورت زیر تولید کنند:

$$\hat{\theta}_k^c = \frac{a_k + N_k^c}{a + N^c}$$

Cover و Thomas در کتاب "Elements of Information Theory" [۵] روشی را برای توصیف خصوصیات توزیع‌های تجربی شرح می‌دهند. فرض کنید X یک متغیر تصادفی با پارامترهای θ_k^c باشد و p_x نشان‌گر توزیع پارامترهای X و p_y نشان‌گر توزیع احتمال هنگامی که N^c نمونه برداشته شده از X دارای تعداد تکرار $\{N_k^c\}$ باشند. در این حالت $\hat{\theta}_k^c = \frac{a_k + N_k^c}{a + N^c}$ خواهد بود. احتمال مشاهده یک چنین $\{N_k^c\}$ و در نتیجه ایجاد تخمین $\hat{\theta}$ برابر است با:

$$p(\hat{\theta}^c | \theta^c) = \frac{N^c!}{\prod_k N_k^c!} \prod_k (\theta_k^c)^{N_k} = \frac{N^c!}{\prod_k N_k^c!} 2^{-N(H(p_Y) + D(p_Y || p_X))}$$

⁶¹Consistent

⁶²Bias

⁶³Variance

میانگین تخمین به روش بالا برابر با $\frac{a_k + N^c \theta_k^c}{a + N^c}$ خواهد بود و با بزرگ شدن $N^c \rightarrow \infty$ به θ_k^c میل می‌کند. همچنین انحراف از معیار تخمین نیز از رابطه زیر بدست می‌آید:

$$\delta_{c,k}^2 = \frac{N^c \theta_k^c (1 - \theta_k^c)}{(a + N^c)^2}$$

با بزرگ شدن N^c به سمت صفر میل می‌کند. پس تخمین حاصل از روش‌های فوق پایدار هستند زیرا در حالت حدی، به سمت مقدار دقیق میل می‌کند و دارای انحراف از معیار صفر است. در نتیجه با افزایش تعداد داده‌های عضو مجموعه آموزشی، تخمین ما به مقدار واقعی نزدیک‌تر می‌شود.

۲-۴ جابه‌جایی

از آن جایی که در کاربردهای عادی، همیشه اندازه مجموعه آموزشی متناهی است، درک چگونگی رفتار روش‌های تخمین پارامترهای Naive Bayes مهم است. جابه‌جایی به صورت اختلاف بین مقدار تخمین زده شده و مقدار واقعی تعریف می‌شود یعنی

$$bias(\hat{\theta}_k^c) = \hat{\theta}_k^c - \theta_k^c = \frac{a_k + N^c \theta_k^c}{a + N^c} - \theta_k^c = \frac{a_k - a \theta_k^c}{a + N^c}$$

در رابطه بالا، اگر $\theta_k^c > \frac{a_k}{a}$ ، امید ریاضی مقدار تخمین زده شده از مقدار واقعی کوچک‌تر خواهد بود و اگر $\theta_k^c < \frac{a_k}{a}$ ، امید ریاضی مقدار تخمین زده شده از مقدار واقعی بزرگ‌تر خواهد بود. ولی به هر حال با رشد N^c ، مقدار bias به سمت صفر میل می‌کند.

۳-۴ انحراف از معیار

فرض کنید که N_k^d تعداد تکرار ویژگی k -ام در داده‌آزمون (d) باشد و $N^d = \sum_k N_k^d$ آن‌گاه:

$$z_c = -N^d \log(a + N^c) + \sum N_k^d \log(a_k + N_k^c)$$

z_c برابر است با لگاریتم امتیازی که داده‌آزمون d از دسته c کسب می‌نماید. d به دسته‌ای که در آن بیشترین امتیاز را آورده باشد منصوب می‌شود. جمع‌های بکار رفته در عبارت بالا از هم مستقل هستند. با توجه به این که $\{N_k^d\}$ برای یک داده‌آزمون ثابت بوده ولی $\{N_k^c\}$ ممکن است در دسته‌های مختلف متفاوت باشد خواهیم داشت:

$$var(z_i) = \sum (N_k^d)^2 var(\log(a_k + N_k^c))$$

از طرفی می‌دانیم:

$$var(\log(a_k + N_k^i)) = E[(\log(a_k + N_k^i))^2] - E[\log(a_k + N_k^i)]^2$$

با در نظر گرفتن هر یک از N_k ها به صورت توزیع دوجمله‌ای با پارامتر θ_k^c خواهیم داشت:

$$E(\log(a_k + N_k^i)) = \sum_n \log(a_k + n) \binom{N}{n} (\theta_k^i)^n (1 - \theta_k^i)^{(N^i - n)}$$

رابطه بالا با استفاده از تخمین N_k^i بوسیله توزیع پواسن^{۶۴} با $\Lambda = \theta^c$ و استفاده از فرمول استرلینگ^{۶۵} برای $n!$ بدست آمده است.

بر اساس تحلیل Rennie در [۲۱]، مقدار $var(\log(1 + N_k))$ هنگامی که امید ریاضی تعداد رخ داده‌های f_k (کامین ویژگی) در داده آزمون بین ۱ تا ۲ برابر مجموعه آموزشی باشد، به مقدار بیشینه خود می‌رسد.

۴-۴ گسترش Naive Bayes

در روش Naive Bayes فرض بر این است که ویژگی‌های بکار رفته در مساله از هم مستقل هستند و در نتیجه احتمال عضویت داده در دسته برابر با حاصل ضرب احتمال هر یک از اعضای آن داده در دسته مورد نظر است. ولی در بسیاری از موارد مساله پیچیده‌تر از این حالت ساده است. برای حل چنین مسائلی از فرآیندهای مارکف گسسته به عنوان فرآیند تصادفی تولید کننده داده‌ها استفاده شده و از روش‌های ML و MAP برای تخمین پارامترهای مجهول مدل (با روشی همانند آنچه در بخش قبل توضیح داده شده بود) استفاده می‌شود. برای دسته‌بندی براین اساس، احتمال تولید شدن داده مورد پرسش را توسط مدل مفروض و پارامترهای تخمین زده شده محاسبه می‌شوند و داده در دسته‌ای که با بیشترین احتمال آن داده را تولید می‌کند دسته‌بندی می‌شود.

با روشی مشابه Naive Bayes می‌توان ثابت کرد که در صورتی که مدل درست انتخاب شود (مدل مفروض بتواند تخمین خوبی برای مدل غیرقابل دسترس ما باشد) این گسترش از Naive Bayes نیز به صورت بهینه کار می‌کند.

⁶⁴Poisson

⁶⁵Stirling's Formula

۵ روش دسته‌بندی داده‌ها بر اساس تئوری اطلاعات

برای دسته‌بندی داده‌ها در هوش مصنوعی روش‌های متنوعی وجود دارد که در فصل سوم به چند نمونه از آن‌ها اشاره شد. در این فصل تلاش داریم تا این کار را با استفاده از تئوری اطلاعات انجام دهیم.

۱-۵ علت نیاز به روش دیگر هنگامی که روش Bayes وجود دارد

همان‌طور که در ۳-۴ توضیح داده شد روش Bayes در حالت کلی دقیق‌ترین نوع دسته‌بندی کننده‌ها است. حال سوالی که در ابتدا باید بدان پاسخ داد این است که در صورت وجود یک روش دقیق مانند روش Bayes آیا نیازی به روش دیگری مبتنی بر احتمالات است؟ در چه مواقعی نمی‌توان از روش Bayes استفاده کرد؟ در چه مواقعی روش Bayes از کارآیی خوبی برخوردار نیست؟ در چهار حالت زیر، روش Bayes جواب مناسب تولید نمی‌کند:

۱. اگر مدل درست انتخاب نشده باشد: در بسیاری از موارد اطلاع دقیقی از جزئیات مدل در دست نیست و مدل باید با توجه به داده‌های تولید شده توسط مدل تخمین زده شود. در این حالت انتخاب پارامتر تعداد حالات و چگونگی ارتباط حالات با یکدیگر پارامتر تعیین کننده‌ای است. در مواردی وضعیت از این مورد هم بدتر است، بدین معنا که در حقیقت منابع تولیدکننده داده‌ها واقعا یک فرآیند تصادفی (با تعداد حالات متناهی) نیستند، اما به دلیل استفاده از روش Bayes ناگزیر به مدل کردن آن با یک فرآیند تصادفی هستیم.

۲. اگر مجموعه آموزش جامعیت کافی نداشته باشد: در این حالت حتی اگر مدل درست انتخاب شده باشد به دلیل جامع نبودن مجموعه آموزش به کار رفته احتمالات جابه‌جایی بین دو حالت درست محاسبه نمی‌شود. حالت وخیم‌تر هنگامی است که یک یا چند ویژگی از مجموعه ویژگی‌ها، در عبور از یک حالت به حالت دیگر در مجموعه آموزش رخ نداده باشند.

۳. احتمال وجود اختلال در داده‌های مجموعه آموزشی و داده‌های آزمون: اختلال^{۶۶} معمولا به علت دقیق نبودن ابزارهای نمونه‌گیری رخ می‌دهد و در برخی از کاربردها

⁶⁶Noise

غیرقابل اجتناب هستند. به عنوان مثال مساله تشخیص چهره افراد را در نظر بگیرید. در این مساله، مجموعه آموزشی از تعدادی عکس از افراد مختلف که با کد صاحب عکس مشخص شده تشکیل می‌شود. در کاربرد، عکسی که از چهره یک فرد گرفته شده به عنوان ورودی به سامانه داده می‌شود و سامانه باید شبیه‌ترین فرد را پیدا کند. از معمول‌ترین مشکلات در این چنین مساله‌هایی پراکندگی رنگ‌ها، تغییر رنگ‌ها، تداخل رنگ‌ها و مانند این‌ها است. شاید قسمتی از این اختلالات را می‌توان با پیش‌پردازش^{۶۷} حذف کرد ولی باز هم قسمتی از اختلال باقی می‌ماند. در این صورت، اگر اندازه مجموعه آموزش نسبتاً بزرگ باشد (معمولاً چنین است) احتمالاً پارامترهای مدل مورد نظر با دقت خوبی قابل تخمین خواهند بود. ولی در صورتی که در داده‌آزمون هم اختلال پیش آید آن‌گاه به علت کوچک بودن احتمال خطا (این یک فرض معمول است) احتمال عضویت داده در هر یک از دسته‌ها بسیار کوچک خواهد بود و چه بسا برابر با صفر باشد.

۴. اگر مجموعه در حال تغییر باشد: گاهی از اوقات دسته‌ها در حال تغییر هستند. بدین معنا که به مرور زمان برخی از احتمالات در فرآیند تصادفی متناظر با دسته‌ها، تغییر می‌کنند. به طور خاص می‌توان این مساله را در مورد خبرها دید.

به عنوان مثال، تا چندی پیش «انرژی هسته‌ای» احتمالاً در حوزه اخبار سیاسی ایران اصلاً ظاهر نشده بود ولی به تازگی این واژه به صورت بسیار زیادی در این حوزه خبری (دسته) ظاهر می‌شود. به عنوان مثالی دیگر، تا چندی پیش، واژه «قانا» با احتمال بیشتری در حوزه اخبار گردشگری یافت می‌شد اما برای مدتی این واژه به صورت فراوان در اخبار سیاسی ظاهر شد و اکنون وجود آن در حوزه‌های مختلف به صورت قبل بازگشته است.

در صورتی که تغییر احتمالات دسته‌ها، بسیار زیاد باشد احتمالاً نمی‌توان روش موثری برای بهبود روش‌های موجود ارائه داد ولی در مثال‌های بالا، واژه‌های کلیدی دسته‌ها (خبرهای سیاسی) مانند «سازمان ملل»، «سیاست» و «ایران» تغییر نکرده است ولی احتمال یک یا چند واژه در این حوزه تغییر کرده است.

مورد ۱، به علت عدم دانش کافی کاربر از مساله رخ می‌دهد و نمی‌توان راه حل موثری برای این مشکل ارائه داد. مورد ۲، در حالت کلی یک مشکل اساسی در مساله دسته‌بندی است زیرا تنها دانش ما از داده‌های عضو یک دسته، مجموعه آموزشی است اما اگر مجموعه آموزشی شامل اطلاعات نسبتاً جامعی از فرآیند باشد ولی بخشی از جزئیات مربوط به دسته

⁶⁷Preprocessing

را در خود نداشته باشد می‌توان راه‌حلی برای رفع این مشکل ارائه داد. یکی از مشکلات اساسی روش‌های Bayes در رویارویی با اختلال‌هاست [۲۲]. به دلیل این‌که روش‌های Bayes بسیار وابسته به احتمالات تخمین‌زده شده هستند معمولاً نمی‌توانند در مورد ۴، پاسخ مناسبی ارائه دهند.

یکی دیگر از مواردی که روش‌های Bayes در حل آن به مشکل برمی‌خورند حالت زیر است: حالتی را در نظر بگیرید که در آن یک دنباله $T = f_1, f_2, \dots, f_n$ داده شده باشد. ماشین M به عنوان مدل دسته C فراخوانی می‌شود و در گام i تولید T در مجموعه حالات $\{S_{i_1}, S_{i_2}, \dots, S_{i_k}\}$ قرار داریم و هیچ یک از S_i ها با ویژگی (کارکتر) f_i به حالت دیگری خروجی نداشته باشند. در این حالت باید چه تصمیمی در مورد مقدار احتمال تولید T توسط ماشین گرفته شود؟ آیا این احتمال واقعا صفر است؟ آیا عدم وجود یال به علت ناقص بودن و یا ناکافی بودن مجموعه آموزشی است؟ به این سوالات به سه طریق می‌توان پاسخ داد:

۱. ساده‌ترین روش برخورد با این مساله صرف‌نظر کردن از f_i است. به این ترتیب که f_i را از T حذف کنیم و به جای محاسبه احتمال T ، احتمال T' (دنباله جدید) را محاسبه کنیم. همان‌طور که گفته شد این راه‌حل ساده‌ترین روش برخورد با این مشکل است و احتمالاً جواب مناسبی برای این مساله نیست (هیچ توجیه علمی نیز برای آن وجود ندارد).

۲. روش دیگر برخورد با این مساله، صفر در نظر گرفتن احتمال f_i در آن حالت است. این کار را می‌توان این‌گونه توجیه کرد که فرض کردیم مجموعه آموزشی ما یک مجموعه کامل و پوشا بوده و عدم وجود یک ویژگی در مجموعه آموزشی به معنای صفر بودن امکان آن ویژگی است. این روش در صورت کامل بودن مجموعه آموزشی درست است ولی چگونه می‌توان از جامع بودن مجموعه آموزشی اطمینان حاصل کرد؟ اگر مجموعه کلاس در حال تغییر باشد چطور؟ در این صورت ممکن است T شامل یک ویژگی جدید که تا به حال در مجموعه آموزشی وجود نداشته است باشد. در فصل بعدی این مورد را به همراه کاربرد عملی آن به طور کامل‌تر توضیح می‌دهیم.

۳. روش دیگر استفاده از یک راه‌حل بینابینی است. به این ترتیب که در صورت عدم وجود یال خروجی از یک حالت مشخص، احتمال آن را برابر با عددی از پیش تعیین‌شده مانند v قرار دهیم. مقدار v باید متناسب با مقدار جامعیت مجموعه آموزشی باشد. با قرار دادن $v = 0$ راه اول و $v = 1$ راه دوم بدست می‌آید. نکته قابل ذکر این است که در این روش، روش Bayes به صورت درست و کامل پیاده شده است زیرا در روش Bayes، در حقیقت باید $P(T|M, H)$ محاسبه شود که در آن H مجموعه‌ای از فرض‌هاست. H می‌تواند شامل این فرض باشد که ”در صورت عدم امکان رخ دادن

یک ویژگی در یک حالت خاص با فلان احتمال به فلان حالت برو". مقدار v را می‌توان با تنظیم پارامترهای a_i و a معرفی شده در بخش سوم بدست آورد. در حقیقت مشکل اساسی این روش چگونگی تعیین این پارامتر است.

در ادامه یک روش جایگزین ارائه می‌شود و ثابت می‌کنیم که در صورت وجود شرایط خاصی (که این شرایط، شرایط محدود کننده‌ای نیستند) مانند روش Bayes کار می‌کند و علاوه بر این می‌تواند مشکلات اشاره شده در ۲ و ۳ و ۴ را حل کند و علاوه بر آن در مورد بالا نیز، نیاز به تنظیم پارامتر خاصی ندارد.

۲-۵ مدل احتمالات

همان‌طور که در معرفی روش Bayes گفته شد، روش Bayes و سایر روش‌های احتمالاتی نیاز به یک فرآیند تصادفی برای مدل کردن تولید رشته‌ها دارند. معمولاً از مدل‌های مارکف (زنجیره‌های مارکف، پردازش‌های مارکف گسسته) برای این منظور استفاده می‌شود. در روش‌های احتمالاتی دسته‌بندی داده‌ها، هر داده توسط یک دنباله از ویژگی‌ها^{۶۸} بیان می‌شود. این دنباله‌ها الزاماً برگشت پذیر و یک به یک نیستند. به ازای هر دسته، یک سامانه در نظر گرفته می‌شود که در سامانه همه داده‌های عضو آن دسته با بیشترین احتمال تولید می‌شوند. در این روش‌ها مفهوم عضویت با مفهوم احتمال عضویت جایگزین شده است و پس از محاسبه احتمال‌های عضویت، در مورد هر دسته فرآیند تصمیم‌گیری انجام می‌شود. از روش‌های معمول برای این کار می‌توان به آستانه‌گیری^{۶۹} [۲۵] و یا ماکزیمم‌گیری اشاره کرد. ماشین توصیف کننده یک دسته، دنباله‌های مربوط به داده‌های آن دسته را با احتمال بیشتری تولید می‌کند.

سامانه موردنظر برای هر دسته معمولاً یک پردازش مارکف گسسته است که آن را به صورت زیر بیان کنیم:

۱. الفبا: ماشین مورد استفاده در این روش‌ها، باید دنباله‌ای از ویژگی‌ها را تولید کند. در نتیجه، باید به ازای هر ویژگی از مجموعه ویژگی‌ها یک حرف الفبا اختصاص یابد. مجموعه الفبا را با $F = \{f_1, f_2, \dots, f_L\}$ نشان می‌دهیم.

۲. حالت‌ها: مجموعه حالات را با $S = \{S_1, S_2, \dots, S_n\}$ نشان می‌دهیم.

⁶⁸feature

⁶⁹Thresholding

۳. انتقال‌ها: هر حالت مانند S_i با احتمالی بزرگ‌تر یا مساوی صفر به تمام حالت‌های دیگر S_j متصل می‌باشد و در نتیجه استفاده از این یال، کارکتر اختصاص داده شده به آن یال در خروجی پدیدار می‌شود. احتمال انتقال از حالت i به حالت j و تولید کارکتر f را با $p_{i,j}^f$ نشان می‌دهیم. با تثبیت تعداد حالات و نحوه اتصال آن‌ها (وجود یا عدم وجود یال حامل کارکتر f بین S_j و S_i)، می‌توان احتمال انتخاب شدن یال مورد نظر ($p_{i,j}^f$) را برحسب مجموعه آموزشی بدست آورد. $n_{i,j}^f$ را تعداد کارکترهای f موجود در مجموعه آموزش که برای تولید آن‌ها ماشین باید از حالت S_i و S_j با کارکتر f برود می‌نامیم. $n_i^f = \sum_j n_{i,j}^f$ را برابر با تعداد دفعاتی که ماشین برای تولید مجموعه آموزشی، از حالت S_i با کارکتر f خارج می‌شود قرار می‌دهیم. به همین ترتیب $n_i = \sum_f n_i^f$ را برابر با تعداد دفعاتی که ماشین برای تولید مجموعه آموزش از S_i خارج می‌شود قرار می‌دهیم. کل تعداد کارکترهای مجموعه آموزش را با $N = \sum_i n_i$ نشان می‌دهیم. تعداد رویداد کارکتر f در مجموعه آموزش را با n_M^f و تعداد رخداد‌های f در رشته d با n_d^f نشان می‌دهیم.

در ادامه در برخی از موارد از نمادهای معرفی در این بخش به عنوان یک تابع استفاده شده است که آرگومان ورودی آن می‌تواند یک مجموعه از بردارها و یا یک بردار باشد.

محاسبه Entropy یک پردازنده مارکف گسسته

Entropy یک سامانه تولید کننده مجموعه‌ای از رشته‌ها با احتمالات مختلف برابر با Entropy تمام رشته‌های قابل تولید آن سامانه با در نظر گرفتن احتمال هر رشته است. در حالت کلی هر پردازنده می‌تواند تعداد نامتناهی رشته تولید کند و در نتیجه نمی‌توان به تولید تمام رشته‌ها پرداخت و سپس Entropy را محاسبه کرد. روش مستقیمی برای محاسبه Entropy مجموعه تمامی رشته‌های قابل تولید توسط یک پردازنده مارکف گسسته به صورت زیر وجود دارد:

$$E = - \sum_i \pi_i \sum_j \sum_c p_{i,j}^c \log p_{i,j}^c$$

که در آن $p_{i,j}^c$ احتمال انتقال از حالت S_i به حالت S_j و تولید کارکتر c است و π_i احتمال قرار داشتن در حالت S_i است. این پارامترها در بخش ۲ توضیح داده شده بود. برای راحت‌تر شدن اثبات‌ها، فرض کرده‌ایم هر حالت M با تولید کارکتر c فقط به یک حالت دیگر می‌تواند منتقل می‌شود، به راحتی می‌توان تحقیق کرد که این محدودیت هیچ کاهشی در قدرت این ماشین‌ها ایجاد نمی‌کند.

چند نکته در مورد چگونگی استفاده از پرده‌های مارکف گسسته به عنوان منبع

همان‌گونه که اشاره شد هدف اصلی در مساله دسته‌بندی، بدست آوردن یک تعمیم بر اساس داده‌های موجود در مجموعه آموزشی برای کل دامنه مساله می‌باشد، به همین دلیل، هنگامی که می‌خواهیم یک مساله را به صورت احتمالاتی و با استفاده از منابع اطلاعاتی حل کنیم از یک منبع اطلاعاتی که توانایی تولید تعداد نامتناهی (شمارا) رشته را دارد استفاده می‌کنیم. از این پس منابع مورد استفاده ما هیچ رشته متناهی تولید نمی‌کنند و فقط رشته‌های نامتناهی تولید می‌کنند. این بدان معناست که هر راس موجود در گراف نمایش منبع در حداقل یک دور قرار دارد.

معمولا برای منابع اطلاعاتی، حالت شروعی در نظر گرفته نمی‌شود و حالت i با احتمال Π_i (احتمال ایستایی در حالت i که در بخش ۲ تعریف شد) می‌تواند حالت شروع منبع مورد نظر باشد. منبع S برای تولید رشته نامتناهی T ، از یک حالت به صورت تصادفی و با توزیع احتمال Π شروع می‌کند و با توجه به توزیع احتمال انتقال در حالت کنونی فرآیند، یکی از انتقال‌های ممکن (انتقالی که مبدا آن حالت کنونی باشد) را انتخاب می‌کند و پس از اضافه کردن کارکتر مرتبط با آن انتقال به رشته T ، حالت کنونی فرآیند را به روز می‌کند.

گوییم رشته (متناهی و یا نامتناهی) T توسط منبع S تولید می‌شود اگر در هر رشته (نامتناهی) تولیدی توسط S مانند T' ، رشته T به عنوان زیررشته بی‌نهایت بار تکرار شده باشد. به صورت مشابه تولید شدن مجموعه‌ای از رشته‌ها توسط یک منبع به این معناست که هر عضو مجموعه مورد نظر در T' بی‌نهایت بار تکرار شود.

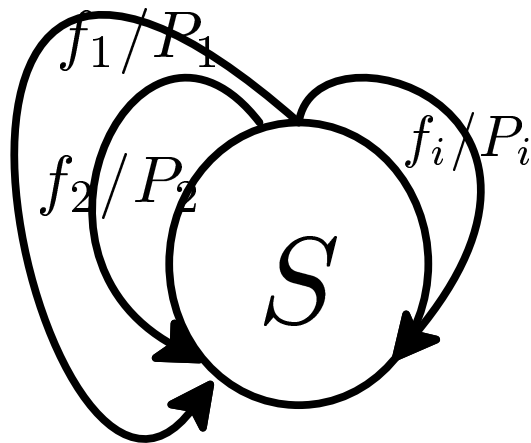
احتمال تولید شدن داده (بردار) d توسط ماشین M را برابر با احتمال تولید شدن $|d|$ کارکتر اول رشته d توسط M ($P(d|M)$) تعریف می‌کنیم. برای محاسبه آن، نیاز به تعاریف زیر داریم:

- $n_{i,j,d}^c$: تعداد دفعاتی که در یک مسیر از پیش تعیین شده در یک فرآیند برای تولید d از حالت i به حالت j رفته‌ایم.

- T_d^M : تمامی دنباله‌های $|d| + 1$ عنصری که هر عنصر آن برابر با نام یک حالت است و با دنبال کردن این دنباله بتوان $|d|$ کارکتر اول d را تولید کرد. به علت فرض ما مبنی بر این که در هر حالت فقط یک انتقال با کارکتر c از آن حالت خارج می‌شود می‌توان هر عضو $t \in T_d$ را به صورت یکتا توسط حالت اولیه t تعریف کرد.

- $P(d|t)$: که در آن $t \in T_d^M$ است برابر با احتمال تولید شدن $|d|$ کارکتر اول d در مسیر t است و برابر است با:

$$P(d|t) = \prod_i p_{t_i t_{i+1}}^{d_i}$$



شکل ۵-۱: فرآیند مربوط به حالتی که در آن ترتیب و چگونگی قرارگرفتن ویژگی‌ها مهم نباشد و فقط وجود یا عدم وجود ویژگی‌ها مهم باشد

است $P(d|M)$ از رابطه زیر بدست می‌آید:

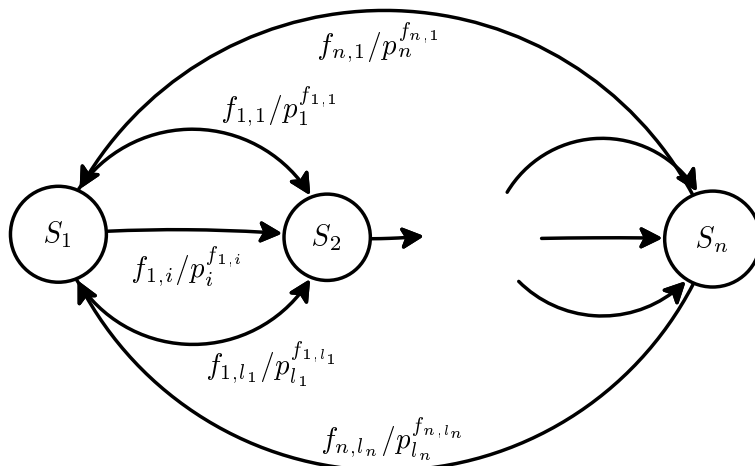
$$p(d|M) = \sum_{t \in T_d^M} \Pi_{t_0} p(d|t)$$

در ادامه به توضیح چند حالت خاص از پرده‌های مارکف گسسته که برای کار ما مناسب است می‌پردازیم:

۱. حالتی که در آن ترتیب و چگونگی قرارگرفتن ویژگی‌ها مهم نباشد و فقط وجود یا عدم وجود ویژگی‌ها مهم باشد: در این حالت به جای دنباله‌ای از ویژگی‌ها، با یک مجموعه از ویژگی‌ها سروکار داریم. پرده‌مارکف زیر برای بیان این حالت مناسب است.

در این حالت مدل دارای یک راس S و L یال است و باید مقدار L پارامتر یعنی احتمال‌های مجهول، $p_{S,S}^L$ تعیین شود. باید نشان دهیم که این مدل توانایی نمایش مساله‌هایی با خصوصیت مورد بحث را دارد. برای این کار کافی است نشان دهیم هر مجموعه از ویژگی‌ها (که به صورت یک دنباله نشان داده شده است) و در نتیجه هر مجموعه آموزشی در این فرآیند قابل نمایش است. به علت وجود تنها یک حالت، بدیهی است که به ازای هر دنباله متناهی از داده‌ها می‌توان آن را در صورت وجود عناصر آن دنباله در فرآیند تولید کرد. از کاربردهای این مدل می‌توان به دسته‌بندی متون^{۷۰} اشاره کرد. این گونه فرآیندها برای مدل‌سازی مساله‌هایی که در آن ویژگی‌ها به صورت مستقل از هم اتفاق می‌افتند دقیقاً هم‌خوانی دارد همچنین می‌توان در

⁷⁰Text Classification



شکل ۵-۲: حالتی که در آن ترتیب عناصر مهم است و هر ویژگی بسته به مکان آن ویژگی احتمال به خصوصی دارد

مواردی که تعداد ویژگی‌ها زیاد است و یا اندازه مجموعه آموزشی به اندازه کافی بزرگ نیست که بتوان بر اساس آن همه احتمالات لازم را به صورت دقیق پیدا کرد از آن‌ها استفاده نمود.

برای تخمین زدن پارامترهای این گونه فرآیندها، می‌توان به صورت زیر عمل کرد:

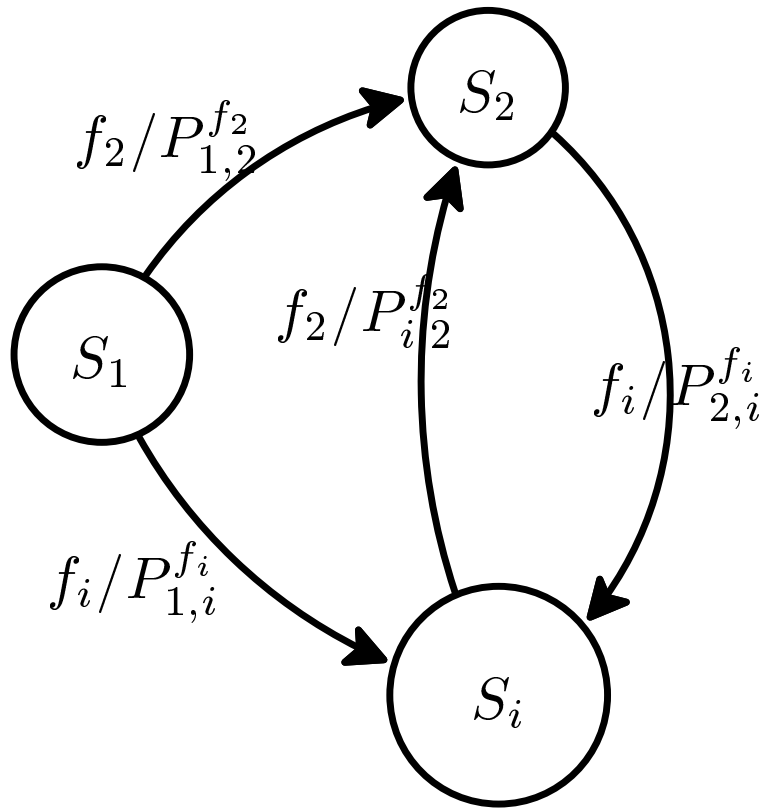
اگر n_f تعداد تکرار ویژگی f در مجموعه آموزشی و $n = \sum_f n_f$ باشد داریم $p_f = \frac{n_f}{n}$

۲. حالتی که در آن ترتیب عناصر مهم است و هر ویژگی بسته به مکان آن ویژگی احتمال به خصوصی دارد: در این حالت واقعا با مفهوم دنباله‌ای از ویژگی‌ها (با طول ثابت l) سروکار داریم که برد مولفه i -ام دنباله (F_i) ، همیشه در مجموعه $F_i = \{f_{i,j}\}$ است. در حقیقت F_i مجموعه تمام ویژگی‌های ممکن در مکان i -ام است. برای سادگی فرض می‌کنیم $\forall i \neq j : F_i \cap F_j = \emptyset$.

پردازه شکل ۵-۲ برای بیان این حالت پیشنهاد می‌شود:

در این حالت، مدل دارای l راس و $\sum |F_i|$ یال است و باید مقدار $|F_i|$ پارامتر احتمالات مجهول p_i^f تعیین شود. برای نشان دادن توانایی این گونه از پدازه‌ها برای توصیف مساله‌هایی به صورت بالا نشان می‌دهیم که هر دنباله $T = \langle f_{1,t_1}, f_{2,t_2}, \dots, f_{l,t_l} \rangle$ توسط یک پدازه از نوع بالا قابل تولید است.

به دلیل وجود تنها یک دور در گراف فرآیند بالا، پس از حداکثر تولید l کارکتر به حالت S_1 می‌رسیم. از S_1 می‌توان با تولید f_{1,t_1} به S_2 رفت و به همین ترتیب از S_i با تولید f_{i,t_i} به حالت S_{i+1} می‌رویم.



شکل ۳-۵: حالتی که در آن ترتیب عناصر مهم است و هر ویژگی بسته به عنصر قبلی دارای احتمال متفاوتی است

برای محاسبه احتمال‌های مجهول پردازش می‌توان از روش زیر استفاده کرد:

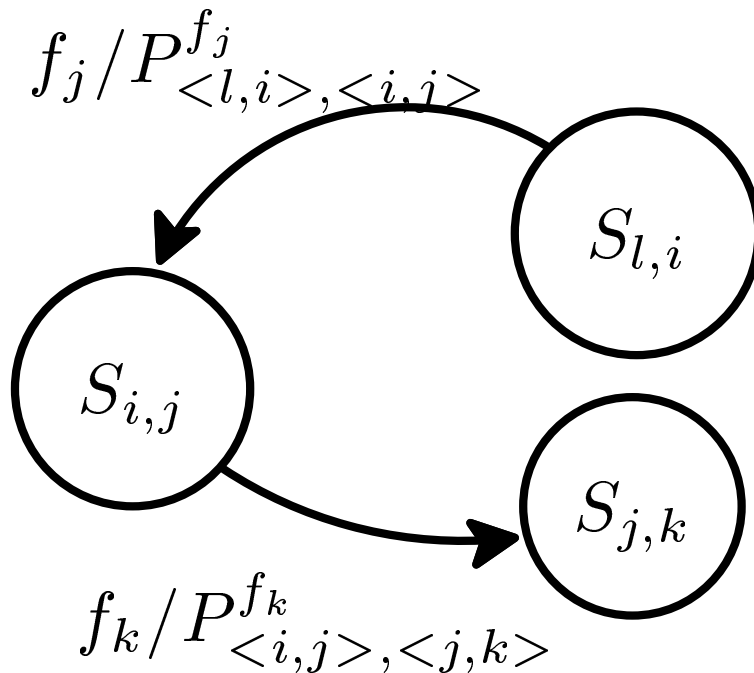
اگر n_i^f تعداد تکرار ویژگی f در مجموعه آموزشی در مکان i -ام باشد و $n_i = \sum_f n_i^f$ باشد داریم $p_i^f = \frac{n_i^f}{n_i}$

فرض اشتراک تهی برای بردهای مولفه‌های مختلف یک شرط ضروری نیست ولی به خوانایی و درک بهتر از مدل منتهی می‌شود.

۳. حالتی که در آن ترتیب عناصر مهم است و هر ویژگی بسته به عنصر قبلی دارای احتمال متفاوتی است: در این حالت نیز واقعا با مفهوم دنباله‌ای از ویژگی‌ها سروکار داریم، اما این بار هر عنصر به جای وابستگی به مکان خود به عنصر قبلی خود وابسته است. زنجیره مارکف زیر برای بیان این حالت مناسب است.

در این حالت مدل دارای L راس (هر راس نشان‌گر آخرین کارکتر تولید شده است) و L^2 یال است و باید مقدار L^2 پارامتر احتمالات مجهول $p_{i,j}^f$ تعیین شود.

در گراف شکل ۳-۵ همه راس‌ها در حداقل یک دور قرار دارند و هر راسی از هر راس



شکل ۴-۵: حالتی که در آن ترتیب عناصر مهم است و هر ویژگی بسته به دو عنصر قبلی دارای احتمال متفاوتی است

دیگر قابل دسترسی است. احتمال انتقال از حالت i به حالت j (و در نتیجه تولید کارکتر f_j) را با $p_{i,j}^{f_j}$ نشان می‌دهیم. بدیهی است که در حالت کلی $\forall i, j, k : k \neq j \ p_{i,j}^{f_j} \neq p_{i,k}^{f_k}$ به راحتی می‌توان مشاهده کرد که هر دنباله‌ای از ویژگی‌ها (کارکترها) توسط فرآیند فوق قابل تولید است

روش مشابهی مانند آنچه در بالا بدان اشاره شد را می‌توان برای این نوع ماشین‌ها به کار برد و به این ترتیب عمل کرد:

برای تولید هر دنباله توسط فرآیند فوق فقط یک مسیر وجود دارد و در نتیجه می‌توان به راحتی $n_{i,j}^{f_j}$ را که تعداد دفعاتی است که در دنباله پس از ویژگی f_i ، ویژگی f_j آمده است روی تمام دنباله‌های موجود در مجموعه آموزشی محاسبه کرد و سپس $n_i = \sum_j n_{i,j}^{f_j}$ که برابر با تعداد دفعاتی است که فرآیند برای تولید مجموعه آموزشی وارد حالت i می‌شود. حال می‌توان $p_{i,j}^{f_j} = \frac{n_{i,j}^{f_j}}{n_i}$ را تعریف کرد.

۴. حالتی که در آن ترتیب عناصر مهم است و هر ویژگی بسته به دو عنصر قبلی دارای احتمال متفاوتی است: زنجیره مارکف زیر برای بیان این حالت مناسب است.

در این حالت، مدل دارای L^2 راس و L^3 یال است و باید مقدار L^3 پارامتر احتمالات مجهول $p_{<l,i>, <i,j>}^{f_j}$ تعیین شود.

روش مشابهی مانند آنچه در بالا بدان اشاره شد را می‌توان برای این نوع ماشین‌ها به کار برد.

یک نکته قابل توجه این است که می‌توان از هر یک از فرآیندهای معرفی شده در بالا برای توصیف یک مساله استفاده کرد اما در صورت استفاده از یک مدل نامناسب احتمالات بدست آمده دیگر کاربردی نیستند. به عنوان مثال، یک مجموعه آموزشی شامل ۱۰۰ نمونه از دنباله‌ای به صورت $\langle f_1, f_2, \dots, f_L \rangle$ را در نظر بگیرید. اگر بخواهیم این دنباله را براساس فرآیندی مانند فرآیند ۱ توصیف کنیم خواهیم داشت $p_f = \frac{1}{L}$. و اگر آن را براساس فرآیندی مانند فرآیند ۲ ($l = L$) توصیف کنیم بردار توزیع احتمال p به صورت زیر خواهد بود:

$$p_j^{f_i} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

و اگر آن را با فرآیندی مانند فرآیند ۳ تخمین بزنیم خواهیم داشت

$$p_{i,j}^{f_j} = \begin{cases} 1 & \text{if } j = i + 1 \text{ or } (j = 0 \text{ and } i = L) \\ 0 & \text{otherwise} \end{cases}$$

۳-۵ چگونگی محاسبه Entropy برای چند فرآیند خاص

طبق تعریف ارائه شده در بخش ۱ برای فرآیندهای مارکف و تعمیم آن فرآیندهای مارکف گسسته، می‌توان Entropy یک فرآیند را در حالت کلی می‌توان به صورت زیر محاسبه کرد:

$$E_M = - \sum_{s \in States} \pi_s \sum_{s' \in States} \sum_{f \in F} p_{s,s'}^f \log p_{s,s'}^f$$

در برخی از فرآیندهای خاص می‌توان، یک رابطه ساده‌تر برای محاسبه E_M بدست آورد که در ادامه به بررسی آن‌ها می‌پردازیم:

حالت اول

در این گونه فرآیندها، فقط یک حالت موجود است و بنابراین بردار π برابر با $\langle 1 \rangle$ خواهد بود و با توجه به رابطه‌های مربوط به محاسبه p_f داریم:

$$\begin{aligned} E_M &= - \sum_{s \in States} \pi_s \sum_{s' \in States} \sum_{f \in F} p_{s,s'}^f \log p_{s,s'}^f \\ &= - \sum_{f \in F} p_f \log p_f = - \sum \frac{n_f}{n} \log \frac{n_f}{n} = \log n \sum \frac{n_f}{n} - \sum \frac{n_f}{n} \log n_f \end{aligned}$$

$$= \log n - \frac{1}{n} \sum n_f \log n_f$$

$$\Rightarrow E_M = \log n - \frac{1}{n} \sum n_f \log n_f \quad (5-3)$$

حالت دوم

در این گونه فرآیندها، l حالت موجود است و از هر حالت فقط به حالت بعدی می‌رویم و در نتیجه احتمال حضور در هر حالت برابر با احتمال حضور در سایر حالات است یعنی $\pi_i = \frac{1}{l}$ و با توجه به رابطه‌های مربوط به محاسبه p_i^f داریم:

$$\begin{aligned} E_M &= - \sum_{s \in States} \pi_s \sum_{s' \in States} \sum_{f \in F} p_{s,s'}^f \log p_{s,s'}^f \\ &= - \sum_{s \in States} \pi_s \sum_{f \in F} p_i^f \log p_i^f = -\frac{1}{l} \sum \sum \frac{n_i^f}{n_i} \log \frac{n_i^f}{n_i} \\ &= \frac{1}{l} \sum_i \left[\sum_f \frac{n_i^f}{n_i} \log n_i - \sum_f \frac{n_i^f}{n_i} \log n_i^f \right] \\ &= \frac{1}{l} \sum_i \left[\log n_i - \sum_f \frac{n_i^f}{n_i} \log n_i^f \right] \end{aligned}$$

همان‌طور که در توصیف این گونه فرآیندها گفته شد $\forall i, j : n_i = n_j = n$ که در آن برابر با تعداد اعضای مجموعه‌ای است که فرآیند بر اساس آن ساخته شده است، و در نتیجه رابطه بالا به صورت زیر ساده می‌شود:

$$E_M = \log n - \frac{1}{l} \sum_i \sum_f \frac{n_i^f}{n_i} \log n_i^f \quad (5-4)$$

حالت سوم

در این گونه فرآیندها، L حالت موجود است و به علت کلی بودن ماتریس احتمالات P تعیین مقدار ویژه این ماتریس به صورت یک فرمول بسته امکان‌پذیر نیست:

$$\begin{aligned}
E_M &= - \sum_{s \in States} \pi_s \sum_{s' \in States} \sum_{f \in F} p_{s,s'}^f \log p_{s,s'}^f \\
&= - \sum_i \pi_i \sum_j \sum_f \frac{n_{i,j}^f}{n_i} \log \frac{n_{i,j}^f}{n_i} = \sum_i \pi_i \sum_j \sum_f \left[\frac{n_{i,j}^f}{n_i} \log n_i - \frac{n_{i,j}^f}{n_i} \log n_{i,j}^f \right] \\
&= \sum_i \pi_i \sum_j \sum_f \frac{n_{i,j}^f}{n_i} \log n_i - \sum_i \pi_i \sum_j \sum_f \frac{n_{i,j}^f}{n_i} \log n_{i,j}^f \\
&= \sum_i \frac{\pi_i \log n_i}{n_i} \sum_j \sum_f n_{i,j}^f - \sum_i \frac{\pi_i}{n_i} \sum_j \sum_f n_{i,j}^f \log n_{i,j}^f \\
&= \sum_i \pi_i \log n_i - \sum_i \frac{\pi_i}{n_i} \sum_j \sum_f n_{i,j}^f \log n_{i,j}^f \\
E_M &= \sum_i \pi_i \log n_i - \sum_i \frac{\pi_i}{n_i} \sum_j \sum_f n_{i,j}^f \log n_{i,j}^f \quad (5-5)
\end{aligned}$$

۴-۵ رابطه بین دسته‌بندی و Entropy

هدف نهایی در دسته‌بندی داده‌ها این است که روشی برای تشخیص تعلق و یا عدم تعلق یک داده به هر یک از دسته‌ها ارائه شود. راه‌های مختلفی برای این منظور وجود دارد که به بخشی از آن‌ها در فصل دوم اشاره شد. در دسته‌بندی معمولاً تعدادی از اعضای یک دسته در اختیار است. یکی از روش‌های تعیین عضویت می‌تواند استفاده از میزان شباهت بین داده با اعضای هر دسته باشد. به این ترتیب که، با استفاده از یک معیار به سنجش میزان شباهت داده با اعضای هر یک از دسته‌ها پردازیم و مشابه‌ترین دسته را به عنوان دسته برنده اعلام کنیم. همان‌طور که دیدیم، یکی از راه‌های محاسبه میزان شباهت استفاده از مفهوم مشابه میزان تردید (عدم قطعیت) است. در بخش بعدی به توضیح چگونگی استفاده از این مفهوم برای دسته‌بندی خواهیم پرداخت.

۵-۵ ایده کلی روش

همان‌طور که در بخش ۳-۱ بحث شد یکی از راه‌های اندازه‌گیری ریاضی میزان آشفتگی داده‌ها، استفاده از Entropy است. در صورت داشتن یک معیار برای اندازه‌گیری میزان آشفتگی داده‌ها، می‌توان از این معیار در دسته‌بندی داده‌ها به صورت زیر استفاده کرد:

در مساله دسته‌بندی داده‌ها، دسته c براساس یک مجموعه آموزشی جامع (S_c) تعریف شده است و می‌خواهیم در مورد تعلق داده d به دسته c نظر بدهیم. اگر $S'_c = S_c \cup \{d\}$ باشد با توجه به فرض ما مبنی بر جامع بودن مجموعه آموزشی، منطقی است که انتظار داشته باشیم اگر d واقعا به c متعلق باشد پارامترهای تعیین شده (در مرحله آموزش الگوریتم‌ها) توسط مجموعه های S_c و S'_c احتمالا بسیار مشابه خواهد بود. برای مقایسه دو توزیع احتمال راه‌های فراوانی موجود است که برای اطلاع از آن‌ها می‌توانید به [۱۶] مراجعه کنید. در کاربرد مورد بحث در این پایان‌نامه، به دلیل رابطه به‌خصوص دو توزیع با یکدیگر، به نظر می‌رسد که استفاده از Entropy یک معیار خوب باشد. به این ترتیب که Entropy فرآیند تولیدکننده S_c و Entropy فرآیند تولیدکننده S'_c را با هم مقایسه کنیم و در حقیقت به مطالعه میزان تغییرات در Entropy در اثر افزوده شدن داده d پردازیم. این تغییر (Δ_E) می‌تواند مثبت، منفی و یا صفر بسته به شرایط باشد که در ادامه به بررسی هر یک خواهیم پرداخت:

۱. $\Delta_E > 0$: بدین معناست که Entropy در حالتی که پارامترهای فرآیند براساس مجموعه آموزش تعیین شده بود از Entropy همان فرآیند هنگامی که پارامترهای آن براساس مجموعه آموزشی و داده d تعیین شده است کمتر است یعنی با اضافه شدن d ، مقدار آشفستگی در بین رشته‌های تولید شده توسط منبع بیشتر گردیده‌است. این می‌تواند بدان معنا باشد که داده d موجب تغییرات اساسی در پارامترهای فرآیند شده است و این تغییرات بیشتر در جهت آشفته‌تر کردن دنباله‌های تولیدی بوده است. به عنوان مثال، فرض کنید برای یک مساله فرآیند ۱ را برای مدل‌سازی انتخاب کرده‌ایم و براساس یک مجموعه آموزشی مقادیر n_f ها به صورت زیر باشد ($L = 4$):

$$n_f = \begin{cases} 100 & \text{if } f \leq 3 \\ 10 & \text{otherwise} \end{cases}$$

Entropy این فرآیند برابر با $1/740$ خواهد بود. اگر توزیع ویژگی‌ها در داده d به صورت زیر باشد:

$$n_f^d = \begin{cases} 0 & \text{if } f \leq 3 \\ 31 & \text{otherwise} \end{cases}$$

در مورد پارامترهای فرآیند حاصل از مجموعه آموزش و d خواهیم داشت که $p = \langle 0/293, 0/293, 0/293, 0/121 \rangle$ و در نتیجه Entropy فرآیند جدید برابر با $1/925$ خواهد بود.

در این حالت احتمالا داده به دسته متعلق نیست و یا احتمال عضویت داده در آن دسته به‌خصوص، کم است.

۲. $\Delta_E = 0$: بدین معناست که Entropy در حالتی که پارامترهای فرآیند بر اساس مجموعه آموزش تعیین شده بود و Entropy فرآیند هنگامی که پارامترهای آن بر اساس مجموعه آموزشی و d تعیین می‌شود با هم برابر بوده است. اگر خوش شانس باشیم این بدان معناست که توزیع احتمال فرآیند پیش و پس از افزوده شدن d اصلاً تغییری نکرده است، یعنی d توزیعی کاملاً مشابه با مجموعه آموزش داشته است. در مثال بالا، اگر توزیع ویژگی‌ها در d به صورت زیر باشد:

$$n_f^d = \begin{cases} 10 & \text{if } f \leq 3 \\ 1 & \text{otherwise} \end{cases}$$

خواهیم داشت $\Delta_E = 0$.

۳. $\Delta_E < 0$: بدین معناست که Entropy در حالتی که پارامترهای فرآیند بر اساس مجموعه آموزش تعیین شده بود از Entropy فرآیند هنگامی که پارامترهای آن بر اساس مجموعه آموزشی و d تعیین می‌شود بیشتر است. یعنی با اضافه شدن d ، مقدار آشفستگی در بین رشته‌های تولید شده توسط منبع کمتر گردیده است و در حقیقت یک نظم جدید با اضافه شدن d به S کشف شده است. این می‌تواند بدان معنا باشد که داده d موجب تغییراتی (اساسی) در پارامترهای فرآیند شده است و این تغییرات بیشتر در جهت منظم‌تر کردن دنباله‌های تولیدی بوده است. اگر در مثال بالا توزیع ویژگی‌ها در d به صورت زیر باشد:

$$n_f^d = \begin{cases} 10 & \text{if } f \leq 3 \\ 0 & \text{otherwise} \end{cases}$$

در مورد پارامترهای فرآیند حاصل از مجموعه آموزش و d خواهیم داشت که $p = \langle 0/322, 0/322, 0/322, 0/032 \rangle$ و در نتیجه Entropy فرآیند جدید برابر با $1/738$ خواهد بود.

خواهیم داشت $\Delta_E < 0$.

چگونگی دسته‌بندی

برای تعیین دسته یک داده مجهول در بین چند دسته داده‌شده می‌توان از روش زیر استفاده کرد:

۱. به ازای هر دسته مانند c یک پدازه مارکف گسسته با مدل مناسب در نظر بگیر و بر اساس مجموعه آموزشی متناظر آن دسته (S_c) ، احتمالات انتقال بین حالت‌ها را

شکل ۵-۵: $\{ < ۱۰۰ >, < ۰۱۱ > \}$ یک فرآیند با ۹ حالت برای مجموعه داده‌های

محاسبه کرده و مقدار Entropy پردازنده (E_{Old}^c) را ذخیره کن.

الف- مجموعه $S'_c = S_c \cup \{d\}$ را تشکیل بده و به محاسبه مجدد احتمالات انتقال بین حالات بدون ایجاد تغییر در مدل پردازنده و سپس Entropy مدل جدید (E_{New}^c) را محاسبه کن.

ب- مقدار Δ_E^c را برابر با $E_{New}^c - E_{Old}^c$ قرار بده.

۲. داده مجهول متعلق به کلاسی است که مقدار Δ_E^c آن کمترین باشد.

چگونگی محاسبه

در الگوریتم بالا، قسمت ۱ فقط یک بار اجرا می‌شود و به ازای هر دسته E_{Old}^c محاسبه می‌شود. در حقیقت این مرحله، همان مرحله آموزش الگوریتم است. قسمت‌های ۱-الف و ۱-ب به ازای هر داده مورد آزمون فراخوانی می‌شود. از آنجا که معمولاً یک سیستم برای دسته‌بندی یک بار آموزش داده شده و پس از آن برای تشخیص استفاده می‌شود توجه اصلی در پیچیدگی این نوع سیستم‌ها، پیچیدگی مرحله آزمون آنهاست [۲۳]. در کاربرد مورد بحث ما، میزان آشفتگی یک مجموعه از داده‌ها به فرآیند تولید کننده داده‌ها وابسته است. برای مثال اگر بخواهیم Entropy مجموعه داده‌های $\{ < ۱۰۰ >, < ۰۱۱ > \}$ را برای یک فرآیند با یک حالت محاسبه کنیم در حقیقت فرآیند مورد نظر یک فرآیند کاملاً تصادفی خواهد بود (با احتمال $\frac{1}{2}$ یک و با $\frac{1}{2}$ صفر تولید می‌کند)، ولی اگر تلاش کنیم که این مجموعه را توسط یک مدل با ۹ حالت به مانند ۵-۵ تولید کنیم Entropy مجموعه بسیار کمتر خواهد بود. باید توجه داشت که علت این اختلاف در این است که ما مجموعه مورد نظر را به عنوان نمونه داده‌های تولید شده توسط یک فرآیند در نظر گرفته‌ایم و در نتیجه تخمین ما از Entropy مجموعه آموزش بسته به چگونگی محدودیت مجموعه و فرآیند تغییر می‌کند.

۵-۶ اثبات هم‌ارزی روش مبتنی بر تئوری اطلاعات با Naive Bayes

همان‌طور که گفته شد دسته‌بندی براساس Naive Bayes، در حالتی که اطلاعات کافی در اختیار داشته باشیم (مساله توسط Naive Bayes قابل حل باشد)، بهترین راه است. در این بخش تلاش خواهیم کرد تا روش ارائه شده در بالا را با روش Naive Bayes مقایسه کنیم و

ثابت کنیم در حالات اول و دوم که مساله می‌تواند توسط Naive Bayes حل شود روش ارائه شده هم مانند Naive Bayes عمل می‌کند. این بدان معناست که این روش در حالات اول و دوم به نوعی یک گسترش برای Naive Bayes است، و در (حداقل برخی از) حالاتی که Bayes، خیلی خوب کار نمی‌کند نتایج خوبی می‌دهد که با مثال‌هایی در دو بخش بعدی به بررسی این موضوع خواهیم پرداخت.

برای این منظور فرض کنید برای هر دسته c یک پردازنده مناسب مانند M_c با استفاده از S_c تشکیل شده باشد. در نتیجه $P(d|S_c) = P(d|M_c)$ و $P(S_c|d) = P(M_c|d)$ همچنین $\forall M_c; n_d \ll n_{S_c}$ یعنی اندازه مجموعه‌های آموزشی بسیار بزرگ‌تر از اندازه یک داده‌است و در نتیجه می‌توانیم فرض کنیم $\forall i, j, c, M : n_{i,j,d}^c \ll n_{i,j,M}^c$. فرض فوق از آن جا ناشی می‌شود که ما در حال مقایسه شرایطی هستیم که در آن‌ها Naive Bayes خوب کار می‌کند. همان‌گونه که اشاره شد در بسیاری از کاربردها، احتمال وقوع یک دسته مشخص نیست و بنابراین از توزیع یکنواخت برای آن استفاده می‌کنند. طبق رابطه Bayes داریم

$$P(M_c|d) = \frac{P(d|M_c)P(M_c)}{P(d)}$$

و به دلیل این که عبارت $P(d)$ در همه $P(M_c|d)$ ظاهر می‌شود و همچنین با برابر در نظر گرفتن همه $P(M_c)$ ها داریم:

$$P(M_c|d) \approx P(d|M_c)$$

فرض می‌کنیم داده d و تعدادی دسته (که توسط پردازنده‌های مارکف توصیف شده‌اند) در اختیار داریم. می‌خواهیم از بین دسته‌های موجود دسته‌ای را به عنوان دسته‌ای که d بدان تعلق دارد اعلام کنیم.

طبق روش Bayes باید احتمال $P(M_c|d)$ را به ازای همه دسته‌ها محاسبه کنیم و سپس دسته‌ای که احتمال آن بیشتر بود را به عنوان جواب اعلام کنیم، اما با توجه به داده‌های مساله داریم $P(M_c|d) \approx P(d|M_c)$ و در نتیجه کافی است به محاسبه $P(d|M_c)$ بپردازیم. روش مبتنی بر Entropy، به محاسبه Δ_E^c می‌پردازد. برای اثبات هم‌ارزی این دو روش کافی است ثابت کنیم

$$\Delta_E^c = f(P(d|M_C))$$

که در آن f تابعی اکیدا نزولی است. در این صورت به ازای هر دو دسته c و c' داریم:

$$\Delta_E^c < \Delta_E^{c'} \iff P(d|M_C) > P(d|M_{C'})$$

و در نتیجه $t = \arg \min(\Delta_E^c) \iff t = \arg \max(P(d|M_C))$

حالت اول

پردازه مارکف گسسته این حالت در شکل ۵-۱، آمده است. در این حالت خاص، پردازه در حقیقت احتمال وقوع هر ویژگی را مستقل از ویژگی‌های قبلی در نظر می‌گیرد و در حقیقت می‌توان از Naive Bayes نیز در این گونه مساله‌ها استفاده کرد. برای محاسبه Δ_E^c نیاز به محاسبه Entropy برای M_c داریم و همان گونه که اشاره شد (در رابطه ۵-۳) از رابطه زیر بدست می‌آید:

$$E_{Old} = E_M(S) = \log n(S) - \frac{1}{n(S)} \sum n_f(S) \log n_f(S)$$

$$\begin{aligned} E_{New} &= E_M(S \cup \{d\}) = \log n(S \cup \{d\}) - \frac{1}{n(S) \cup \{d\}} \sum n_f(S \cup \{d\}) \log n_f(S \cup \{d\}) \\ &= \log(n(S) + n(d)) - \frac{1}{n(S) + n(d)} \sum (n_f(S) + n_f(d)) \log(n_f(S) + n_f(d)) \end{aligned}$$

و در نتیجه $\Delta_E(S, d)$ به صورت زیر محاسبه می‌شود:

$$\begin{aligned} \Delta_E^c(S, d) &= \log(n(S) + n(d)) - \sum_f \frac{n_f(S) + n_f(d)}{n(S) + n(d)} \log(n_f(S) + n_f(d)) - \log n(S) + \sum_f \frac{n_f(S)}{n(S)} \log n_f(S) \\ &= \log \frac{n(S) + n(d)}{n(S)} - \sum_f \left[\frac{n_f(S) + n_f(d)}{n(S) + n(S)} \log(n_f(S) + n_f(d)) - \frac{n_f(S)}{n(S)} \log n_f(S) \right] \end{aligned}$$

با فرض این که $n(d) \ll n(S)$ داریم:

$$\frac{n(S) + n(d)}{n(S)} = 1 + \frac{n(d)}{n(S)} \approx 1$$

$$\Rightarrow \log \frac{n(S) + n(d)}{n(S)} \approx 0$$

$$\Rightarrow \Delta_E^c(S, d) \approx - \sum_f \left[\frac{n_f(S) + n_f(d)}{n(S) + n(d)} \log(n_f(S) + n_f(d)) - \frac{n_f(S)}{n(S)} \log n_f(S) \right]$$

با توجه به فرض ما مبنی بر $N = n(S) \gg n(d)$ اگر داده و مجموعه آموزشی توسط یک مدل یکسان تولید شده باشند باید $n_f(S) \gg n_f(d)$ (این فرض در حالاتی است که Naive Bayes خوب کار می‌کند برقرار است) و در نتیجه می‌توان $\log(n_f(S) + n_f(d))$ را با $\log n_f(S)$ تخمین زد.

$$\begin{aligned}
\Rightarrow \Delta_E^c(S, d) &\approx - \sum_f \left[\left(\frac{n_f(S) + n_f(d)}{n(S) + n(d)} - \frac{n_f(S)}{n(S)} \right) \log n_f(S) \right] \\
&= - \sum_f \left[\frac{n_f(d)n(S) - n_f(S)n(d)}{n_S(n(S) + n(d))} \log n_f(S) \right] \\
&= - \sum_f \left[\frac{n_f(d)n(S)}{n_S(n(S) + n(d))} \log n_f(S) - \frac{n_f(S)n(d)}{n_S(n(S) + n(d))} \log n_f(S) \right] \\
&= - \sum_f \left[\frac{n_f(d)}{n(S) + n(d)} \log n_f(S) - \frac{n_f(S)n(d)}{n_S(n(S) + n(d))} \log n_f(S) \right]
\end{aligned}$$

در عبارت دوم در رابطه بالا یعنی $\sum_f \frac{n_f(S)n(d)}{n_S(n(S)+n(d))} \log n_f(S)$ مقادیر $n(S)$ و $n(d)$ مستقل از مقدار f هستند و در نتیجه

$$\begin{aligned}
&\sum_f \frac{n_f(S)n(d)}{n_S(n(S) + n(d))} \log n_f(S) \\
&= \frac{n(d)}{n_S(n(S) + n(d))} \sum_f n_f(S) \log n_f(S)
\end{aligned}$$

مقدار عبارت $\sum_f n_f(S) \log n_f(S)$ با شرط $\sum_f n_f(S) = n(S)$ در رابطه زیر صدق می‌کند:

$$n(S) \log n(S) - n(S) \log L \leq \sum_f n_f(S) \log n_f(S) \leq n(S) \log n(S)$$

$$\Rightarrow \frac{\log n(S) - \log L}{n(S) + n(d)} \leq \frac{\sum_f n_f(S) \log n_f(S)}{n_S(n(S) + n(d))} \leq \frac{\log n(S)}{n(S) + n(d)}$$

فرض کردیم که $n(S)$ ، عدد بزرگی است. اگر $n(S)$ طوری بزرگ باشد که $\log n(S) \ll n(S)$ آن‌گاه خواهیم داشت:

$$\begin{aligned}
\frac{\log n(S) - \log L}{n(S) + n(d)} &\leq \frac{\sum_f n_f(S) \log n_f(S)}{n_S(n(S) + n(d))} \leq \frac{\log n(S)}{n(S) + n(d)} \\
\Rightarrow 0 &\leq \frac{\sum_f n_f(S) \log n_f(S)}{n_S(n(S) + n(d))} \leq 0 \\
\frac{\sum_f n_f(S) \log n_f(S)}{n_S(n(S) + n(d))} &\approx 0
\end{aligned}$$

در کاربردهای عملی (به عنوان مثال مباحث الکترونیک) گوئیم $a \ll b$ اگر $a < \frac{b}{10}$ و در نتیجه اگر $n > 64$ باشد شرط بالا برقرار است و در نتیجه می‌توان Δ_E^c را به صورت زیر تخمین زد:

$$\begin{aligned}\Delta_E^c(S, d) &\approx C - \frac{\sum_f n_f(d) \log n_f(S)}{n(S) + n(d)} + \frac{\sum_f n_f(S) \log n_f(S) n(d)}{n_S (n(S) + n(d))} \\ &\approx C - \frac{\sum_f n_f(d) \log n_f(S)}{n(S) + n(d)}\end{aligned}$$

و همچنین برای $P(d|M)$ داریم:

$$p(d|M) = \sum_{t \in T_d^M} \pi_{t_1} p(d|t)$$

از آن جا که در این نوع فرآیندها، فقط یک مسیر برای تولید بردار $d = \langle f_1, f_2, \dots, f_l \rangle$ وجود دارد خواهیم داشت:

$$\begin{aligned}p(d|t) &= \prod_i p_{f_i}(S) = \prod_f p_f(S)^{n_f(d)} \\ \Rightarrow \log p(d|t) &= \sum_f n_f(d) \log p_f(S) \\ &= \sum_f n_f(d) \log \frac{n_f(S)}{n(S)} = \sum_f n_f(d) \log n_f(S) - \log n(S) \sum_f n_f(d) \\ &= \sum_f n_f(d) \log n_f(S) - \log n(S) n(d) = \sum_f n_f(d) \log n_f(S) - C' \\ \Rightarrow \Delta_E^c(S, d) &\approx C_1 - C_2 \log(p(d|t))\end{aligned}$$

پس تابع f که به دنبال آن بودیم بدست آمده است و در نتیجه روش ارائه شده در پایان نامه در شرایط اشاره شده می تواند مانند Naive Bayes رفتار کند.

حالت دوم

فرآیند مارکف گسسته این مدل در شکل ۵-۲، آمده است. در این حالت اگر احتمال وقوع هر ویژگی مستقل از ویژگی قبلی باشد و فقط به مکان آن ویژگی بستگی داشته باشد، می توان از Naive Bayes در این گونه مساله ها استفاده کرد. به این صورت که احتمال رخ دادن هر یک از دنباله ویژگی های با طول l را در هر فرآیند محاسبه کرد و سپس از Naive Bayes روی فرآیندهای حاصل (که در آن ها احتمال وقوع دنباله های l -تایی از هم مستقل هستند) استفاده کرد.

برای محاسبه Δ_E^c نیاز به محاسبه Entropy برای M_c داریم و همان گونه که اشاره شد از رابطه زیر بدست می آید:

$$E_{Old} = E_M(S) = \log n(S) - \frac{1}{l} \sum_i \sum_f \frac{n_i^f(S)}{n_i(S)} \log n_i^f(S)$$

$$\begin{aligned} E_{New} &= E_M(S \cup \{d\}) = \log n(S \cup \{d\}) - \frac{1}{l} \sum_i \sum_f \frac{n_i^f(S \cup \{d\})}{n_i(S \cup \{d\})} \log n_i^f(S \cup \{d\}) \\ &= \log (n(S) + n(d)) - \frac{1}{l} \sum_i \sum_f \left(\frac{n_i^f(S) + n_i^f(d)}{n_i(S) + n_i(d)} \right) \log (n_i^f(S) + n_i^f(d)) \end{aligned}$$

با توجه به این که d باید یک بردار l -تایی باشد ($n(d) = n_i(d) = 1$) داریم:

$$E_{New} = \log (n(S) + 1) - \frac{1}{l} \sum_i \sum_f \left(\frac{n_i^f(S) + n_i^f(d)}{n_i(S) + 1} \log (n_i^f(S) + n_i^f(d)) \right)$$

و در نتیجه $\Delta_E(S, d)$ به صورت زیر محاسبه می شود:

$$\begin{aligned} \Delta_E^c(S, d) &= \\ &= \log (n(S) + 1) - \frac{1}{l} \sum_i \sum_f \left(\frac{n_i^f(S) + n_i^f(d)}{n(S) + 1} \log (n_i^f(S) + n_i^f(d)) \right) - \log n(S) + \\ &\quad \frac{1}{l} \sum_i \sum_f \frac{n_i^f(S)}{n(S)} \log n_i^f(S) \\ &= \log \frac{n(S) + 1}{n(S)} - \frac{1}{l} \sum_i \sum_f \left[\frac{n_i^f(S) + n_i^f(d)}{n(S) + 1} \log (n_i^f(S) + n_i^f(d)) - \frac{n_i^f(S)}{n(S)} \log n_i^f(S) \right] \\ &= \log \frac{n(S) + 1}{n(S)} - \frac{1}{l} \sum_i \sum_f \left[\frac{n_i^{f_i}(S) + 1}{n(S) + 1} \log (n_i^{f_i}(S) + 1) - \frac{n_i^{f_i}(S)}{n(S)} \log n_i^{f_i}(S) \right] \end{aligned}$$

اگر دنباله $d = \langle f_{d_1}, f_{d_2}, \dots, f_{d_l} \rangle$ باشد داریم:

$$\begin{aligned} \Delta_E^c(S, d) &= \\ &= \log \frac{n(S) + 1}{n(S)} - \frac{1}{l} \sum_i \left[\frac{n_i^{f_{d_i}}(S) + 1}{n(S) + 1} \log (n_i^{f_{d_i}}(S) + 1) - \frac{n_i^{f_{d_i}}(S)}{n(S)} \log n_i^{f_{d_i}}(S) \right] \end{aligned}$$

با توجه به فرض ما مبنی بر بزرگ بودن $n(S)$ خواهیم داشت:

$$\log \frac{n(S) + 1}{n(S)} \approx \log 1 = 0$$

و مانند بخش قبل اگر بتوان $\log(n_i^f(S) + 1)$ را با $\log n_i^f(S)$ تخمین بزینم می توان Δ_E^c را به صورت زیر تخمین زد:

$$\begin{aligned}\Delta_E^c &= \log \frac{n(S) + 1}{n(S)} + \frac{1}{l} \sum_i \left[\frac{n_i^{f_{d_i}}(S)}{n(S)} \log n_i^{f_{d_i}}(S) - \frac{n_i^{f_{d_i}}(S) + 1}{n(S) + 1} \log(n_i^{f_{d_i}}(S) + 1) \right] \\ &\approx 0 + \frac{1}{l} \sum_i \left[\frac{n_i^{f_{d_i}}(S)}{n(S)} - \frac{n_i^{f_{d_i}}(S) + 1}{n(S) + 1} \right] \log n_i^{f_{d_i}}(S) \\ &= \frac{1}{l} \sum_i \frac{n_i^{f_{d_i}}(S) - n(S)}{n(S)(n(S) + 1)} \log n_i^{f_{d_i}}(S) \\ &= \frac{\sum_i n_i^{f_{d_i}}(S) \log n_i^{f_{d_i}}(S)}{l * n(S)(n(S) + 1)} - \frac{\sum_i \log n_i^{f_{d_i}}(S)}{l(n(S) + 1)}\end{aligned}$$

می دانیم $n_i^{f_{d_i}}(S) \leq n_i(S) = n(S)$ و در نتیجه $\sum_i n_i^{f_{d_i}}(S) \log n_i^{f_{d_i}}(S)$ مقداری در بازه $[n(S)(\log n(S) - \log l), n(S) \log n(S)]$ خواهد بود که در مقایسه با $n(S)^2$ قابل صرف نظر است. پس

$$\Delta_E^c \approx 0 - \frac{\sum_i \log n_i^{f_{d_i}}(S)}{l(n(S) + 1)}$$

از طرفی $P(d|M)$ به علت فرض ما مبنی بر متفاوت بودن مجموعه ویژگی های ممکن در هر حالت $(F_i \cap F_j = \emptyset)$ فقط یک مسیر $t = \langle t_1, t_2, \dots, t_l \rangle$ برای تولید d وجود دارد صورت زیر محاسبه می شود:

$$\begin{aligned}P(d|M) &= P(d|t) = \prod_i p_i^{f_{d_i}}(S) \\ \Rightarrow \log P(d|M) &= \sum \log \frac{n_i^{f_{d_i}}(S)}{n_i(S)} = \sum \log n_i^{f_{d_i}}(S) - \sum \log n_i(S) \\ &= \sum \log n_i^{f_{d_i}}(S) - l \log n(S) = \sum \log n_i^{f_{d_i}}(S) - C\end{aligned}$$

و در نتیجه:

$$\Delta_E^c \approx -C_1 \log P(d|M)$$

تابع f اکیدا نزولی به صورت $f(x) = -C_1 x$ می باشد.

تعمیم حالت دوم

گاهی اوقات به علت ساختار دسته‌ها، نمی‌توان از ساختار مشابه حالت دوم استفاده کرد در این‌گونه موارد احتمال وقوع برخی ویژگی‌ها تنها وابسته به مکان آن ویژگی نیست بلکه می‌تواند بسته به ویژگی‌های دیگر تغییر کند. به عنوان مثال اگر بخواهیم در مورد عضویت بردارهای $l = 5$ عنصری در یک دسته که به صورت زیر تعریف شده است فرآیندی مانند حالت دوم ارائه دهیم احتمالاً به مشکل برخوردیم خورد:

$$l \in C \leftrightarrow (l_1 = f_1) \vee (l_2 = f_2)$$

احتمال رخ دادن ویژگی f_2 در مکان دوم یک بردار در C وابسته به مقدار l_1 است بدین معنا که اگر $f_1 \neq l_1$ باشد آن‌گاه باید داشته باشیم $l_2 = f_2$ (ویژگی f_2 با احتمال ۱ رخ خواهد داد) ولی در صورت $l_1 = f_1$ نمی‌توان در مورد احتمال وقوع ویژگی f_2 در مکان دوم ادعایی کرد. برای رفع این مشکل می‌توان از ساختاری مشابه حالت دوم ولی با کمی تغییرات استفاده کنیم، به این صورت که فرآیند دارای دوره‌هایی با طول l باشد و فقط یک راس مانند (S_0) موجود باشد که در همه دوره‌ها مشترک است. در صورت انتقال از هر یال (عبور از هر یال) کاراکتر (ویژگی منسوب به آن یال) در رشته ظاهر می‌شود. در حقیقت راس S_0 به عنوان نقطه شروع هر بردار در نظر گرفته شود و هر مسیر شامل l یال با شروع از هر حالتی باید از S_0 عبور می‌نماید. برای تعیین احتمال انتقال بین حالات بر اساس مجموعه آموزشی نیز به صورت زیر عمل می‌کنیم:

- از حالت S_0 با توجه به کارکتر مربوط به هر یال به حالت مناسب بعدی برو.

- به علت شروع از حالت S_0 و وجود دوره‌هایی به طول l در انتها به راس S_0 بازگردیم

اگر $n_{i,j}$ را برابر با تعداد دفعاتی که از حالت S_i به حالت S_j برای تولید کل مجموعه آموزشی منتقل شویم و در صورتی که فرآیند مارکف گسسته‌ای با خصوصیات بالا در نظر بگیریم خواهیم داشت:

$$\forall i : \sum_j n_{i,j} = \sum_j n_{j,i}$$

در این حالت خاص خواهیم داشت:

$$\pi_i = \frac{n_i}{\sum_j n_j}$$

همچنین گوییم دنباله $V = \langle v_1, \dots, v_l \rangle$ در مکان i -ام رشته نامتناهی T تولید شده توسط یک فرآیند ظاهر شده است اگر بردار $\langle T_i, T_{i+1}, \dots, T_{i+l-1} \rangle$ برابر با یکی از l بردار $\langle v_1, v_2, \dots, v_l \rangle, \langle v_2, v_3, \dots, v_l, v_1 \rangle, \dots, \langle v_l, v_1, \dots, v_{l-1} \rangle$ باشد. بنابراین تعریف فوق احتمال تولید شدن دنباله V توسط یک فرآیند مارکف را می‌توان به صورت زیر محاسبه کرد (با فرض این که ساختار فرآیند به گونه‌ای باشد که هر بردار با شروع از هر حالت فقط به یک طریق قابل تولید باشد):

$$P(V|M) = \sum_{t \in T} \pi_{t_0} \prod_i p_{t_i, t_{i+1}}^{v_i}$$

به علت فرض بالا مبنی بر یکتا بودن مسیر تولید یک رشته در صورت شروع از یک حالت و همچنین فرض قبلی ما در حالت دوم مبنی بر متفاوت بودن مجموعه ویژگی‌های قابل قبول در هر مکان خواهیم داشت:

۱. اگر V توسط M قابل تولید باشد، الزاما یک مسیر با شروع از S_0 نیز وجود دارد که V را تولید کند (با توجه به تعریف). این مسیر را با $T(V, M)$ نشان می‌دهیم.
۲. سایر اعضای T را می‌توان با جابه‌جا کردن بردار T_M^V به صورت دایره‌ای بدست آورد.

$$\begin{aligned} P(V|M) &= \sum_i \pi_{T(V,M)_i} \prod_j p_{T(V,M)_{(i+j) \bmod l}, T(V,M)_{(i+j+1) \bmod l}}^{V_j} \\ &= \prod_j p_{T(V,M)_j, T(V,M)_{j+1}}^{V_j} \sum_i \pi_i \\ &= \prod_j p_{T(V,M)_j, T(V,M)_{j+1}}^{V_j} \end{aligned}$$

با در دست بودن Π می‌توان به صورتی مشابه حالت دوم ثابت کرد که $P(V|M)$ با $\Delta_E(S, V)$ نسبت معکوس دارد.

در حقیقت این نوع فرآیندها برای هماهنگ سازی کار با روش‌های تعمیم Bayes که در بخش ۴-۴ معرفی شده بود می‌باشد. در فصل ۶ به یک نمونه از این فرآیندها در کاربرد عملی اشاره می‌شود.

حالت‌های سوم و چهارم

در این حالت، به علت تغییر احتمال رخ دادن ویژگی‌ها بر حسب ویژگی قبلی نمی‌توان از روش Naive Bayes انتظار عمل کرد خوبی داشت (به علت فرض Naive Bayes مبنی بر مستقل

بودن ویژگی‌ها) اما می‌توان حدس زد که الگوریتم ارائه شده در این پایان‌نامه، بتواند در این موارد نیز از کارآیی خوبی برخوردار باشد. علت این حدس این است که این روش بر خلاف Naive Bayes می‌تواند با مدل‌ها و فرآیندهای پیچیده‌تری منطبق شود. در حقیقت می‌توان گفت $\Delta_E(S, d)$ یک تخمین برای $P(d|M)$ است.

۷-۵ کارآیی

No Free Lunch Theorem یک قضیه در زمینه بهینه‌سازی ترکیبیاتی^{۷۱} است که توسط دو فیزیک‌دان با نام‌های David H. Wolpert و Willian G. Macready ثابت شده است [۱۵]. در این قضیه ادعا شده است:

کارآیی همه الگوریتم‌هایی که یک نقطه کمینه و یا بیشینه^{۷۲} از یک تابع را جستجو می‌کنند هنگامی که روی همه توابع در نظر گرفته شوند با هم برابر است [۳۲]. در مساله دسته‌بندی داده‌ها، در حقیقت به دنبال یافتن یک نقطه کمینه برای تابع Err تعریف شده در بخش ۱-۳ هستیم و در نتیجه NFLT برای مساله دسته‌بندی قابل اعمال است، و در نتیجه هر الگوریتم دلخواهی، از جمله الگوریتمی که همیشه دسته ۱ را به عنوان برنده اعلام می‌کند و یا الگوریتمی که به ازای هر ورودی به صورت تصادفی یک دسته را اعلام می‌کند در روی کلیه توابع ممکن از کارآیی یکسانی برخوردار هستند. از NFLT می‌توان نتیجه گرفت که هیچ دو الگوریتم دسته‌بندی قابل مقایسه با هم نیستند یعنی نمی‌توان گفت که الگوریتم A در حالت کلی بهتر از الگوریتم B است. این بدان معناست که در صورت وجود یک مساله به طوری که A بهتر از B کار کند می‌توان مساله‌ای (شاید به صورت مصنوعی) ایجاد کرد که در آن B بهتر از A کار کند. البته توجه به این نکته ضروری است که مسائل عملی و کاربردی قسمت بسیار کوچکی از کل مسائل ممکن را تشکیل می‌دهند و بنابراین معمولاً برای مقایسه دو الگوریتم از چند کاربرد عملی استفاده می‌شود.

برای نشان دادن کارآیی روش مبتنی بر تئوری اطلاعات (در حالت‌های اول و دوم)، با استفاده از دو بانک اطلاعاتی استاندارد موجود در وب به مقایسه نتایج Naive Bayes و الگوریتم فوق پرداختیم. بانک اول شامل مجموعه‌ای از نامه‌های الکترونیکی بیست گروه خبری^{۷۳} که در آدرس <http://people.csail.mit.edu/jrennie/20Newsgroups/> موجود است [۲۱] و به عنوان یک مجموعه آزمون استاندارد در زمینه دسته‌بندی متون^{۷۴} شناخته شده است بود.

⁷¹Combinatorics Optimization

⁷²Extremum

⁷³20 news groups

⁷⁴Text Classification

بانک دوم مورد آزمون، مجموعه‌ای از داده‌های مصنوعی بود که با نام MONK معروف هستند و به عنوان یک مجموعه آزمون استاندارد در زمینه یادگیری ماشین شناخته شده است [۲۸]. در فصل‌های ۶ و ۷ به توضیح صورت مساله‌ها و نتایج بدست آمده روی هر یک از بانک‌ها پرداخته‌ایم.

۶ یک مثال

۱-۶ معرفی مساله

مساله MONK یک مساله دسته‌بندی داده‌هاست که برای دسته‌بندی ربات‌های مصنوعی بکار می‌رود. در این مساله ربات‌ها بر اساس شش ویژگی متفاوت خود توصیف می‌شوند [۲۸]:

$x_1 : \text{Head_Shape} \in \{\text{Round}, \text{Square}, \text{Octagon}\}$

$x_2 : \text{Body_Shape} \in \{\text{Round}, \text{Square}, \text{Octagon}\}$

$x_3 : \text{Is_Smiling} \in \{\text{Yes}, \text{No}\}$

$x_4 : \text{Holding} \in \{\text{Sword}, \text{Balloon}, \text{Flag}\}$

$x_5 : \text{Jacket_Color} \in \{\text{Red}, \text{Yellow}, \text{Green}, \text{Blue}\}$

$x_6 : \text{Has_Tie} \in \{\text{Yes}, \text{No}\}$

این مساله شامل سه زیر مساله است که به صورت مستقل از هم و با یک عبارت منطقی مطرح می‌شوند. هر ربات به دسته توصیف شده توسط عبارت منطقی هر زیر مساله، متعلق است و یا نیست. در مساله به جای بیان عبارت منطقی که بر اساس آن دسته‌بندی صورت می‌گیرد شماری از ۴۳۲ ربات ممکن و دسته مربوط به هر یک را به ما داده‌اند و از ما خواسته شده تا تمام ۴۳۲ ربات را در یکی از دو دسته مناسب قرار دهیم. سه زیر مساله به صورت زیر هستند.

مساله ۱

در این مساله اعضای دسته + ربات‌هایی هستند که در عبارت زیر صدق می‌کنند:

$$(x_1 = x_2) \vee (x_5 = \text{Red})$$

از ۴۳۲ ربات ممکن، ۱۲۴ تای آنها به صورت تصادفی انتخاب شده‌اند و در مجموعه آموزش قرار گرفته‌اند. در این مساله هیچ‌گونه آشوب و یا دسته‌بندی اشتباه در مجموعه آموزش وجود ندارد.

مساله ۲

در این مساله اعضای دسته + ربات‌هایی هستند که دقیقاً دو تا از شش ویژگی (x_1, \dots, x_6) آن مقدار نخست خود را اختیار کنند. از ۴۳۲ ربات ممکن، این بار ۱۶۹ تای آنها به صورت

تصادفی انتخاب شده‌اند و در مجموعه آموزش قرار گرفته‌اند. باز هم در این مساله هیچ‌گونه آشوب و یا دسته‌بندی اشتباه در مجموعه آموزش وجود ندارد.

مساله ۳

در این مساله اعضای دسته +، ربات‌هایی هستند که در عبارت زیر صدق کنند:

$$(x_5 = Green \wedge x_6 = Sword) \vee (x_5 \neq Blue \wedge x_6 \neq Octagon)$$

از ۴۳۲ ربات ممکن، ۱۲۲ تای آنها به صورت تصادفی انتخاب شده‌اند و در مجموعه آموزش قرار گرفته‌اند. این بار ۵٪ از اعضای مجموعه آموزش به اشتباه دسته‌بندی شده‌اند.

مساله ۱، یک مساله استاندارد به صورت DNF^{۷۵} است و احتمالاً همه الگوریتم‌های یادگیری نمادین مانند درخت‌های تصمیم باید به خوبی از پس آن برآیند. برخلاف مساله ۱، مساله ۲ یک مساله کنترل‌کردن زوجیت است اما توصیف آن به صورت یک فرمول بسته DNF و یا CNF^{۷۶} کمی دشوار است. مساله ۳، را به راحتی می‌توان به صورت رابطه DNF نوشت ولی این بار الگوریتم باید قابلیت دسته‌بندی براساس آشوب را داشته باشد. [۲۸]

۲-۶ تبدیل مساله به مدل مارکف مربوطه

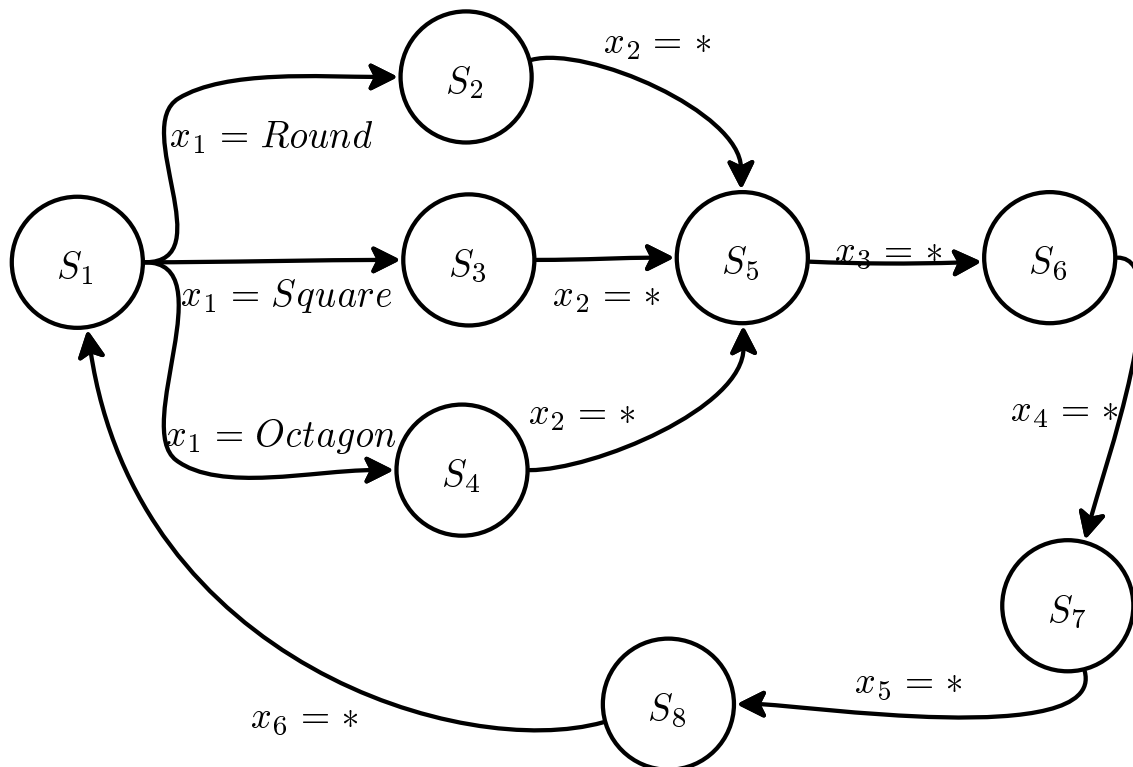
برای حل هر یک از زیر مساله‌های مطرح شده در بالا از یک فرآیند نوع ۲ با جزئیات متفاوت (بسته به رابطه دسته‌بندی) استفاده شد.

مساله ۱

برای حل این مساله از فرآیند نشان داده شده در شکل ۱-۶ استفاده شد: علت پیشنهاد داده شدن فرآیند ۱-۶ این بوده است که اعضای دسته +، دارای یک رابطه بین ویژگی‌های ۱ و ۲ خود (x_1 و x_2) و ویژگی ۵ (x_5) هستند و در فرآیند پیشنهادی اگر داده‌ای در این دسته نباشد پس قطعاً دارای x_1 و x_2 متفاوتی خواهد بود و این موضوع باعث می‌شود از یکی از راس‌های S_1 ، S_2 و یا S_4 یال جدیدی خارج شود که می‌توان انتظار افزایش

⁷⁵Disjunctive Normal Form

⁷⁶Conjunctive Normal Form

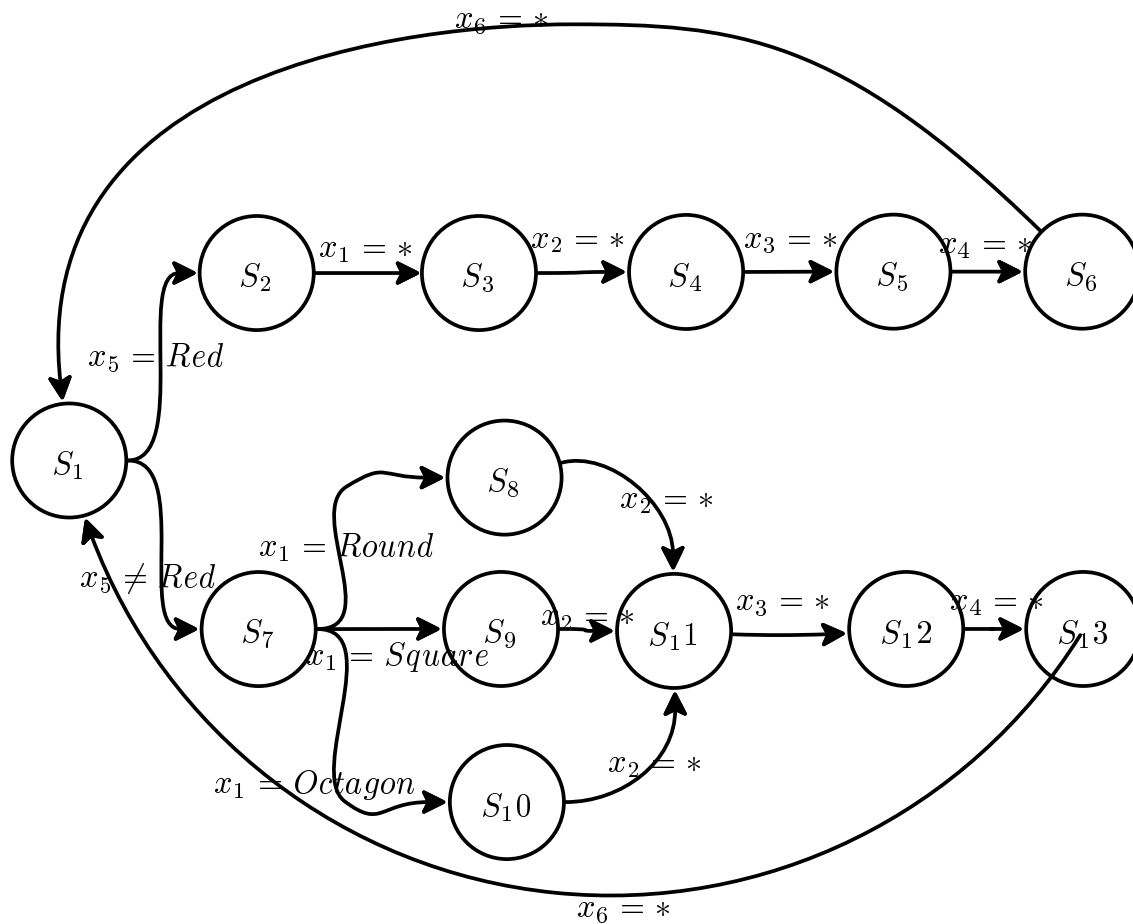


شکل ۶-۱: یک فرآیند پیشنهادی برای مساله ۱

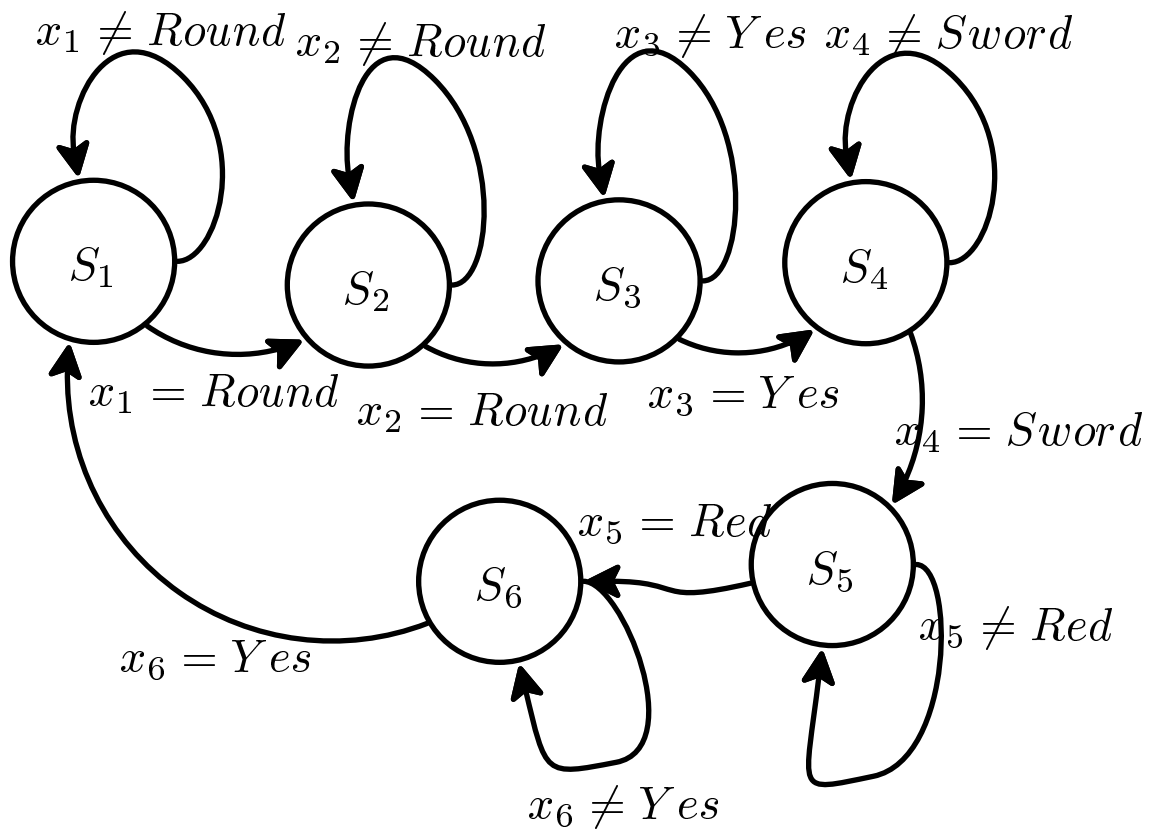
آشفته‌گی در بین احتمالات باشد اما به دلیل وجود داده‌هایی در این دسته که دارای x_1 و x_2 متفاوت هستند این فرآیند ممکن است کارایی خوبی نداشته باشد. فرآیند ۶-۲ فرآیند دقیق‌تری برای توصیف این دسته‌ها می‌باشد.

مساله ۲

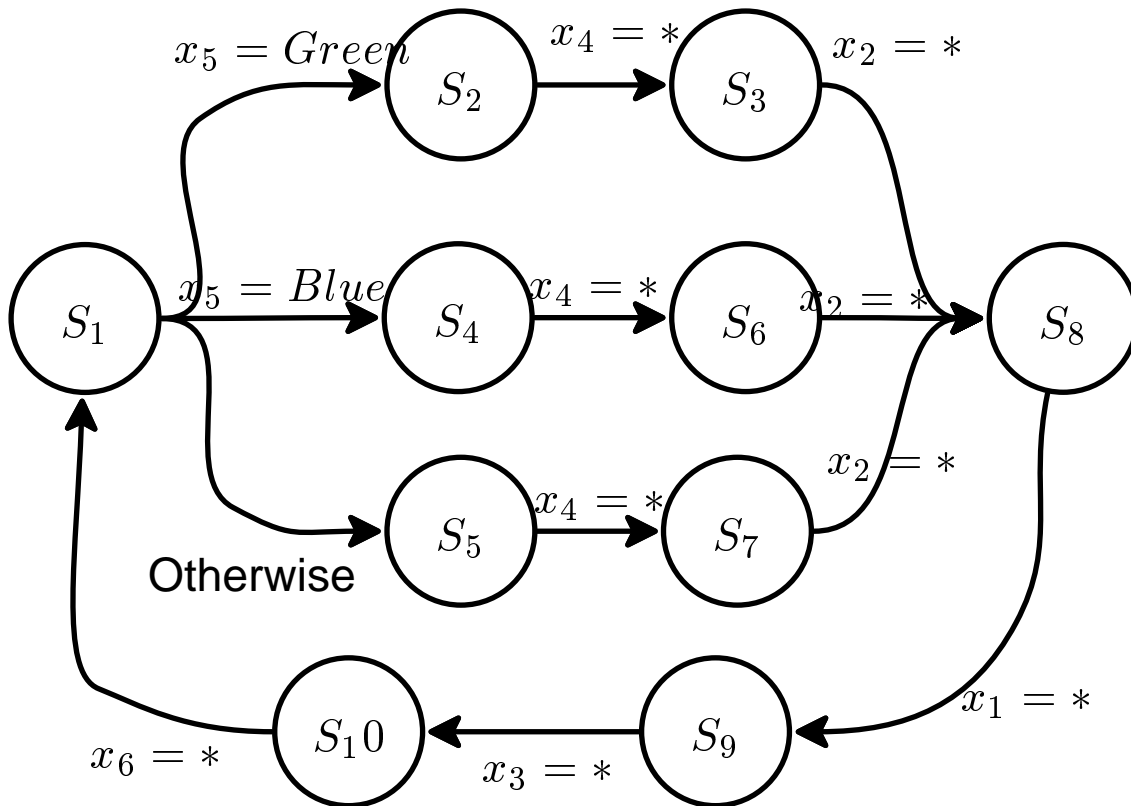
برای حل این مساله از فرآیند نشان داده شده در شکل ۶-۳ استفاده شد: علت پیشنهاد داده شدن این فرآیند این بوده است که اعضای دسته +، دارای دقیقا دو ویژگی با مقدار نخست خود هستند و در نتیجه تعداد دفعاتی که فرآیند برای تولید اعضای مجموعه آموزش وارد حالات S_4 ، S_5 و S_6 شده برابر با صفر است و در نتیجه با ورود داده‌ای که سه و یا بیشتر ویژگی آن دارای ویژگی نخست باشد Entropy فرآیند را بالا می‌برد همچنین اگر داده‌ای صفر و یا یک ویژگی با مقدار نخست داشته باشد تعداد ورود در حالت S_1 و یا S_2 (n_1 و n_2) تغییر نامتناسبی نسبت به داده‌های مجموعه آموزش خواهد کرد و در نتیجه Entropy تغییرات مثبتی خواهد داشت. برای این مساله فرآیندهای پیچیده‌تری را نیز می‌توان ارائه داد ولی هیچ تضمینی برای بهتر



شکل ۶-۲: یک فرآیند دیگر برای مساله ۱



شکل ۶-۳: یک فرآیند پیشنهادی برای مساله ۲



شکل ۶-۴: یک فرآیند پیشنهادی برای مساله ۳

کارکردن این فرآیندها وجود ندارد.

مساله ۳

برای حل این مساله از فرآیند نشان داده شده در شکل ۶-۴ استفاده شد: علت پیشنهاد داده شدن این فرآیند این بوده است که اعضای دسته +، دارای نوعی محدودیت روی ویژگی های ۲، ۴ و ۵ خود (x_2 ، x_4 و x_5) خود است و در فرآیند پیشنهادی اگر داده ای در این دسته نباشد احتمالاً آشفتگی زیادی در بین احتمالات ایجاد می کند و در نتیجه Entropy را افزایش می دهد.

۶-۳ بحث در مورد تخمین ها

همان گونه که در فصل ۵ شرح داده شد در صورت وجود شرایط خاص، می توان از ITDC انتظار کارایی خوبی داشت (یعنی به خوبی Naive Bayes عمل کند). به همین دلیل در این

بخش به بررسی برقرار بودن آن شرایط در هر یک از سه مساله بالا می‌پردازیم. فرض اولیه در فصل ۵، بسیار بزرگ‌تر بودن اندازه مجموعه آموزش نسبت به داده آزمون بوده است ($n(s) \gg n(d)$). همان‌طور که شرح داده شد $n(s) = 124 * 6$ بوده در حالی که $n(d) = 1 * 6$ و در نتیجه می‌توان با دقت خوبی فرض کرد که شرط فوق برقرار است.

مساله ۱

در این گونه از فرآیندهای مارکف، شرط اشاره شده در بخش ۶-۵، باید شرط زیر برقرار باشد: $\log n_c(S) = \log n_c(d)$: این شرط نیز با توجه به داده‌های موجود در مجموعه آموزشی و با توجه به این واقعیت که $n_c(d) \leq 1$ قابل ارضا می‌باشد.

مساله ۲

برای فرآیند مارکف از نوع دوم که در فصل ۵ بدان‌ها اشاره شد، باید دو شرط زیر برقرار باشد:

۱. $\log n_c(S) = \log n_c(d)$: این شرط نیز با توجه به داده‌های موجود در مجموعه آموزشی و با توجه به این واقعیت که $n_c(d) \leq 1$ قابل ارضا می‌باشد.

۲. $\log \frac{(n(S)+1)}{n(S)} \approx 0$: با توجه به $n(S) = 169$ خواهیم داشت:

$$\log \frac{170}{169} = 0.0085 \approx 0$$

مساله ۳

برای فرآیند مارکف مورد استفاده در این مساله نیز مانند مساله دوم، باید دو شرط زیر برقرار باشد:

۱. $\log n_c(S) = \log n_c(d)$: این شرط نیز با توجه به داده‌های موجود در مجموعه آموزشی و با توجه به این واقعیت که $n_c(d) \leq 1$ قابل ارضا می‌باشد.

۲. $\log \frac{(n(S)+1)}{n(S)} \approx 0$: با توجه به $n(S) = 122$ خواهیم داشت:

$$\log \frac{123}{122} = 0.0118$$

که اگر ۵٪ خطای موجود در مجموعه آموزشی را در نظر بگیریم این عدد بزرگ‌تر نیز می‌شود. اما با این حال ITDC در این مساله نیز جواب قابل قبولی داده است.

۴-۶ نتایج حاصل

ستون اول جدول زیر نشان دهنده میزان دقت الگوریتم دسته‌بندی براساس تئوری اطلاعات (ITDC) است که در این پایان‌نامه معرفی شد و در ستون‌های بعدی دقت چند الگوریتم دیگر آمده است. این اطلاعات از [۲۸] استخراج شده است. باید توجه داشت که احتمالاً فرآیندهای پیشنهادی بهترین فرآیندهای ممکن نبودند. یکی از کارهای پیشنهادی در ادامه این پایان‌نامه ایجاد خودکار فرآیند با استفاده از مجموعه آموزش است که در فصل ۸ بدان اشاره شده است.

Problem	ITDC	Decision Trees	Naive Bayes
M1	100%	100%	75%
M2	94%	81%	61%
M3	100%	100%	97%

جدول بالا نشان می‌دهد که ITDC از Naive Bayes حداقل در برخی موارد از کارایی بهتری برخوردار است. دلیل این برتری در این مساله، ذاتا تصادفی نبودن و همچنین صدق نکردن مساله در شرط استقلال ویژگی‌ها است که از کارایی Naive Bayes کاسته است.

۷ یک مثال

۱-۷ معرفی مساله

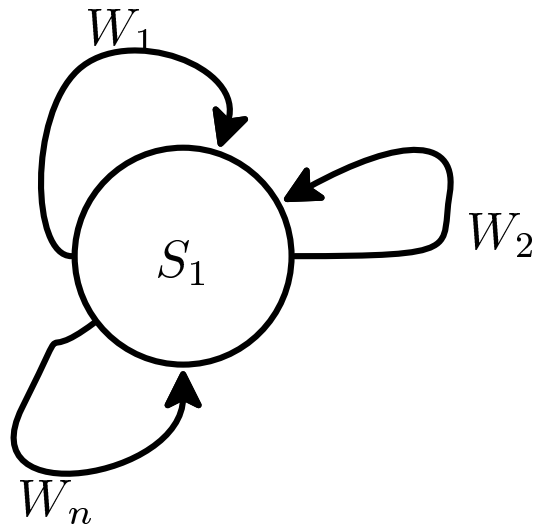
بانک 20 NewsGroup از حدود ۲۰۰۰۰۰ متن e-mail که از بین ۲۰ گروه خبری متفاوت جمع آوری شده‌اند تشکیل شده‌است. این بانک توسط Ken Lang به عنوان یک مجموعه آزمون برای مقاله خود با نام Newsweeder: Learning to filter netnews تهیه شده است. این مجموعه امروزه به عنوان یک مجموعه استاندارد در آزمایشات و اندازه‌گیری کیفیت الگوریتم‌ها استفاده می‌شود. کاربرد اصلی این مجموعه در دسته‌بندی متن‌ها^{۷۷} و کاربردهای آن در هوش مصنوعی است. ۲۰ گروه خبری مورد استفاده در این مجموعه، موضوعات متفاوت و متنوعی را در برمی‌گیرد. بعضی از موضوعات به صورت زیادی مشابه هستند مانند comp.sys.ibm.pc.hardware و comp.sys.mac.hardware، و برخی دیگر کاملاً متفاوت به نظر می‌رسند مانند misc.forsale و soc.religion.christian. در جدول زیر این ۲۰ گروه بر حسب ارتباط در موضوعات مورد بحث در آن‌ها آورده شده‌اند [۲۲].

comp.graphics	rec.autos	sci.cryp
comp.os.ms-windows.misc	rec.motorcycles	sci.electronics
comp.sys.ibm.pc.hardware	rec.sport.baseball	sci.med
comp.sys.mac.hardware	rec.sport.hockey	sci.space
comp.windows.x		
misc.forsale	talk.politics.guns	talk.religion.misc
	talk.politics.mideast	alt.atheism
	talk.politics.misc	soc.religion.christian

۲-۷ تبدیل مساله به مدل مارکف مربوطه

به علت گسترده بودن دامنه واژگان موجود در مجموعه (حدود واژه) از یک فرآیند بدون حافظه مانند شکل ۱-۷ برای مدل کردن فرآیند تولید این داده‌ها استفاده شد. در پیاده‌سازی هر واژه به عنوان یک ویژگی در نظر گرفته شد و برای راحت‌تر کردن پیاده‌سازی به هر واژه یک عدد طبیعی به عنوان کد آن واژه اختصاص داده شد. برای محدود کردن تعداد

⁷⁷Text Classification



شکل ۷-۱: فرآیند پیشنهادی برای مساله دسته‌بندی خبرها

ویژگی‌ها، واژه‌هایی که دارای طولی بیشتر از یک مقدار از پیش تعیین شده بود (در پیاده‌سازی ما ۱۵) با صفر (یک ویژگی که برای این منظور در نظر گرفته شد) نمایش داده شد.

۷-۳ نتایج حاصل

برای بدست آوردن نتیجه بهتر، از روش چند سطحی استفاده شد. به این ترتیب که عمل تشخیص در دو مرحله انجام می‌گیرد. در مرحله اول، داده به یک گروه از دسته‌ها بین چندین گروه تخصیص داده می‌شود و سپس در مرحله بعد در مورد دسته نهایی بین اعضای آن گروه تصمیم‌گیری می‌شود. رده بندی گروه‌ها به صورت نمایش داده شده در جدول زیر بوده است (مانند [۳۳]):

Level 1	Level 2	Level3
alt.atheism misc.forsale soc.religion.christian		
comp.	graphics os.ms-windows.misc windows.x	
	sys.	ibm.pc.hardware mac.hardware
rec.	autos motor-cycle	
	sport.	baseball hockey
sci.	crypt electronics med space	
talk	religion.misc	
	politics.	guns mideast misc

متوسط دقت دسته‌بندی در جدول پایین با چند روش دیگر از [۲۱] مقایسه شده است:

Method	Accuracy
Information Theoretic	88.0%
H-SVN	89.1%
MNB	84.8%
TWCNB	86.1%
SVM	86.2%

در این جدول میزان دقت در دسته‌بندی روی هر یک از دسته‌ها را نشان می‌دهد:

Group Name	Accuracy
alt.theism	74%
comp.graphics	79%
comp.os.ms-windows.misc	70%
comp.sys.ibm.pc.hardware	67%
comp.sys.mac.hardware	80%
comp.windows.x	80%
misc.forsale	85%
rec.autos	63%
rec.motorcycles	89%
rec.sport.baseball	92%
rec.sport.hockey	95%
sci.crypt	83%
sci.electronics	80%
sci.med	91%
sci.space	96%
soc.religion.christian	70%
talk.politics.guns	93%
talk.politics.mideast	90%
talk.politics.misc	84%
talk.religion.misc	87%

۸ کارهای آتی

به علت محدود بودن تعریف این پایان نامه به دسته بندی داده های متنی، تمرکز اصلی روی منابع بدون حافظه و Naive Bayes بود و الگوریتم بدست آمده در طول تحقیق فقط در همین زمینه ها آزمایش شد. تعمیم ریاضی این روش در حالاتی که منبع دارای حافظه است و مقایسه آن با روش های Bayes و سایر روش ها و همچنین انجام یک سری مقایسه عملی برای مشاهده کارایی الگوریتم ارائه شده در حالات حافظه دار را می توان به عنوان اصلی ترین کار آینده معرفی کرد.

یکی از مشکلاتی که روش پیشنهادی (ITDC) با آن روبرو است نیاز به طراحی یک فرآیند دقیق و کارا برای هر مساله (هر دسته) است. همان طور که در فصل ۶ دیدیم کارایی ITDC وابسته به فرآیند تغییر می کند. یکی دیگر از کارهای مرتبط با این پایان نامه را می توان ایجاد یک ساختار در مرحله آموزش برای پیشنهاد دادن فرآیندهایی بر اساس مجموعه آموزشی تعریف کرد. بدین صورت که در ابتدا فرآیند یک فرآیند ساده باشد و با افزوده شدن داده های جدید آموزشی به یک دسته فرآیند پیچیده تر و در عین حال اختصاصی تر برای هر دسته گردد ولی باید توجه داشت که فرآیند بسیار وابسته به مجموعه آموزشی نگردد.

References

- [1] J. ACZEL, B. FORTE, AND C. NG, *Why the Shannon and Hartley Entropies Are 'Natural'*, Advances in Applied Probability, 6 (1974), pp. 131–146.
- [2] P. AHAMMAD, K. DASKALAKIS, O. ETESAMI, AND A. FROME, *Claude Shannon and A Mathematical Theory of Communication*.
- [3] P. BERGNER, *Stochastic processes and applications in biology and medicine. Part 1 Theory. Part 2 Models*, Bulletin of Mathematical Biology, 36 (1974), pp. 607–610.
- [4] S. BOUCHERON, O. BOUSQUET, AND G. LUGOSI, *Theory of classification: a survey of some recent advances*, ESAIM P AND S, 9 (2005), p. 323.
- [5] T. COVER AND J. THOMAS, *Elements of information theory*, Wiley New York, 1991.
- [6] L. DEVROYE, L. GYORFI, AND G. LUGOSI, *A Probabilistic Theory of Pattern Recognition*, Springer, 1996.
- [7] M. D.J.C, *Information Theory, Inference and Learning Algorithm*, Cambridge University Press, 2005.
- [8] B. FORTE AND C. SASTRI, *Is something missing in the Boltzmann entropy?*, Journal of Mathematical Physics, 16 (2006), pp. 1453–1456.
- [9] K. FUKUNAGA AND D. HUMMELS, *Bayes error estimation using Parzen and k-NN procedures*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 9 (1987), pp. 634–643.
- [10] D. GAO, M. MADDEN, D. CHAMBERS, AND G. LYONS, *Bayesian ANN classifier for ECG arrhythmia diagnostic system: a comparison study*, Neural Networks, 2005. IJCNN'05. Proceedings. 2005 IEEE International Joint Conference on, 4 (2005).
- [11] S. GOLDSTEIN, *Remarks on the Global Markov Property*, Commun. math. Phys, 74 (1980), pp. 223–234.
- [12] R. GRAY, *Entropy and information theory*, Springer-Verlag New York, Inc. New York, NY, USA, 1990.

- [13] S. HAYKIN, *Neural Networks: A Comprehensive Foundation*, Prentice Hall PTR Upper Saddle River, NJ, USA, 1998.
- [14] C. HILLMAN, *An Entropy Primer*, tech. rep., Mimeo retrieved from URL: <http://www.math.washington.edu/hillman>, 1996.
- [15] Y. HO AND D. PEPYNE, *Simple Explanation of the No-Free-Lunch Theorem and Its Implications*, Journal of Optimization Theory and Applications, 115 (2002), pp. 549–570.
- [16] R. MASON, D. LIND, AND W. MARCHAL, *Statistics: an introduction*, Harcourt Brace Jovanovich, 1983.
- [17] D. MICHIE, D. SPIEGELHALTER, AND C. TAYLOR, *Machine learning of rules and trees*, Machine Learning, Neural and Statistical Classification, (1994).
- [18] T. MITCHELL, *Machine Learning*, McGraw-Hill Higher Education, 1997.
- [19] M. NELSON AND J. GAILLY, *The Data Compression Book*, MIS: Press New York, NY, USA, 1995.
- [20] J. QUINLAN, *Bagging, Boosting, and C4. 5*.
- [21] J. RENNIE, *Improving Multi-class Text Classification with Naive Bayes*, (2001).
- [22] J. RENNIE, L. SHIH, J. TEEVAN, AND D. KARGER, *Tackling the Poor Assumptions of Naive Bayes Text Classifiers*, Proceedings of the Twentieth International Conference on Machine Learning, 41 (2003), p. 18.
- [23] O. B. S. BOUCHERON AND G. LUGOSI, *Theory of classification: A survey of some recent advances*, (2001).
- [24] S. SAFAVIAN AND D. LANDGREBE, *A survey of decision tree classifier methodology*, IEEE Transactions on Systems, Man, and Cybernetics, 21 (1991), pp. 660–673.
- [25] P. SAHOO, S. SOLTANI, A. WONG, AND Y. CHEN, *A survey of thresholding techniques*, Computer Vision, Graphics, and Image Processing, 41 (1988), pp. 233–260.
- [26] C. SHANNON, *A mathematical theory of communication*, The Bell Syst Tech J, 27 (1948), pp. 379–423.

- [27] F. THOLLARD, P. DUPONT, AND C. DE LA HIGUERA, *Probabilistic DFA inference using Kullback-Leibler divergence and minimality*, Proc. Int. Conf. on Machine Learning, (2000), pp. 975–982.
- [28] S. THRUN ET AL., *The MONK’s Problems: A Performance Comparison of Different Learning Algorithms*, School of Computer Science, Carnegie Mellon University, 1991.
- [29] S. VAN DE GEER, *Empirical Processes in M-Estimation*, Cambridge University Press, 2000.
- [30] A. WEHRL, *General properties of entropy*, Reviews of Modern Physics, 50 (1978), pp. 221–260.
- [31] D. WILLIAMS, *Weighing the Odds: A Course in Probability and Statistics*, Cambridge University Press, 2001.
- [32] D. WOLPERT, W. MACREADY, I. CENTER, AND C. SAN JOSE, *No free lunch theorems for optimization*, Evolutionary Computation, IEEE Transactions on, 1 (1997), pp. 67–82.
- [33] Y. YOON, C. LEE, AND G. LEE, *Systematic Construction of Hierarchical Classifier in SVM-based Text Categorization*.
- [34] J. ZURADA, *Introduction to artificial neural systems*, West Publishing Co. St. Paul, MN, USA, 1992.