

مقایسه رگرسیون لوژستیک با برخی از الگوریتم‌های یادگیری ماشین در رده‌بندی داده‌ها

طیبه کرمی^۱، محی‌الدین ایزدی^۲، مهرداد نیاپرست^۳

تاریخ دریافت: ۱۴۰۰/۰۲/۰۴

تاریخ پذیرش: ۱۴۰۰/۰۷/۱۶

چکیده:

یکی از مسائل مهم در علوم مختلف موضوع رده‌بندی است. رگرسیون لوژستیک یکی از روش‌های آماری برای رده‌بندی داده‌ها است که در آن توزیع داده‌ها معلوم فرض می‌شود. امروزه، محققان علاوه بر روش‌های آماری از روش‌های دیگری مانند الگوریتم‌های یادگیری ماشین که در آن نیاز به معلوم بودن توزیع داده‌ها نیست برای رده‌بندی داده‌ها استفاده می‌کنند. در این مقاله، رگرسیون لوژستیک و برخی از الگوریتم‌های یادگیری ماشین شامل CART، جنگل تصادفی، Bagging و تقویت در حوزه‌ی یادگیری با نظارت توضیح داده می‌شود. با استفاده از ۴ مجموعه داده واقعی و همچنین یک مثال شبیه‌سازی شده کارایی رگرسیون لوژستیک با الگوریتم‌های یادشده، بر اساس معیارهای دقت، حساسیت و صحت مورد مقایسه قرار می‌گیرند.

واژه‌های کلیدی: جنگل تصادفی، درخت تصمیم، یادگیری با نظارت، یادگیری گروهی

۱ مقدمه

نظارت مجموعه داده‌ها دارای برجسب هستند. به عبارت دیگر، یک مجموعه داده شامل n مشاهده (نمونه) از ورودی‌ها (ویژگی‌ها، متغیرهای پیشگو) و خروجی (ویژگی هدف، متغیر پاسخ) است که ورودی‌ها با X_1, X_2, \dots, X_p و خروجی با Y نشان داده می‌شود. در یک مجموعه داده اگر متغیر پاسخ، رسته‌ای باشد یادگیری با نظارت به مسئله رده‌بندی می‌پردازد. درخت تصمیم یکی از پرکاربردترین و قدیمی‌ترین الگوریتم‌های یادگیری با نظارت است که ساختاری شبیه فلوچارت دارد و جزء روش‌های ناپارامتری به حساب می‌آید ([۳]). هر درخت تصمیم شامل ریشه، گره‌های درونی، برگ و شاخه‌ها است. یک درخت تصمیم از ریشه به سمت پایین رشد می‌کند و در نهایت به برگ می‌رسد. تعریف اجزاء درخت تصمیم در زیر ارائه می‌شود.

ریشه: بالاترین گره در درخت تصمیم ریشه نام دارد و سایر گره‌ها زیر آن قرار می‌گیرند. ریشه مهم‌ترین جزء در ساختار درخت تصمیم است.

شاخه: در نمودار درخت تصمیم خطوطی که گره‌ها را به یکدیگر وصل می‌کنند شاخه نام دارند.

برگ: گره پایانی در ساختار درخت تصمیم برگ نامیده می‌شود. به عبارت دیگر برگ همان ویژگی هدف یا متغیر پاسخ است.

گره درونی: سایر گره‌ها به جز ریشه و برگ را گره درونی گویند.

یکی از مسائل مهم در علوم مختلف از جمله بیمه، اقتصاد، علوم پزشکی و ... رده‌بندی است. معمولاً در مسائل رده‌بندی، مشاهدات موجود شامل p متغیر پیشگو و یک متغیر پاسخ است که می‌تواند مقادیر ۱، ۲، ...، k را اختیار کند. هدف در مسائل رده‌بندی این است که با استفاده از مشاهدات موجود، رده یک داده‌ی جدید پیش‌بینی شود. رگرسیون لوژستیک یکی از روش‌های آماری برای رده‌بندی داده‌ها است که در آن توزیع داده‌ها معلوم فرض می‌شود. پژوهشگران امروزه علاوه بر روش‌های آماری از روش‌های دیگری که در آن نیاز به معلوم بودن توزیع داده‌ها نیست؛ مانند یادگیری ماشین برای رده‌بندی داده‌ها استفاده می‌کنند. یکی از شاخه‌های وسیع و پرکاربرد در هوش مصنوعی یادگیری ماشین است. پژوهشگران در داده‌کاوی، کنترل روبات‌ها، تشخیص گفتار و بسیاری از مسائل دیگر یادگیری ماشین را به کار می‌گیرند. در یادگیری ماشین، کامپیوتر مجموعه داده را دریافت کرده و داده الگوریتم را مشخص کرده و به آن آموزش می‌دهد. یادگیری ماشین شامل یادگیری با نظارت، یادگیری بدون نظارت، یادگیری نیمه نظارتی و یادگیری تقویتی است. در این مقاله، به مطالعه برخی از الگوریتم‌های یادگیری با نظارت می‌پردازیم. در یادگیری با

^۱دانش‌آموخته گروه آمار، دانشگاه رازی، کرمانشاه، (نویسنده مسئول: tayebeh.karami20@gmail.com)

^۲هیئت علمی گروه آمار، دانشگاه رازی، کرمانشاه

^۳هیئت علمی گروه آمار، دانشگاه رازی، کرمانشاه

۲ رگرسیون لوژیستیک

یک حالت از رگرسیون خطی تعمیم‌یافته، رگرسیون لوژیستیک است که ارتباط بین متغیرهای پیشگو و متغیر پاسخ را بیان می‌کند. در رگرسیون لوژیستیک متغیر پاسخ یک متغیر رسته‌ای با دو یا بیش از دو مقدار است و متغیرهای پیشگو می‌توانند متغیرهای پیوسته و رسته‌ای باشند.

۱۰.۲ رگرسیون لوژیستیک دودویی

متغیر پاسخ در رگرسیون لوژیستیک دودویی یک متغیر رسته‌ای است که شامل دو مقدار ۰ و ۱ است. برای مثال، دو رده بیماری یا سلامت، مرگ یا زندگی و... را می‌توان نام برد. فرض کنید Y متغیر پاسخ و همچنین X_1, X_2, \dots, X_p متغیرهای پیشگو با مشاهدات متناظر x_1, x_2, \dots, x_p باشند. آنگاه در رگرسیون لوژیستیک دودویی،

$$Y \sim \text{Bin}(1, \pi(\beta, x))$$

که در آن $\pi(\beta, x) = P(Y = 1)$ دلالت بر احتمال موفقیت (متعلق بودن به رده ۱) به ازای بردار مشاهدات $x = (x_1, \dots, x_p)$ است. تابع لوجیت به‌عنوان تابع ربط (تبدیل) در رگرسیون لوژیستیک استفاده می‌شود و به‌صورت

$$\begin{aligned} \text{logit}(\pi(\beta, x)) &= \log \frac{\pi(\beta, x)}{1 - \pi(\beta, x)} \\ &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \end{aligned}$$

است. پارامترهای رگرسیون لوژیستیک با استفاده از روش ماکسیمم درستنمایی برآورد می‌شوند. فرض کنید $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ برآورد پارامترهای رگرسیون لوژیستیک باشند. آنگاه مدل برازش شده به‌صورت

$$\text{logit}(\hat{\pi}(\hat{\beta}, x)) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

و یا

$$\hat{\pi}(\hat{\beta}, x) = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p)}$$

است. برای مشاهده جدید اگر $\hat{\pi}(\hat{\beta}, x) > \frac{1}{2}$ باشد، آنگاه x متعلق به رده $y = 1$ است ([۲]).

درخت‌های تصمیم مختلفی از جمله درخت‌های ID3، C4.5، C5، CART⁴، CHAID و QUEST توسط پژوهشگران معرفی شده است ([۱۰]). الگوریتم CART یکی از پرکاربردترین درخت‌های تصمیم است که در بخش‌های بعد به‌صورت کامل توضیح داده می‌شود. از مزایای درخت تصمیم می‌توان به تفسیرپذیری بالا، فهم ساده، مقاومت در برابر نویز و سرعت مناسب اشاره کرد. زمانی که اندازه و عمق درخت زیاد باشد، بیش‌برازش اتفاق می‌افتد که می‌توان از روش هرس کردن برای رفع آن استفاده کرد. در هرس کردن درخت تصمیم، بعضی از شاخه‌ها و گره‌هایی که تأثیر ناچیزی بر رده‌بندی داده‌ها دارند حذف می‌شود. درخت تصمیم معیابی هم به شرح زیر دارد ([۱۰]).

۱- هرس درخت تصمیم به هزینه بالا و زمان بیشتری نیاز دارد.
۲- در مواردی که تعداد داده‌های آموزشی کم است ولی تعداد رده‌های متغیر پاسخ زیاد باشد، احتمال خطا در درخت تصمیم زیاد است.
۳- الگوریتم درخت تصمیم جز الگوریتم‌های یادگیری ناپایدار است؛ زیرا با یک تغییر کوچک در داده‌های آموزشی می‌تواند منجر به ایجاد یک درخت بسیار متفاوت می‌شود.

این معایب سبب می‌شود که درخت تصمیم در بعضی از مسائل رده‌بندی دارای دقت پایین باشد. به این منظور الگوریتم‌های یادگیری گروهی پیشنهاد می‌شود. الگوریتم‌های یادگیری گروهی متشکل از دو یا بیش از دو الگوریتم یادگیری ماشین است. مطالعات نشان می‌دهد که در بعضی مواقع عملکرد الگوریتم‌های یادگیری گروهی بهتر از استفاده از یک الگوریتم است ([۱]). یادگیری گروهی الگوریتم‌های مختلفی دارد که در این مقاله، سه الگوریتم Bagging^۵، جنگل تصادفی و تقویت^۶ مورد مطالعه قرار می‌گیرد. در این مطالعه، با استفاده از ۴ مجموعه داده مختلف، به مقایسه مدل رگرسیون لوژیستیک با الگوریتم‌های CART، جنگل تصادفی، تقویت سازوار^۷ CART و Bagging CART می‌پردازیم. ساختار ارائه مطالب در این مقاله به این صورت زیر است.

در بخش ۲، رگرسیون لوژیستیک ارائه می‌شود. در بخش ۳، الگوریتم CART بیان شده است. الگوریتم‌های یادگیری گروهی در بخش ۴ توضیح داده می‌شود. در بخش ۵، به توضیح معیار ارزیابی عملکرد الگوریتم پرداخته شده است و مقایسه الگوریتم‌ها در بخش ۶ ارائه می‌شود.

⁴Classification And Regression Tree

⁵Bootstrap aggregating

⁶Boosting

⁷Adaptive boosting (Adaboosting)

۲.۲ رگرسیون لوژستیک چند رده‌ای

متغیر پاسخ در رگرسیون لوژستیک چند رده‌ای بیش از دو مقدار اختیار می‌کند. فرض کنید Y یک متغیر رسته‌ای با c مقدار ممکن (c رده) باشد و همچنین X_1, \dots, X_p متغیرهای پیشگو باشند. آنگاه مدل رگرسیون لوژستیک بر پایه رابطه

$$\log\left(\frac{\pi_j(\beta, \mathbf{x})}{\pi_c(\beta, \mathbf{x})}\right) = \alpha_j + \beta_{j1}x_1 + \beta_{j2}x_2 + \dots + \beta_{jp}x_p; \quad j = 1, 2, \dots, c$$

است که در آن $\pi_j(\beta, \mathbf{x}) = P(Y = j)$ و $\sum_{j=1}^c \pi_j(\beta, \mathbf{x}) = 1$ به‌طور کلی برای پیش‌بینی رده داده‌ی جدید در رگرسیون لوژستیک چند رده‌ای از شاخصی به نام نسبت بخت‌ها ($^{\wedge}OR$)؛ یعنی

$$OR_{ij} = \frac{\frac{\pi_i(\beta, \mathbf{x})}{1 - \pi_i(\beta, \mathbf{x})}}{\frac{\pi_j(\beta, \mathbf{x})}{1 - \pi_j(\beta, \mathbf{x})}}$$

استفاده می‌شود.

(۱) اگر $OR_{ij} > 1$ باشد، آنگاه مشاهده $\mathbf{x} = (x_1, \dots, x_p)$ متعلق به رده i است.

(۲) اگر $OR_{ij} < 1$ باشد، آنگاه مشاهده $\mathbf{x} = (x_1, \dots, x_p)$ متعلق به رده j است.

(۳) اگر $OR_{ij} = 1$ باشد، آنگاه نمی‌توان پیش‌بینی کرد که مشاهده $\mathbf{x} = (x_1, \dots, x_p)$ به کدام رده تعلق دارد ([۲]).

۳ الگوریتم CART

الگوریتم CART توسط بریمن^۹ [۴] معرفی شده است. این درخت برای مسائل رده‌بندی و رگرسیون به‌کار برده می‌شود. در این مقاله، از الگوریتم CART در مسائل رده‌بندی استفاده می‌شود. خصوصیات الگوریتم CART به‌صورت زیر است:

(۱) متغیرهای پیشگو در یک مجموعه داده برای این درخت می‌توانند به‌صورت متغیرهای رسته‌ای و کمی باشند.

(۲) فقط دو شاخه از هر گره درونی این درخت خارج می‌شود.

(۳) برای هرس کردن الگوریتم CART از معیار پیچیدگی هزینه هرس^{۱۰}، برای هرس کردن استفاده می‌شود.

(۴) این درخت توانایی انجام رده‌بندی با داده گمشده در یک مجموعه داده را دارد.

(۵) در این درخت برای انتخاب ریشه و گره‌ها از شاخص جینی و شاخص دوپاره‌سازی^{۱۱} استفاده می‌شود. در یک مجموعه

داده هنگامی متغیر پاسخ دارای دو رده است از شاخص جینی^{۱۲} و هنگامی که متغیر پاسخ بیش از دو رده داشته باشد از شاخص دوپاره‌سازی استفاده می‌شود ([۱۰]). فرض کنید S یک مجموعه داده با n مشاهده است. شاخص جینی متغیر پاسخ Y به‌صورت

$$\text{Gini}(Y, S) = 1 - \sum_{c_j \in \text{dom}(Y)} \left(\frac{|\sigma_{Y=c_j}(S)|}{|S|} \right)^2 \quad (۱)$$

تعریف می‌شود که در آن

$\sigma_A(S)$: زیرمجموعه‌ای از S است که شرط A برای اعضای آن برقرار است.

$|S|$: تعداد اعضای مجموعه S است.

$\text{dom}(X_i)$: مجموعه مقادیر ممکن متغیر X_i است؛ یعنی

$$\text{dom}(X_i) = \{x_{i,1}, x_{i,2}, \dots, x_{i,|\text{dom}(X_i)}\}$$

$\text{dom}(Y)$: مجموعه مقادیر ممکن متغیر پاسخ Y است؛ یعنی

$$\text{dom}(Y) = \{c_1, c_2, \dots, c_{|\text{dom}(Y)}\}$$

همچنین شاخص دوپاره‌سازی متغیر پاسخ Y به‌صورت

$$\begin{aligned} \text{Twoing}(Y, \text{dom}_1(Y), \text{dom}_2(Y)) &= 0.25 \times \frac{|\sigma_{Y \in \text{dom}_1(Y)}(S)|}{|S|} \\ &\times \frac{|\sigma_{Y \in \text{dom}_2(Y)}(S)|}{|S|} \times \left(\sum_{c_j \in \text{dom}(Y)} \left| \frac{|\sigma_{Y \in \text{dom}_1(Y) \text{ AND } y=c_j}(S)|}{|\sigma_{Y \in \text{dom}_1(Y)}(S)|} \right. \right. \\ &\left. \left. - \frac{|\sigma_{Y \in \text{dom}_2(Y) \text{ AND } y=c_j}(S)|}{|\sigma_{Y \in \text{dom}_2(Y)}(S)|} \right| \right)^2 \end{aligned}$$

تعریف می‌شود که در آن رده‌های متغیر پاسخ Y (دارای هر چند رده که باشد) به دو رده $\text{dom}_1(Y)$ و $\text{dom}_2(Y)$ تقسیم شده است. لازم به ذکر است شاخص جینی متغیر X_i در این مجموعه به‌صورت

$$\text{Gini}(X_i, S) = \frac{|\sigma_{X_{i,1}}(S)|}{|S|} \text{Gini}(X_{i,1}) + \dots + \frac{|\sigma_{X_{i,|\text{dom}(X_i)}}(S)|}{|S|} \text{Gini}(X_{i,|\text{dom}(X_i)})$$

تعریف می‌شود که در آن

$$\text{Gini}(X_{i,j}, S) = 1 - \sum_{x_{i,j} \in \text{dom}(X_i)} \left(\frac{|\sigma_{X_i=x_{i,j}}(S)|}{|S|} \right)^2$$

است.

⁸Odd ratio

⁹Breiman

¹⁰Cost Complexity Pruning

¹¹Twoing index

¹²Gini index

۱.۳ نحوه هرس کردن الگوریتم CART

الگوریتم CART برای هرس کردن از معیار پیچیدگی هزینه هرس استفاده می‌کند. معیار پیچیدگی هزینه هرس در دو مرحله انجام می‌شود. در مرحله اول، دنباله‌ای از درخت‌های T_0, T_1, \dots, T_k بر اساس مجموعه داده S ساخته می‌شوند که T_0 درخت اصلی قبل از هرس کردن و T_k درخت ریشه است. درخت T_1 با هرس کردن T_0 ، یعنی با جایگزینی یک یا چند گره درونی از درخت T_0 با برگ ایجاد می‌شود. به همین ترتیب، درخت T_i با هرس کردن T_{i-1} به دست می‌آید. در مرحله دوم، با استفاده از برآورد خطای تعمیم‌یافته؛ یعنی

$$\varepsilon(T, D) = \sum_{(x,y) \in S} D(x, y) \times L(y, T(x)) \quad (2)$$

که در آن D توزیع مجموعه داده S و نامعلوم است، درختی که دارای کمترین مقدار خطای تعمیم‌یافته است به عنوان درخت هرس شده انتخاب می‌شود.

در مرحله اول، درخت T_{i+1} با جایگزینی یک یا چند گره در درخت T_i با برگ‌های مناسب به دست می‌آید. گره‌هایی که هرس می‌شوند گره‌هایی هستند که با هرس کردن آن‌ها نرخ خطای متناظر آن‌ها؛ یعنی

$$\alpha_i = \frac{\varepsilon(\text{pruned}(T_i, S), S) - \varepsilon(T, S)}{|\text{leaves}(T)| - |\text{leaves}(\text{pruned}(T_i, t))|} \quad (3)$$

کمترین مقدار باشد که در آن

$\varepsilon(T, S)$: نرخ خطای درخت T که بر اساس مجموعه داده S ساخته شده است.

$\varepsilon(\text{pruned}(T_i, S), S)$: نرخ خطای درخت هرس شده T_i که بر اساس مجموعه داده S ساخته شده است.

$\text{leaves}(T)$: تعداد برگ‌های درخت T است.

$\text{leaves}(\text{pruned}(T, t))$: تعداد برگ‌های درخت هرس شده T_i است. برای جزئیات بیشتر در مورد الگوریتم CART و کاربردهای آن، می‌توان به [۴]، [۱۰]، [۱۲] و [۱۵] مراجعه کرد.

۴ الگوریتم‌های یادگیری گروهی

الگوریتم‌های یادگیری گروهی^{۱۳} متشکل از دو یا بیش از دو الگوریتم یادگیری ماشین است. در استفاده برخی از الگوریتم‌های یادگیری گروهی از روش بازنمونه‌گیری خودگردان برای تولید چندین مجموعه داده آموزشی متفاوت استفاده می‌شود که در زیر به اختصار به بیان آن می‌پردازیم.

۱.۴ بازنمونه‌گیری خودگردان

روش بازنمونه‌گیری خودگردان^{۱۴} توسط افرون^{۱۵} [۷] معرفی شده است. فرض کنید Z_1, Z_2, \dots, Z_n یک نمونه تصادفی از توزیع نامعلوم F باشد. بر اساس نمونه تصادفی در نظر گرفته شده، یکی از رایج‌ترین برآوردگرهای تابع توزیع تجمعی نامعلوم جامعه، تابع توزیع تجربی است که به صورت

$$\hat{F}_n(z) = \frac{1}{n} \sum_{i=1}^n I(Z_i \leq z)$$

تعریف می‌شود که در آن $I(A) = 1$ اگر پیشامد A رخ دهد و در غیر این صورت $I(A) = 0$.

اکنون آماره $T = T(\mathbf{Z})$ را در نظر بگیرید که در آن $\mathbf{Z} = (Z_1, \dots, Z_n)$. بازنمونه‌گیری خودگردان روشی برای تقریب توزیع T است. در روش بازنمونه‌گیری خودگردان، با استخراج نمونه‌هایی به حجم n به روش باجایگذاری از نمونه اولیه Z_1, Z_2, \dots, Z_n می‌توان نمونه‌ای تصادفی از T به دست آورد که بر اساس آن، ویژگی‌های توزیع T را مورد بررسی قرار داد. بنابراین اگر B باز نمونه $Z_{(i)}^*$ ، $i = 1, 2, \dots, B$ ، از نمونه اولیه استخراج شود، آنگاه نمونه $T(\mathbf{Z}_{(1)}^*)$ ، \dots ، $T(\mathbf{Z}_{(B)}^*)$ از توزیع T خواهیم داشت. بنابراین می‌توان با روش‌های ناپارامتری موجود، ویژگی‌های توزیع T را مورد مطالعه قرار داد.

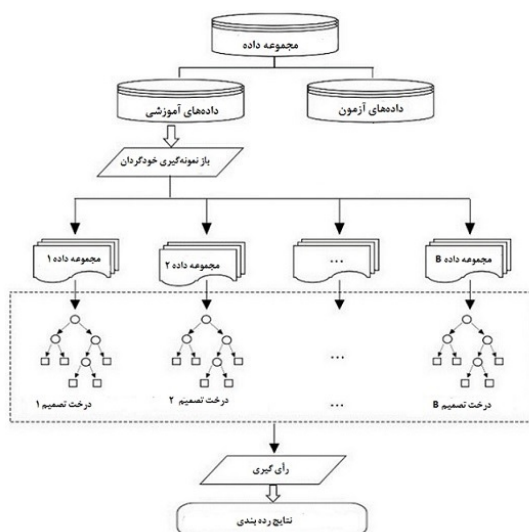
۲.۴ الگوریتم Bagging

یکی از ساده‌ترین و درعین حال موفق‌ترین الگوریتم‌های یادگیری گروهی الگوریتم Bagging در حوزه یادگیری با نظارت است که توسط برینمن [۵] معرفی شده است. این الگوریتم معمولاً برای الگوریتم‌های ناپایدار مانند درخت تصمیم بسیار کارا تر است و باعث کاهش واریانس، افزایش دقت و جلوگیری از بیش‌برازش می‌شود. الگوریتم Bagging به صورت زیر است. فرض کنید S یک مجموعه داده با n مشاهده است و G مجموعه داده آموزشی با m مشاهده است ($m < n$). از مجموعه داده آموزشی G به روش بازنمونه‌گیری خودگردان B مجموعه داده آموزشی تولید می‌شود. یک الگوریتم یادگیری ماشین رده‌بندی یکسان (الگوریتم درخت تصمیم، الگوریتم شبکه عصبی و ...) روی هر کدام از این B مجموعه تولید شده، پیاده‌سازی می‌شود و در نهایت نتایج این الگوریتم با هم ترکیب می‌شوند.

¹³Ensemble learning

¹⁴Bootstrap Resampling

¹⁵Efron



شکل ۱: فرایند ساخت جنگل تصادفی

به‌طور معمول $q = \sqrt{p}$ برای ساخت جنگل تصادفی در نظر گرفته می‌شود.

(۳) بر روی این مجموعه داده آموزشی تولیدشده با q متغیر پیشگو الگوریتم CART پیاده‌سازی می‌شود. مراحل ۱ تا ۳، B بار تکرار می‌شود.

هر درخت تصمیم در جنگل تصادفی پیش‌بینی می‌کند که هر مشاهده به کدام رده تعلق دارد و در نهایت جنگل تصادفی رده‌ای که اکثریت آرا را به دست آورده باشد برای هر مشاهده انتخاب می‌کند.

الگوریتم جنگل تصادفی تا حدودی شبیه الگوریتم Bagging است، با این تفاوت که این الگوریتم q متغیر پیشگو به‌صورت تصادفی انتخاب می‌کند. اگر $q = p$ ، در این صورت جنگل تصادفی همان الگوریتم Bagging است. شکل ۱ فرایند تشکیل جنگل تصادفی بر روی یک مجموعه داده را نشان می‌دهد. برای توضیحات بیشتر در مورد جنگل تصادفی و کاربردهای آن به [۶]، [۱۵] و [۱۷] مراجعه نمایید.

۴.۴ الگوریتم تقویت

الگوریتم تقویت یکی دیگر از الگوریتم‌های یادگیری گروهی در حوزه یادگیری با نظارت است که توسط شاپیر^{۱۶} [۱۰] معرفی شده است. برای آموزش الگوریتم تقویت از کل داده‌های آموزشی استفاده می‌شود، اما در هر بار تکرار وزن داده‌های آموزشی تغییر می‌کند. برای مثال، فرض کنید یک مجموعه داده آموزشی دارای سه مشاهده x_1 ، x_2 و x_3 است که به دو رده تعلق دارند. در ابتدا الگوریتم یادگیری ماشین h_1 روی داده‌های آموزشی پیاده‌سازی می‌شود. فرض کنید نتایج این مرحله به این صورت است که رده‌های

به‌عنوان مثال، فرض کنید x_i یک مشاهده جدید است و هرکدام از این B الگوریتم یک رده خاص برای x_i در نظر گرفته باشند. بر این اساس، رده‌ای برای x_i در نظر گرفته می‌شود که در B نتیجه به‌دست‌آمده دارای بیشترین فراوانی باشد. در این مقاله از روش Bagging CART استفاده می‌شود؛ یعنی بر روی هرکدام از B مجموعه داده، الگوریتم CART اعمال می‌شود. برای توضیحات بیشتر و کاربرد این الگوریتم می‌توان به [۵]، [۱۳] و [۱۴] مراجعه نمود.

۳.۴ جنگل تصادفی

یکی دیگر از الگوریتم‌های یادگیری گروهی در حوزه یادگیری با نظارت جنگل تصادفی است که توسط بریمن [۶] معرفی شده است. جنگل تصادفی شامل مجموعه‌ای از درخت‌های هرس نشده است. در مواردی که مجموعه داده شامل داده‌های گمشده، دورافتاده و نویز باشد از جنگل تصادفی می‌توان استفاده کرد. شایان‌ذکر است الگوریتم جنگل تصادفی در بعضی از مسائل رده‌بندی عملکرد بهتری نسبت به درخت تصمیم دارد ([۱]). فرض کنید S یک مجموعه داده شامل n مشاهده و p متغیر پیشگو و G مجموعه داده آموزشی با m مشاهده باشد ($m < n$). فرایند الگوریتم جنگل تصادفی به‌صورت زیر است:

- (۱) یک مجموعه داده آموزشی به روش بازنمونه‌گیری خودگردان از مجموعه داده آموزشی G تولید می‌شود.
- (۲) از مجموعه داده آموزشی تولیدشده (به روش بازنمونه‌گیری خودگردان) q متغیر پیشگو به روش نمونه‌گیری تصادفی بدون جایگذاری از p متغیر پیشگو انتخاب می‌شود ($q < p$).

است. همان‌گونه که پیش‌تر گفته شد، در ابتدا وزن یکسان به صورت

$$w_r(i) = \frac{1}{m}$$

برای تمام داده‌های آموزشی در نظر گرفته می‌شود. وزن مرحله‌ی r ام از الگوریتم مورد استفاده شده به صورت

$$\alpha_r = \frac{1}{\gamma} \ln \frac{1 - \varepsilon_r(h_r(x))}{\varepsilon_r(h_r(x))}$$

به دست می‌آید که در آن خطای وزنی الگوریتم در مرحله r ام، $r = 1, 2, \dots, B$ به صورت

$$\varepsilon_r(h_r(x)) = \frac{\sum_{i=1}^m w_r(i) \times L(y, h_r(x))}{\sum_{i=1}^m w_r(i)}$$

محاسبه می‌شود و در آن

$$L(y, h_r(x)) = \begin{cases} 0 & \text{if } y = h_r(x) \\ 1 & \text{if } y \neq h_r(x) \end{cases}$$

است و $h_r(x)$ رده پیش‌بینی شده مشاهده x با استفاده از الگوریتم h_r ، $w_r(i)$ وزن مشاهده‌ی i ام در مرحله r ام است. وزن نرمال شده هر داده آموزشی بعد از پیاده‌سازی مرحله r ام الگوریتم به صورت

$$w_{r+1}(i) = \frac{w_r(i)}{Z_r} \times \begin{cases} e^{-\alpha_r} & \text{if } y = h_r(x_i) \\ e^{+\alpha_r} & \text{if } y \neq h_r(x_i) \end{cases}$$

به دست می‌آید که در آن

$$Z_r = \sum_{j=1}^m w_r(j) \times \begin{cases} e^{-\alpha_r} & \text{if } y = h_r(x_j) \\ e^{+\alpha_r} & \text{if } y \neq h_r(x_j) \end{cases}$$

در نهایت، الگوریتم نهایی

$$H(x) = \sum_{r=1}^B \alpha_r h_r(x)$$

که متوسط وزنی از الگوریتم‌های به دست آمده در B مرحله است مورد استفاده قرار می‌گیرد. رده مشاهده x براساس علامت $H(x)$ مشخص می‌شود. اگر علامت $H(x)$ مثبت باشد، مشاهده x به رده (+) تعلق دارد، اما اگر علامت $H(x)$ باشد مشاهده به رده (-) تعلق دارد.

لازم به ذکر است، در این مطالعه، الگوریتم CART بر روی B مجموعه داده آموزشی با وزن‌های متفاوت پیاده‌سازی شده است که به اصطلاح به آن تقویت سازوار CART گفته می‌شود. برای توضیحات بیشتر و کاربرد این الگوریتم، می‌توان به [۱۳] و [۱۴] مراجعه نمود.

x_1 و x_2 توسط الگوریتم h_1 به درستی پیش‌بینی شده است، اما رده x_2 به اشتباه پیش‌بینی شده است. در تکرار دوم، به دلیل اینکه رده مشاهده x_2 در مرحله قبل به اشتباه رده‌بندی شده است، بنابراین وزن بیشتری به x_2 اختصاص داده می‌شود. اکنون الگوریتم h_2 روی مشاهدات با وزن‌های مشخص شده پیاده‌سازی می‌شود. فرض کنید در این مرحله رده‌ی x_2 به اشتباه پیش‌بینی شده است. در تکرار سوم، وزن x_2 افزایش و وزن x_1 کاهش می‌یابد و مشاهده x_1 نیز کمترین وزن را به خود اختصاص می‌دهد. سپس، الگوریتم h_3 روی این داده‌های آموزشی که وزن‌های آن‌ها تغییر کرده است، پیاده‌سازی می‌شود. در نهایت نتایج الگوریتم‌های h_1 ، h_2 و h_3 باهم ترکیب می‌شوند. لازم به ذکر است الگوریتم‌های h_1 ، h_2 و h_3 همگی از یک نوع الگوریتم هستند. به عنوان مثال، همگی درخت تصمیم یا همگی شبکه عصبی و... می‌باشند. شایان ذکر است که الگوریتم تقویت شامل الگوریتم‌های متفاوت است. در ادامه، یکی از پرکاربردترین آن‌ها به نام الگوریتم تقویت سازوار را بیان خواهیم کرد.

۱۰.۴.۴ الگوریتم تقویت سازوار

الگوریتم تقویت سازوار توسط شاپیر و فروند^{۱۷} [۸] معرفی شد. فرض کنید S یک مجموعه داده شامل n مشاهده از p متغیر پیشگو و یک متغیر پاسخ (ویژگی هدف) دو رده‌ای است که m مشاهده آن داده‌های آموزشی را تشکیل می‌دهند. برای انجام الگوریتم تقویت سازوار، در ابتدا از بین p متغیر پیشگو، یک متغیر که دارای کمترین مقدار شاخص جینی یا بیشترین مقدار معیار بهره اطلاع (IG) باشد، انتخاب می‌شود. معیار بهره اطلاع متغیر پیشگو X_i در این مجموعه داده به صورت

$$IG(X_i, S) = E(Y, S) - \sum_{x_{i,j} \in \text{dom}(X_i)} \frac{|\sigma_{X_i=x_{i,j}}(S)|}{|S|} \quad (۴)$$

تعریف می‌شود که در آن

$$E(Y, S) = - \sum_{c_j \in \text{dom}(Y)} \frac{|\sigma_{Y=c_j}(S)|}{|S|} \log_2 \frac{|\sigma_{Y=c_j}(S)|}{|S|}$$

$E(Y, B)$: آنتروپی متغیر پاسخ Y (ویژگی هدف) براساس مجموعه داده B است که با $E(Y|B)$ نیز نشان داده می‌شود.

فرایند الگوریتم تقویت سازوار بر روی یک متغیر انتخاب شده از بین p متغیر پیشگو به شرح زیر است. فرض کنید D یک مجموعه داده آموزشی شامل m مشاهده از یک متغیر پیشگو و یک متغیر پاسخ به صورت

$$D = \{(x_i, y_i)\}_{i=1}^m \quad x_i \in X \quad y_i \in \{+1, -1\}$$

¹⁷Freund

¹⁸Information Gain

۵ معیار ارزیابی عملکرد الگوریتم

به دست می‌آید. با توجه به معیار دقت می‌توان نتیجه گرفت که یک الگوریتم به چه میزان به درستی آموزش داده شده است و یا کارایی الگوریتم به‌طور کلی چگونه است.

در معیارهای ارزیابی یک الگوریتم از ماتریس درهم‌ریختگی استفاده می‌شود که در زیر به اختصار به تعریف آن می‌پردازیم.

ماتریس درهم‌ریختگی

ماتریس درهم‌ریختگی^{۱۹} یک ماتریس مربع $k \times k$ است که k تعداد رده‌های متغیر پاسخ در مسائل رده‌بندی است. به منظور سهولت کار، فرض کنید $k = 2$ ؛ یعنی متغیر پاسخ شامل دو رده مثبت و منفی است. مانند جدول ۱، مؤلفه‌های روی قطر فرعی ماتریس درهم‌ریختگی تعداد مشاهداتی است که الگوریتم رده‌ی آن‌ها را به درستی پیش‌بینی کرده است؛ یعنی

TP ^{۲۰}: تعداد مشاهداتی است که به رده مثبت تعلق دارند و الگوریتم به درستی رده مثبت برای آن‌ها پیش‌بینی کرده است.

TN ^{۲۱}: تعداد مشاهداتی است که به منفی تعلق دارند و الگوریتم به درستی رده منفی برای آن‌ها پیش‌بینی کرده است.

قطر اصلی ماتریس درهم‌ریختگی تعداد مشاهداتی است که الگوریتم رده آن‌ها را به اشتباه پیش‌بینی کرده است؛ یعنی

FN ^{۲۲}: تعداد مشاهداتی است که به رده مثبت تعلق دارند، اما الگوریتم به اشتباه رده منفی برای آن‌ها پیش‌بینی کرده است.

FP ^{۲۳}: تعداد مشاهداتی است که به رده منفی تعلق دارند، اما الگوریتم به اشتباه رده مثبت برای آن‌ها پیش‌بینی کرده است.

جدول ۱. ماتریس درهم‌ریختگی

پیش‌بینی		واقعی
منفی	مثبت	
FN	TP	مثبت
TN	FP	منفی

برای ارزیابی عملکرد الگوریتم‌ها معیارهای متفاوتی وجود دارد که در ادامه به توضیح معیارهای دقت، حساسیت و صحت می‌پردازیم.

۱.۵ دقت

دقت یکی از متداول‌ترین معیارهای ارزیابی الگوریتم‌ها در مسائل رده‌بندی است. این معیار بیانگر این است که رده چه درصدی از مشاهدات توسط الگوریتم به‌کاررفته به درستی پیش‌بینی شده است. با توجه به جدول ۱، دقت یک الگوریتم به صورت

$$Accuracy = \frac{TP + TN}{TP + FN + TN + FP}$$

۲.۵ حساسیت

در برخی از مسائل رده‌بندی، پیش‌بینی درست مشاهدات مربوط به رده مثبت اهمیت بیشتری دارد. معیار حساسیت که به صورت زیر تعریف می‌شود، دقت پیش‌بینی در رده مثبت را اندازه‌گیری می‌کند.

$$Sensitivity = \frac{TP}{TP + FN}$$

۳.۵ صحت

معیار صحت، که مقدار پیشگویی مثبت^{۲۴} هم گفته می‌شود، میزان صحت رده پیش‌بینی شده مثبت توسط الگوریتم برای یک مشاهده را اندازه‌گیری می‌کند. معیار صحت به صورت زیر تعریف می‌شود.

$$Precision = \frac{TP}{TP + FP}$$

در حالت کلی، برای محاسبه معیار حساسیت یک الگوریتم، هر بار یکی از رده‌ها به‌عنوان رده مثبت در نظر گرفته شده و معیار حساسیت برای آن رده به دست می‌آید. میانگین مقادیرهای به‌دست‌آمده برای تمام رده‌ها به‌عنوان حساسیت الگوریتم در نظر گرفته می‌شود. به شیوه مشابه، معیار صحت یک الگوریتم محاسبه می‌شود. برای جزئیات بیشتر در مورد معیارهای ارزیابی می‌توان به [۱۸] مراجعه کرد.

۶ مقایسه الگوریتم‌ها

در این بخش، بر اساس معیارهای دقت، حساسیت و صحت و با استفاده از ۴ مجموعه داده واقعی و یک مثال شبیه‌سازی، به مقایسه کارایی الگوریتم‌های رگرسیون لوزستیک، CART، جنگل تصادفی، Bagging CART و تقویت سازوار CART می‌پردازیم. همچنین برای اعتبارسنجی آن‌ها از روش اعتبارسنجی ۱۰ گروهی استفاده می‌شود.

لازم به ذکر است برای اجرای الگوریتم‌های ذکر شده از نرم‌افزار Python 3.7 در محیط Anaconda 3 (2020.07) استفاده شده است.

¹⁹Confusion Matrix

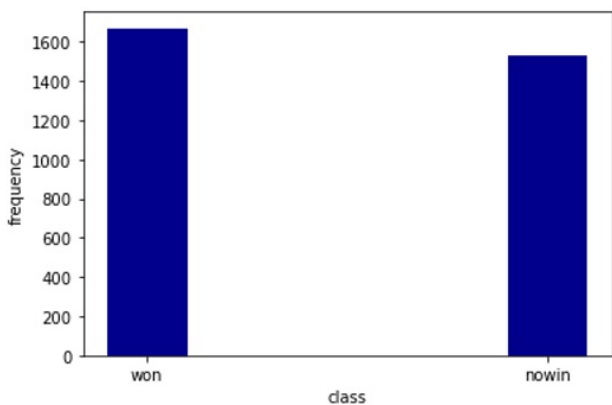
²⁰True Positive

²¹True Negative

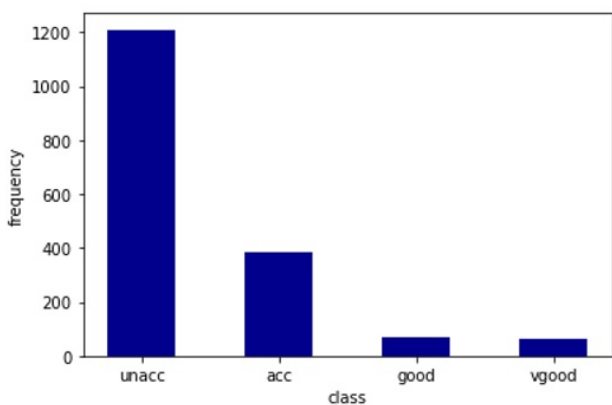
²²False Negative

²³False Positive

²⁴Positive predictive value



شکل ۳. نمودار میله‌ای متغیر پاسخ برای مجموعه داده بازی شطرنج

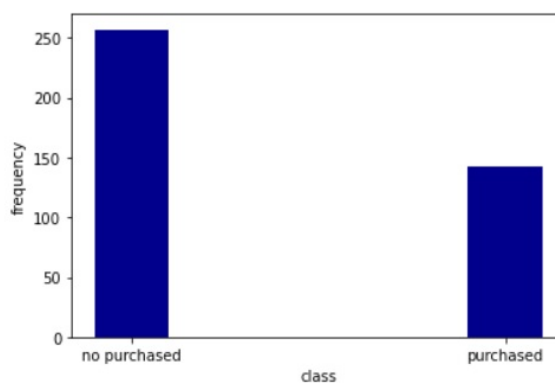


شکل ۴. نمودار میله‌ای متغیر پاسخ برای مجموعه داده ارزیابی اتومبیل

ویژگی‌های ۴ مجموعه داده مورداستفاده به‌طور اختصار در جدول ۲ آورده شده است.

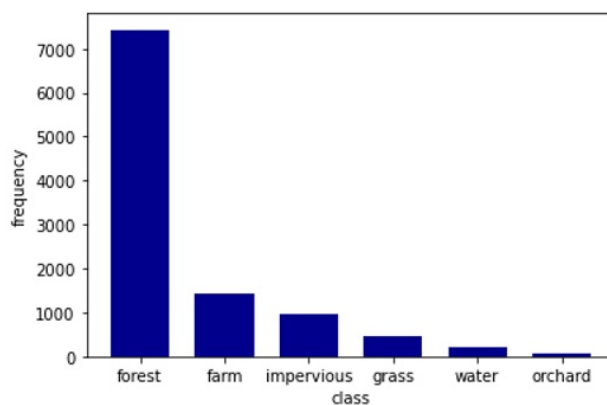
جدول ۲. معرفی مجموعه داده‌ها

مجموعه داده	تعداد مشاهده	تعداد متغیرهای پیشگو	نوع متغیرها	تعداد رده متغیر پاسخ
تلیغات شبکه اجتماعی	۴۰۰	۳	کمی و رستهای	۲
تصاویر ماهواره‌ای	۱۰۵۴۵	۲۸	کمی و رستهای	۶
بازی شطرنج	۳۱۹۶	۳۶	رستهای	۲
ارزیابی اتومبیل	۱۷۲۸	۴	رستهای	۴



شکل ۱. نمودار میله‌ای متغیر پاسخ برای مجموعه داده تبلیغات شبکه اجتماعی

مجموعه داده تبلیغات شبکه اجتماعی^{۲۵} مربوط به مشتری‌های سایت trusted است که شامل ۴۰۰ مشاهده از ۳ متغیر پیشگو و یک متغیر پاسخ با دو رده خرید به‌صورت اینترنتی و عدم خرید است. مجموعه داده تصاویر ماهواره‌ای^{۲۶} مربوط به سال‌های ۲۰۱۴ تا ۲۰۱۵ است و شامل ۱۰۵۴۵ مشاهده از ۲۸ متغیر پیشگو و یک متغیر پاسخ با ۶ رده جنگل، مزرعه، مناطق غیرقابل نفوذ، فضای سبز، آب و باغ است. مجموعه داده بازی شطرنج^{۲۷} شامل ۳۱۹۶ مشاهده از ۳۶ متغیر پیشگو و یک متغیر پاسخ با دو رده برنده و بازنده است. مجموعه داده ارزیابی اتومبیل^{۲۸} شامل ۱۷۲۸ مشاهده از ۴ متغیر پیشگو و یک متغیر پاسخ با ۴ رده غیرقابل قبول، قابل قبول، خوب و خیلی خوب است. نمودار میله‌ای متغیر پاسخ در مجموعه داده‌های در نظر گرفته‌شده در شکل‌های ۱، ۲، ۳ و ۴ رسم شده است. الگوریتم‌های رگرسیون لوژستیک، CART، جنگل تصادفی، Bagging CART و تقویت سازوار CART روی این چهار مجموعه



شکل ۲. نمودار میله‌ای متغیر پاسخ برای مجموعه داده تصاویر ماهواره‌ای

^{۲۵}مجموعه داده Social Network Ads در سایت www.superdatascience.com/machine-learning موجود است.

^{۲۶}مجموعه داده Crowdsourced Mapping در سایت archive.ics.uci.edu/ml/index.php موجود است.

^{۲۷}مجموعه داده Chess (King-Rook -King-Pawn) در سایت archive.ics.uci.edu/ml/index.php موجود است.

^{۲۸}مجموعه داده Car Evaluation در سایت archive.ics.uci.edu/ml/index.php موجود است.

داده بازی شطرنج

معیار	رگرسیون لوژیستیک	CART	Bagging CART	سازوار تقویت CART	جنگل تصادفی
دقت	۰/۸۸۱۴	۰/۹۴۹۳	۰/۷۶۶۰	۰/۹۰۲۱	۰/۸۹۴۲
حساسیت	۰/۹۵۶۱	۰/۹۸۲۱	۰/۹۱۹۵	۰/۹۶۲۰	۰/۹۷۹۵
صحت	۰/۹۵۶۳	۰/۹۸۲۲	۰/۹۸۴۴	۰/۹۶۳۷	۰/۹۷۵۶
مدت زمان اجرا (ثانیه)	۴/۱۶۹۰	۰/۳۳۰۶	۱/۱۷۶۳	۲۴/۱۴۹۰	۰/۸۶۳۵

با توجه به جدول ۶، برای مجموعه داده ارزیابی اتومیل الگوریتم CART دارای بیشترین مقدار دقت، حساسیت و صحت به ترتیب با مقدارهای ۰/۷۴۱۷، ۰/۷۸۰۶ و ۰/۸۲۵۹ است. همچنین الگوریتم جنگل تصادفی دارای کارایی نزدیک (رقابت پذیر) به الگوریتم CART است. برای این مجموعه داده، الگوریتم تقویت سازوار CART و سپس رگرسیون لوژیستیک دارای ضعیف ترین کارایی هستند. الگوریتم تقویت سازوار CART نسبت به بقیه الگوریتم ها به مدت زمان بیشتری برای اجرا نیاز دارد.

جدول ۶. معیارهای ارزیابی الگوریتم ها برای مجموعه داده

ارزیابی اتومیل

معیار	رگرسیون لوژیستیک	CART	Bagging CART	سازوار تقویت CART	جنگل تصادفی
دقت	۰/۶۹۱۹	۰/۷۴۱۷	۰/۷۰۰۰	۰/۲۴۳۷	۰/۷۳۳۰
حساسیت	۰/۷۳۳۸	۰/۷۸۰۶	۰/۷۰۰۲	۰/۲۸۴۰	۰/۷۷۵۴
صحت	۰/۷۱۳۲	۰/۸۲۵۹	۰/۴۹۰۸	۰/۳۲۵۳	۰/۸۱۸۹
مدت زمان اجرا (ثانیه)	۲/۰۶۹۳	۰/۰۹۰۱	۰/۳۹۱۱	۷/۵۶۴۰	۰/۴۷۰۲

مثال ۱۰۶ (شبیه سازی).

مدل رگرسیون لوژیستیک دودویی به صورت

$$\log \frac{\pi(\mathbf{X})}{1 - \pi(\mathbf{X})} = 1 + 2X_1 + 3X_2$$

را در نظر بگیرید که در آن X_1 و X_2 متغیرهای تصادفی نرمال استاندارد هستند. از این مدل ۱۰۰۰ مشاهده تولید و الگوریتم های رگرسیون لوژیستیک، CART، جنگل تصادفی، Bagging CART و تقویت سازوار CART روی داده ها اجرا می شود. نتایج کارایی الگوریتم ها براساس معیارهای دقت، صحت و حساسیت در جدول ۷ آورده شده است. از نتایج به دست آمده، مشاهده می شود با وجود اینکه داده ها از مدل رگرسیون لوژیستیک تولید شده است، اما الگوریتم های یادگیری ماشین به ویژه جنگل تصادفی و تقویت سازوار CART دارای کارایی نزدیک به رگرسیون لوژیستیک می باشند.

جدول ۷. معیارهای ارزیابی الگوریتم ها برای مجموعه

داده اجرا و دقت پیش بینی آن ها براساس معیارهای دقت، حساسیت و صحت در جدول های ۳، ۴، ۵ و ۶ محاسبه شده است. با توجه به جدول ۳، برای مجموعه داده تبلیغات شبکه اجتماعی، الگوریتم تقویت سازوار CART دارای بیشترین دقت با مقدار ۰/۸۸۰ و جنگل تصادفی دارای بیشترین حساسیت و صحت به ترتیب با مقدارهای ۰/۸۹۷۵ و ۰/۹۰۰۸ است. رگرسیون لوژیستیک دارای ضعیف ترین کارایی براساس معیارهای ارزیابی است. الگوریتم تقویت سازوار CART نسبت به بقیه الگوریتم ها به مدت زمان بیشتری برای اجرا نیاز دارد.

جدول ۳. معیارهای ارزیابی الگوریتم ها برای مجموعه داده

تبلیغات شبکه اجتماعی

معیار	رگرسیون لوژیستیک	CART	Bagging CART	سازوار تقویت CART	جنگل تصادفی
دقت	۰/۶۹۵۰	۰/۸۴۷۴	۰/۷۸۷۵	۸۸۰۰	۰/۸۷۷
حساسیت	۰/۶۵	۰/۸۵۲۵	۰/۸۱۲۵	۰/۸۱۱۵	۰/۸۹۷۵
صحت	۰/۴۱۹۴	۰/۸۶۰۴	۰/۸۰۶۷	۰/۸۸۶۳	۰/۹۰۸۷
مدت زمان اجرا (ثانیه)	۰/۲۲۵۴	۰/۰۳۷۳	۰/۲۶۴۲	۲/۸۷۷۲	۰/۳۷۶۳

با توجه به جدول ۴، برای مجموعه داده تصاویر ماهواره ای الگوریتم جنگل تصادفی دارای بیشترین مقدار دقت، حساسیت و صحت به ترتیب با مقدارهای ۰/۸۳۸۷، ۰/۹۳۷۶ و ۰/۹۳۵۰ است. الگوریتم تقویت سازوار CART و سپس رگرسیون لوژیستیک دارای ضعیف ترین کارایی براساس معیارهای ارزیابی هستند. الگوریتم رگرسیون لوژیستیک نسبت به بقیه الگوریتم ها به مدت زمان بیشتری برای اجرا نیاز دارد.

جدول ۴. معیارهای ارزیابی الگوریتم ها برای مجموعه داده

تصاویر ماهواره ای

معیار	رگرسیون لوژیستیک	CART	Bagging CART	سازوار تقویت CART	جنگل تصادفی
دقت	۰/۷۵۴۳	۰/۷۷۶۵	۰/۸۲۸۴	۰/۵۵۹۱	۰/۸۲۸۷
حساسیت	۰/۸۸۰۵	۰/۸۸۹۰	۰/۹۲۶۲	۰/۶۶۶۷	۰/۹۳۷۶
صحت	۰/۸۷۴۶	۰/۸۹۲۶	۰/۹۲۳۵	۰/۷۹۹۸	۰/۹۳۵۰
مدت زمان اجرا (ثانیه)	۱۸۱/۳۳۱۹	۹/۸۸۳۵	۱۴/۰۷۵۸	۱۱۶/۶۱۲۸	۸/۷۱۸۸

با توجه به جدول ۵، برای مجموعه داده بازی شطرنج الگوریتم CART دارای بیشترین دقت، حساسیت و صحت به ترتیب با مقدارهای ۰/۹۴۹۳، ۰/۹۸۲۱ و ۰/۹۸۲۲ است. الگوریتم های جنگل تصادفی و تقویت سازوار CART دارای کارایی نزدیک (رقابت پذیر) به الگوریتم CART هستند. الگوریتم Bagging CART و سپس رگرسیون لوژیستیک دارای ضعیف ترین کارایی هستند. رگرسیون لوژیستیک نسبت به بقیه الگوریتم ها به مدت زمان بیشتری برای اجرا نیاز دارد.

جدول ۵. معیارهای ارزیابی الگوریتم ها برای مجموعه

۷ بحث و نتیجه‌گیری

در این مقاله، در بحث رده‌بندی، مدل رگرسیون لوژیستیک با الگوریتم‌های یادگیری ماشین CART، جنگل تصادفی، Bag-ging و تقویت سازوار CART مورد مقایسه قرار گرفت. برای انجام مقایسه، ۴ مجموعه داده واقعی و همچنین یک مثال شبیه‌سازی شده مورد استفاده قرار گرفت. براساس نتایج به دست آمده از معیارهای ارزیابی (دقت، حساسیت و صحت)، مشاهده گردید که برای هر مجموعه داده حداقل یک الگوریتم یادگیری ماشین وجود دارد که دارای کارایی بهتر از مدل رگرسیون لوژیستیک است. همچنین ضمن اینکه انتخاب بهترین الگوریتم بستگی به نوع داده از نظر تعداد مشاهدات، تعداد رده‌های متغیر پاسخ دارد، اما می‌توان گفت الگوریتم‌های جنگل تصادفی و CART دارای کارایی قابل قبول‌تر با توجه به انواع داده‌ها هستند.

داده شبیه‌سازی شده

معیار	رگرسیون لوژیستیک	CART	Bagging CART	سازوار تقویت CART	جنگل تصادفی
دقت	۰/۸۶۱۰	۰/۸۰۸	۰/۷۷۰	۰/۸۴۲	۰/۸۳۰۰
حساسیت	۰/۸۸۶۲	۰/۸۳۸۳	۰/۷۵۲۲	۰/۸۸۰۸	۰/۸۵۰۷
صحت	۰/۸۸۱۷	۰/۸۴۰۰	۰/۷۶۳۳	۰/۸۵۶۳	۰/۸۶۱۷
مدت زمان اجرا (ثانیه)	۰/۱۲۹۱	۰/۰۶۳۳	۰/۷۴۷۳	۳/۶۹۶۳	۰/۴۱۲۸

با توجه به نتایج به دست آمده، از مقایسه مدل رگرسیون لوژیستیک با الگوریتم‌های یادگیری ماشین با استفاده از داده‌های واقعی می‌توان گفت که الگوریتم‌های جنگل تصادفی و CART دارای عملکرد بهتر از رگرسیون لوژیستیک هستند.

مراجع

- [1] Aggarwal, C. C. (2014). *Data classification: algorithms and applications*. CRC press, New York, USA.
- [2] Agresti, A. (2003). *Categorical data analysis*. John Wiley and Sons, Florida, United States.
- [3] Alpaydin, E. (2014). *Introduction to machine learning*. MIT press, London, England.
- [4] Breiman, L., Friedman, J., Stone, C. J. and Olshen, R. A. (1984). *Classification and regression trees*. CRC press, London, New York, Washington, D.C.
- [5] Breiman, L. (1996). Bagging predictors. *Machine learning*, **24(2)**, 123-140.
- [6] Breiman, L. (2001). Random forests. *Machine learning*, **45(1)**, 5-32.
- [7] Efron, B. (1979). Bootstrap Methods: Another Look at the Jackknife. *The Annals of Statistics*, **7(1)**, 1-26.
- [8] Freund, Y., and R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, **55(1)**, 119-139.
- [9] Kass, G. V. (1980). An exploratory technique for investigating large quantities of categorical data. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **29(2)**, 119-127.
- [10] Maimon, O. Z., Rokach, L. (2014). *Data mining with decision trees: theory and applications*. World scientific. Singapore.
- [11] Schapire, R. E. (1990). The strength of weak learn ability. *Machine learning*, **5(2)**, 197-227.
- [12] Soleimanpour, S. M., Mesbah, S. H., and Hedayati, B. (2018). Application of CART decision tree data mining to determine the most effective drinking water quality factors (case study: Kazeroon plain, Fars province). *Iranian Journal of Health and Environment*, **11(1)**, 1-14.
- [13] Syarif, I., Zaluska, E., Prugel-Bennett, A. and Wills, G. (2012) Application of bagging, boosting and stacking to intrusion detection. In *International Workshop on Machine Learning and Data Mining in Pattern Recognition* (pp. 593-602). Springer, Berlin, Heidelberg.

- [14] Taser, P. Y. (2021). Application of bagging and boosting approaches using decision tree-based algorithms in diabetes risk prediction. *In Multidisciplinary Digital Publishing Institute Proceedings* **74**, 6.
- [15] Magidi, J., Nhamo, L., Mpandeli, S., and Mabhaudhi, T. (2021). Application of the random forest classifier to map irrigated areas using google earth engine. *Remote Sensing*, **13(5)**, 876.
- [16] Mitchell, T. M. (1997). *Machine Learning*. McGraw-Hill, Ithaca.
- [17] You, J., van der Klein, S. A., Lou, E. and Zuidhof, M. J. (2020). Application of random forest classification to predict daily oviposition events in broiler breeders fed by precision feeding system. *Computers and Electronics in Agriculture*, **175**, 105526.
- [18] Zhu, W., Zeng, N. and Wang, N. (2010). Sensitivity, specificity, accuracy, associated confidence interval and ROC analysis with practical SAS implementations. *NESUG proceedings: health care and life sciences, Baltimore, Maryland*, **19**, 67.

پیوست

```
# code of computed accuracy of Logistic Regression for Social Network Ads data set  
with a python 3.7 (anaconda 3) Software:
```

```
from pandas import read_csv  
from sklearn.model_selection import KFold  
from sklearn.model_selection import cross_val_score  
from sklearn.linear_model import LogisticRegression  
df=read_csv(r"C:Desktop\data\Social_Network_Ads.csv")  
array=df.values  
x=array[:,0:3]  
y=array[:,3]  
kfold=KFold(n_splits=10, random_state=7)  
model=LogisticRegression(max_iter=5000,solver='lbfgs')  
results=cross_val_score(model, x, y, cv=10,scoring='accuracy')  
results.mean()
```

```
-----  
# code of computed accuracy of CART for Social Network Ads data set
```

```
from pandas import read_csv  
from sklearn.model_selection import KFold  
from sklearn.model_selection import cross_val_score  
a=read_csv(r"C:Desktop\data\Social_Network_Ads.csv")  
array=a.values  
x=array[:,0:3]  
y=array[:,3]  
kfold=KFold(n_splits=10, random_state=7)  
from sklearn.model_selection import cross_val_score  
from sklearn.datasets import make_blobs  
from sklearn.tree import DecisionTreeClassifier  
cl1 = DecisionTreeClassifier(max_depth=None, min_samples_split=2, random_state=7)  
results=cross_val_score(model, x, y, cv=10,scoring='accuracy')  
results.mean()
```

```
# code of computed accuracy of Bagging CART for Social Network Ads data set
```

```
from pandas import read_csv  
from sklearn.model_selection import KFold  
from sklearn.model_selection import cross_val_score  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.ensemble import BaggingClassifier  
df=read_csv(r"C:Desktop\data\Social_Network_Ads.csv")  
array=df.values
```

```
x=array[:,0:3]
y=array[:,3]
kfold=KFold(n_splits=10, random_state=7)
bagging = BaggingClassifier( DecisionTreeClassifier(), max_samples=0.5, max_features=0.5)
results=cross_val_score(bagging, x, y, cv=10,scoring='accuracy')
results.mean()
```

```
# code of computed accuracy of Adaboosting CART for Social Network Ads data set
```

```
from pandas import read_csv
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import AdaBoostClassifier
df=read_csv(r"C:Desktop\data\Social_Network_Ads.csv")
array=df.values
x=array[:,0:3]
y=arry[:,3]
kfold=KFold(n_splits=10, random_state=7)
boosting = AdaBoostClassifier(DecisionTreeClassifier(),n_estimators=100)
results=cross_val_score(boosing, x, y, cv=10,scoring='accuracy')
results.mean()
```

```
# code of computed accuracy of Random Forest for Social Network Ads data set
```

```
from pandas import read_csv
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
a=read_csv(r"C:Desktop\data\Social_Network_Ads.csv")
array=a.values
x=array[:,0:3]
y=arry[:,3]
kfold=KFold(n_splits=10, random_state=7)
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestClassifier
cl1 = RandomForestClassifier(n_estimators=10, max_depth=None,
min_samples_split=2, random_state=7)
results=cross_val_score(cl1, x, y, cv=10,scoring='accuracy')
results.mean()
```

Comparison of logistic regression with some machine learning methods in classifying data

Tayebeh Karami¹, Muhyiddin Izadi¹, Mehrdad Niaparast¹

Abstract:

One of the most important issues in various sciences is classification. Logistic regression is one of the statistical methods for data classification in which the distribution of data is supposed to be known. In addition to statistical methods, researchers are now using other methods to classify data such as machine learning algorithms that do not require the data distribution to be known. In this paper, logistic regression and some machine learning algorithms including CART, random forest, Bagging and Boosting of supervising learning are discussed. Using four real data sets and a simulation example, we compare the performance of the logistic regression with machine learning algorithms in terms of accuracy, sensitivity and precision measures.

Keywords: Decision Tree, Ensemble Learning, Random Forest and Supervised Learning.