

DATA MINING HOMEWORK 2

TAHSIN ILKHAS ZADEH STD.NO: 400422034

Datasets:

1. Mobile Price
2. Apartment Rental Offers in Germany



Shahid Beheshti University, Department of computer science

Data Mining Course: Dr.Farahani, Dr.Parand

Teacher Assistant: Ali Sharifi

May 18, 2022

۲	سوال ۶. بررسی استفاده از تبدیل log transform و سایر تبدیلات
۴	سوال ۸. cross validation and bootstrapping
۶	سوال ۹. 5x2 cross validation
۷	سوال ۱۰. الگوریتم های متداول برای ساخت درخت تصمیم Decision tree در داده کاوی
۱۰	سوال ۱۳. هرس درخت تصمیم ، Prepruning و Postpruning
۱۱	سوال ۱۴. The elbow method
۱۴	سوال ۲ تسک امتیازی : چگونه می توان با استفاده از statistical significance tests به مقایسه مدل ها پرداخت؟
۱۹	سوال ۳ تسک امتیازی : معیار Matthews Correlation Coefficient (MCC) چیست و کجا استفاده می شود؟

سوال ۶. بررسی استفاده از تبدیل log transform و سایر تبدیلات

چرا به تبدیل ویژگی و مقیاس بندی نیاز داریم؟

اغلب، ما مجموعه داده هایی داریم که در آن ستون های مختلف، واحدهای متفاوتی دارند - مثلاً یک ستون می تواند بر حسب کیلوگرم باشد، در حالی که ستون دیگر می تواند بر حسب سانتی متر باشد. علاوه بر این، می توانیم ستون هایی مانند درآمد داشته باشیم که می تواند از ۲۰۰۰۰ تا ۱۰۰۰۰۰ و حتی بیشتر باشد. در حالی که یک ستون سنی که می تواند از ۰ تا ۱۰۰ (حداکثر) متغیر باشد. بنابراین، درآمد حدود ۱۰۰۰ برابر بیشتر از سن است. اما چگونه می توان مطمئن بود که مدل با هر دو متغیر به طور یکسان رفتار می کند؟

وقتی این feature را به مدل آن طور که هست می دهیم، احتمال اینکه درآمد به دلیل ارزش بزرگ تر بر نتیجه تأثیر بگذارد، وجود دارد. اما این لزوماً به این معنی نیست که به عنوان یک عامل پیش بینی کننده، فیچر درآمد، مهم تر است. بنابراین، برای اهمیت دادن به سن و درآمد، به مقیاس بندی ویژگی نیاز داریم.

در بیشتر نمونه های مدل های یادگیری ماشینی، مقیاس کننده استاندارد یا مقیاس کننده MinMax را مشاهده کرده اید. با این حال، کتابخانه قدرتمند sklearn بسیاری از تکنیک های مقیاس بندی تبدیل ویژگی های دیگر را نیز ارائه می کند که بسته به داده هایی که با آن ها سروکار داریم، می توانیم از آنها استفاده کنیم.

تبدیلات دیگری نیز وجود دارد که در موارد لزوم می توان از آنها استفاده کرد مثلاً:

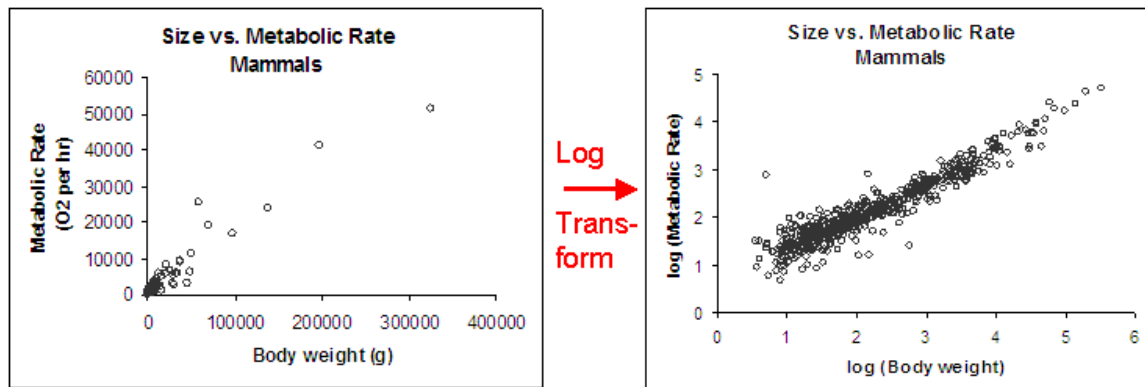
- MinMax Scaler : مقیاس کننده MinMax یکی از ساده ترین مقیاس کننده ها برای درک است که همه داده ها را بین ۰ و ۱ مقیاس می کند.
- Standard Scaler : برای هر ویژگی، مقیاس کننده استاندارد مقادیر را طوری مقیاس می دهد که میانگین ۰ و انحراف استاندارد ۱ (یا واریانس) باشد.
- MaxAbsScaler : مقیاس کننده MaxAbs ابتدا قدر مطلق هر مقدار در ستون را می گیرد و سپس حداکثر مقدار مطلق هر ستون را می گیرد و هر مقدار در ستون را بر حداکثر مقدار تقسیم می کند. این عملیات داده ها را بین محدوده [-۱، ۱] مقیاس می کند.
- Robust Scaler : Robust Scaler، به موارد دورافتاده حساس نیست. این مقیاس کننده، میانه را از مقیاس داده ها حذف می کند به صورت: $IQR = Q3 - Q1$ و $x_scaled = (x - Q1)/(Q3 - Q1)$
- Quantile Transformer Scaler
- Log Transformation
- Power Transformer Scaler
- Unit Vector Scaler/Normalizer

تبدیل log داده ها چیست؟

تبدیل log transform یک روش تبدیل داده است که در آن هر متغیر x را با یک $\log(x)$ جایگزین می کند. انتخاب پایه لگاریتمی معمولاً به تحلیلگر واگذار می شود و به اهداف مدل سازی آماری بستگی دارد. فرم کلی یک تابع لگاریتمی به این صورت است: $f(x) = k + a \log_b(x - h)$ که در آن a, b, k و h اعداد حقیقی هستند به طوری که b یک عدد مثبت $\neq 1$ است و $x - h > 0$. به عنوان مثال:

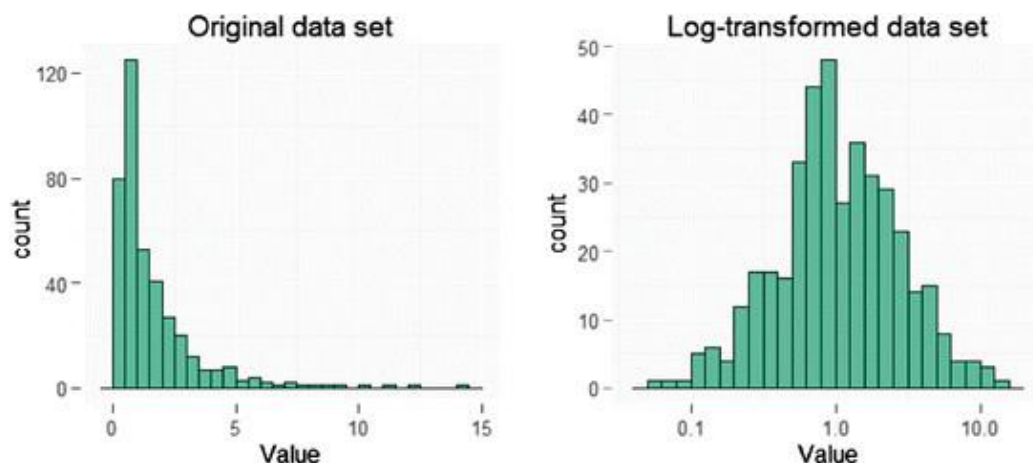
$$f(x) = 4 + 3 \log(x - 5)$$

معمولا در داده کاوی ، داده ها را برای سازماندهی بهتر آنها تغییر شکل می دهند. استفاده از داده های تبدیل شده هم برای انسان و هم برای رایانه آسان تر است. داده های فرمت شده و تأیید شده به درستی کیفیت داده ها را بهبود می بخشد و از برنامه ها در برابر مشکلات بالقوه مانند مقادیر null، duplicate های غیرمنتظره، ایندکس نادرست و فرمت های ناسازگار محافظت می کند.



تبدیل log transform در مهندسی ویژگی چیست؟

Log Transform یکی از محبوب ترین تکنیک های تبدیل است. log Transform نیز تأکید بر نقاط outlier را کاهش می دهد و به ما امکان می دهد به طور بالقوه یک توزیع زنگ شکل به دست آوریم. ایده این است که گرفتن گزارش داده ها می تواند تقارن را به داده ها بازگرداند. تبدیل log همیشه برای تجزیه و تحلیل داده ها ضروری نیست. این می تواند به تجزیه و تحلیل آماری که ما انجام می دهیم بستگی داشته باشد.



چه زمانی باید از تبدیل log استفاده کنیم؟

تبدیل log را می توان برای تبدیل توزیع های highly skewed به less skewed استفاده کرد. این می تواند هم برای تفسیرپذیرتر کردن الگوهای موجود در داده ها و هم برای کمک به برآورده کردن مفروضات آمار استنباطی ارزشمند باشد.

در درجه اول برای تبدیل یک توزیع اریب^۱ به توزیع معمولی^۲/توزیع کم انحراف^۳ استفاده می شود. در این تبدیل، مقادیر یک ستون را می گیریم و به جای آن از این مقادیر به عنوان ستون استفاده می کنیم و عملیات log نقش دوگانه دارد :

- کاهش تاثیر مقادیر خیلی کم
- کاهش تاثیر مقادیر بیش از حد بالا

توجه شود که اگر داده های ما مقادیر منفی یا مقادیری از ۰ تا ۱ داشته باشند، نمی توانیم تبدیل log را مستقیماً اعمال کنیم و از آنجایی که گزارش اعداد منفی و اعداد بین ۰ و ۱ تعریف نشده است، خطا یا مقادیر NaN را در داده ها ی خود دریافت می کنیم. در چنین مواردی، می توانیم یک عدد به این مقادیر اضافه کنیم تا همه آنها بزرگتر از ۱ شوند. سپس، می توانیم تبدیل log را اعمال کنیم.

بطور خلاصه، هر تکنیک مقیاس بندی feature، خصوصیات خاص خود را دارد که می توانیم از آنها برای بهبود مدل خود استفاده کنیم. با این حال، درست مانند سایر مراحل در ساخت یک مدل پیش بینی، انتخاب مقیاس کننده مناسب نیز یک فرآیند آزمون و خطا است و بهترین مقیاس کننده ای وجود ندارد که هر بار کار کند.

سوال ۸. Cross validation and bootstrapping

Cross validation و bootstrapping هر دو از روش های نمونه گیری مجدد هستند.

بوت استرپینگ معیاری است که بر نمونه گیری تصادفی با جایگزینی تکیه دارد. بوت استرپ مجدداً با جایگزینی نمونه برداری می کند (و معمولاً مجموعه های داده «جایگزین» جدیدی با تعداد موارد مشابه مجموعه داده اصلی تولید می کند). این روشی است که در بسیاری از موقعیت ها مانند اعتبارسنجی عملکرد یک مدل پیش بینی کننده، روش های ensemble، تخمین bias و واریانس پارامتر یک مدل و غیره کمک می کند. به دلیل ترسیم با جایگزینی، یک مجموعه داده بوت استرپ ممکن است حاوی چندین نمونه از همان موارد اصلی باشد و ممکن است دیگر موارد اصلی را کاملاً حذف کند.

ما می توانیم این روش را چندین بار تکرار کنیم و میانگین امتیاز را به عنوان تخمین عملکرد مدل خود محاسبه کنیم. همچنین، Bootstrapping به روش های ensemble training مربوط می شود، زیرا می توانیم با استفاده از هر مجموعه داده بوت استرپ یک مدل بسازیم و این مدل ها را در یک مجموعه با استفاده از رأی اکثریت (برای طبقه بندی) یا محاسبه میانگین (برای پیش بینی های عددی) برای همه موارد، «bag» کنیم. این مدل ها به عنوان نتیجه نهایی ما هستند.

Cross validation بدون جایگزینی مجدد نمونه برداری می کند و بنابراین دیتاست های جایگزینی تولید می کند که کوچکتر از اصلی هستند. Cross validation روشی برای اعتبارسنجی عملکرد یک مدل است و با تقسیم داده های آموزشی به k قسمت انجام می شود. ما فرض می کنیم که قسمت های k-1 مجموعه آموزشی است و استفاده از قسمت دیگر مجموعه تست ما است. می توانیم هر بار با نگه داشتن بخش متفاوتی از داده ها، این k بار را به طور متفاوت تکرار کنیم. در نهایت، میانگین نمره k را به عنوان تخمین عملکرد خود در نظر می گیریم. اعتبارسنجی متقاطع می تواند از bias یا واریانس تاثیر پذیرد. با افزایش

¹ skewed distribution

² normal distribution

³ less-skewed

تعداد تقسیم ها، واریانس نیز افزایش می یابد و bias کاهش می یابد. از سوی دیگر، اگر تعداد تقسیم ها را کاهش دهیم، bias افزایش و واریانس کاهش می یابد.

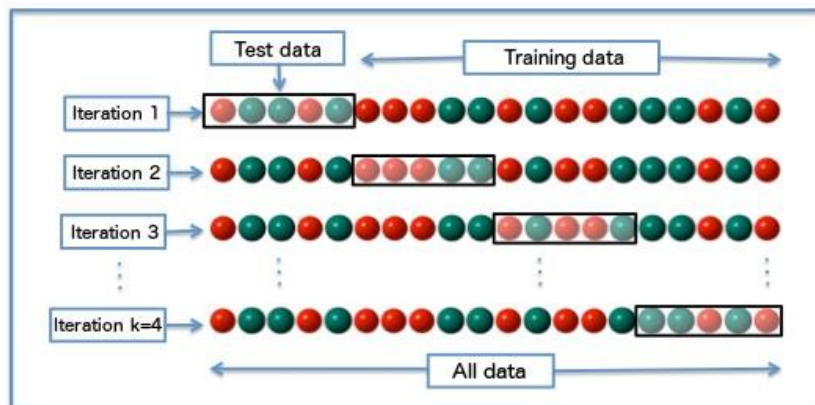
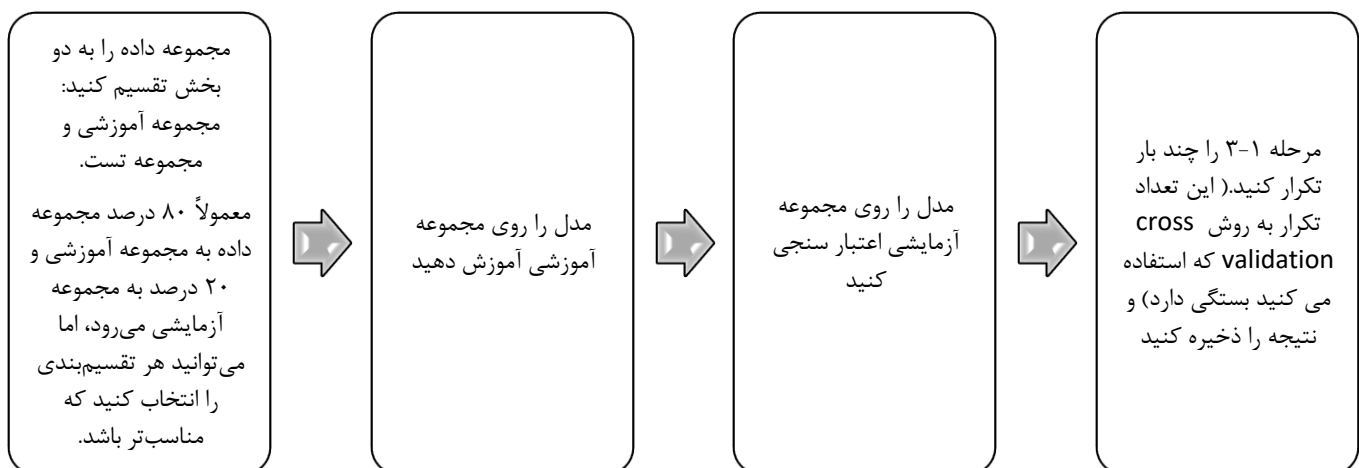


Diagram of k-fold cross-validation with $k=4$.

به طور خلاصه، تکنیک های مختلفی وجود دارد که ممکن است برای Cross validation یک مدل استفاده شود. با این حال، همه آنها یک الگوریتم مشابه دارند:



الگوریتم تکنیک Hold-out cross-validation

Cross validation مجموعه داده موجود را برای ایجاد مجموعه داده های متعدد تقسیم می کند و روش بوت استرپینگ از مجموعه داده اصلی برای ایجاد مجموعه داده های متعدد پس از نمونه برداری مجدد با جایگزینی استفاده می کند. هنگامی که برای اعتبارسنجی مدل استفاده می شود، بوت استرپینگ به اندازه Cross validation قوی نیست. Bootstrapping بیشتر در مورد ساخت ensemble models یا فقط تخمین پارامترها است



ما معمولاً از روش Hold-out در مجموعه داده های بزرگ استفاده می کنیم، زیرا فقط یک بار نیاز به آموزش مدل دارد.

سوال ۹ . 5x2 cross validation

5x2 cross validation عبارت است از ۵ بار تکرار 2-fold و سپس محاسبه میانگین . در سالهای گذشته، در این مقاله^۴ نوشته Dietterich، روش نمونه‌گیری مجدد را به نام اعتبارسنجی متقابل ۲×۵ توصیه می‌کند که شامل ۵ تکرار اعتبارسنجی متقاطع ۲ برابری است. استفاده از آزمون مک نمار^۵ یا اعتبارسنجی متقاطع ۲×۵ در بیشتر ۲۰ سال از زمان انتشار مقاله به توصیه اصلی تبدیل شده است.

در این روش، ابتدا از t تست جفتی ۵ در ۲ استفاده می‌کنیم و این تست برای هر طبقه‌بندی کننده، ۵ بار تکرار 2-fold را اجرا می‌کند (دو تا برای اطمینان از این است که هر مشاهده فقط در مجموعه train یا test برای یک تخمین واحد از مدل انتخاب شود). سپس از تفاوت امتیازات برای محاسبه آمار t از فرمول زیر استفاده می‌شود:

$$t = \frac{p_1^{(1)}}{\sqrt{\frac{1}{5} \sum_{i=1}^5 s_i^2}}$$

where :

- $p_1^{(1)}$ is the classifiers' scores difference for the first fold of the first iteration
- s_i^2 is the estimated variance of the score difference for i^{th} iteration. This variance computes as $(p_i^{(1)} - \bar{p}_i)^2 + (p_i^{(2)} - \bar{p}_i)^2$
- $p_i^{(j)}$ is the classifiers' scores difference for the i^{th} iteration and fold j
- $\bar{p}_i = (p_i^1 + p_i^2) / 2$

از آزمون t-student برای به‌روزرسانی شده نتایج استفاده می‌شود تا درجات محدود آزادی را با توجه به وابستگی بین نمرات مهارت تخمین زده بهتر منعکس کند.

آنچه باید بدانیم این است که بر اساس فرضیه صفر (یعنی هر دو طبقه‌بندی کننده از نظر آماری با هم برابر هستند) فرض می‌شود که تفاوت امتیاز بین دو طبقه‌بندی کننده در هر فولد از توزیع نرمال پیروی می‌کند. با این فرض آمار t از توزیع t با ۵ درجه آزادی پیروی می‌کند.

برای آزمایش فرضیه صفر، مقدار t را محاسبه کرده و بررسی می‌کنیم که آیا توزیع t با ۵ درجه آزادی را برآورده می‌کند. یعنی ما بررسی می‌کنیم که آیا مقدار به نظر یک مقدار پرت است یا نه. اگر مقدار به اندازه کافی به ۰ نزدیک شود، فرضیه صفر برآورده می‌شود و طبقه بندی کننده ها برابر فرض می شوند. توجه شود که بیشتر آزمون‌های آماری برای طبقه‌بندی کننده‌ها از امتیاز دقت استفاده می‌کنند که می‌توان این را به Gini تعمیم داد.

مزیت مهمی که از cross validation تکراری (یا تخمین‌های خارج از بوت استرپ) به دست می‌آید این است که چندین پیش‌بینی برای یک مورد مشابه توسط مدل‌های جایگزین متفاوت دارید. این ها با آموزش در مجموعه های آموزشی کمی متفاوت ، متفاوت هستند. بنابراین، می‌توانید واریانس را به دلیل این تفاوت‌های جزئی در داده‌های train، با توجه به مبادله چند مورد train، اندازه‌گیری کنید. تکرارها به تخمین این واریانس (و کاهش عدم قطعیت واریانس مربوطه در برآورد نهایی) کمک می‌کنند.

⁴ Approximate Statistical Tests For Comparing Supervised Classification Learning Algorithms, Thomas G. Dietterich, department of computer science, Oregon state university, Corvallis, OR 97331, dec 30, 1997

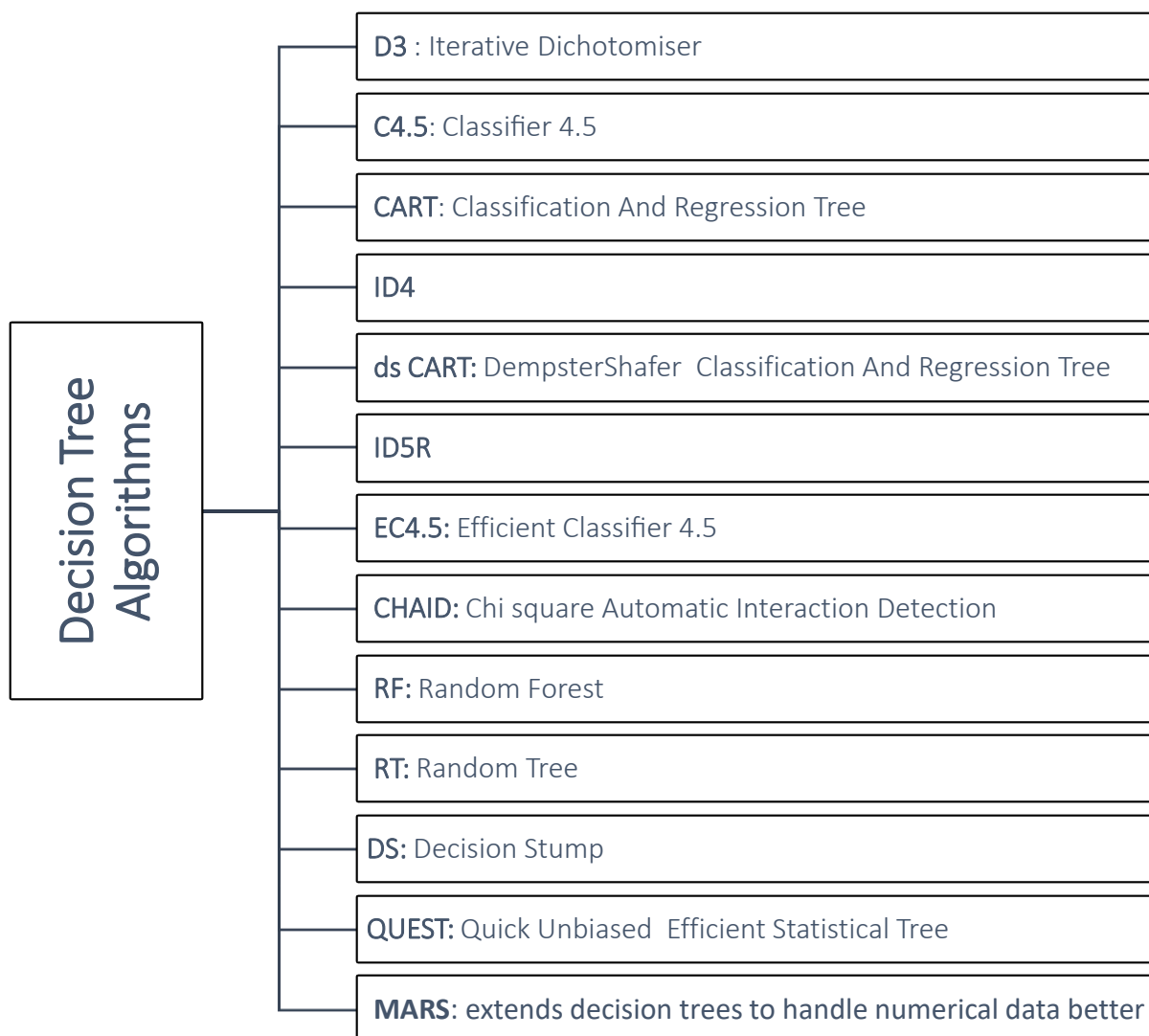
⁵ McNemar's statistical hypothesis test

علاوه بر این، شما انتخاب واقعی موارد آزمایشی را دارید که به واریانس کمک می‌کنند: تخمین‌ها بر اساس تعداد کل کمتر موارد مستقل، قطعیت کمتری دارند. با این نوع واریانس، موارد واقعی بیشتر، کمک کننده است، اما تکرارهای بیشتر هیچ تاثیری بر این منبع عدم قطعیت ندارند. توجه به این نکته ضروری است که این سهم‌های واریانس مستقل از یکدیگر هستند.

سوال ۱۰. الگوریتم‌های متداول برای ساخت درخت تصمیم Decision tree در داده کاوی

Decision Tree یکی از قدرتمندترین و محبوب‌ترین الگوریتم‌ها است. الگوریتم درخت تصمیم در دسته الگوریتم‌های یادگیری نظارت شده قرار می‌گیرد و هم برای متغیرهای خروجی پیوسته و هم برای متغیرهای categorical^۶ کار می‌کند.

درخت تصمیم گره‌ها را روی همه متغیرهای موجود تقسیم می‌کند و سپس تقسیمی را انتخاب می‌کند که منجر به اکثر گره‌های فرعی همگن می‌شود. الگوریتم ID3 درخت‌های تصمیم را با استفاده از رویکرد جستجوی حریصانه از بالا به پایین در فضای شاخه‌های ممکن و بدون پس‌گرد ایجاد می‌کند. الگوریتم‌های درخت تصمیم قابل توجه عبارتند از:



⁶ Supervised

⁷ categorical output variables.

الگوریتم ID3

این الگوریتم یکی از ساده ترین الگوریتم های درخت تصمیم (decision-tree) است. در این الگوریتم درخت تصمیم از بالا به پایین ساخته می شود. این الگوریتم با این سوال شروع می شود: کدام ویژگی باید در ریشه درخت مورد آزمایش، قرار بگیرد؟ برای یافتن جواب از معیار بهره اطلاعات استفاده می شود.

با انتخاب این ویژگی، برای هر یک از مقادیر ممکن آن یک شاخه ایجاد شده و نمونه های آموزشی بر اساس ویژگی هر شاخه مرتب می شوند. سپس عملیات فوق برای نمونه های قرار گرفته در هر شاخه تکرار می - شوند تا بهترین ویژگی برای گره بعدی انتخاب شود.

الگوریتم C4.5

این الگوریتم یکی از تعمیم های الگوریتم ID3 است که از معیار نسبت بهره (Gain ratio) استفاده می - کند. الگوریتم هنگامی متوقف می شود که تعداد نمونه ها کمتر از مقدار مشخص شده ای باشد. این الگوریتم از تکنیک پس هرس استفاده می کند و همانند الگوریتم قبلی داده های عددی را نیز می پذیرد.

از نقاط ضعف الگوریتم ID3 که در C4.5 رفع شده است می توان به موارد زیر اشاره کرد:

الگوریتم C4.5 می تواند مقادیر گسسته یا پیوسته را در ویژگی ها درک کند و الگوریتم C4.5 قادر است با وجود مقادیر گمشده نیز درخت تصمیم (decision tree) خود را بسازد، در حالی که الگوریتمی مانند ID3 و بسیاری دیگر از الگوریتم های طبقه بندی نمی توانند با وجود مقادیر گمشده، مدل خود را بسازند.

سومین موردی که باعث بهینه شدن الگوریتم C4.5 نسبت به ID3 می شود، عملیات هرس کردن جهت جلوگیری از بیش برآزش می باشد. الگوریتم هایی مانند ID3 به خاطر اینکه سعی دارند تا حد امکان شاخه و برگ داشته باشند (تا به نتیجه مورد نظر برسند) با احتمال بالاتری دارای پیچیدگی در ساخت مدل و این پیچیدگی در بسیاری از موارد الگوریتم را دچار بیش برآزش و خطای بالا می کند. اما با عملیات هرس کردن درخت که در الگوریتم ۵ انجام می شود، می توان مدل را به یک نقطه بهینه رساند که زیاد پیچیده نباشد (و البته زیاد هم ساده نباشد) و بیش برآزش یا کم برآزش (Underfitting) رخ ندهد. الگوریتم C4.5 این قابلیت را دارد که وزن های مختلف و غیر یکسانی را به برخی از ویژگی ها بدهد.

الگوریتم CHAID

محققان آمار کاربردی، الگوریتم هایی را جهت تولید و ساخت درخت تصمیم توسعه دادند. الگوریتم CHAID در ابتدا برای متغیرهای اسمی طراحی شده بود. این الگوریتم با توجه به نوع برچسب کلاس از آزمون های مختلف آماری استفاده می کند. این الگوریتم هرگاه به حداکثر عمق تعریف شده ای برسد و یا تعداد نمونه ها در گره جاری از مقدار تعریف شده ای کمتر باشد، متوقف می شود. الگوریتم CHAID هیچگونه روش هرسی را اجرا نمی کند.

الگوریتم درخت تصمیم به این صورت است که یک گره ریشه در بالای آن قرار دارد و برگ های آن در پایین می باشند. یک رکورد در گره ریشه وارد می شود و در این گره یک تست صورت می گیرد تا معلوم شود که این رکورد به کدام یک از گره های فرزند (شاخه پایین تر) خواهد رفت

CHAID. هنگام محاسبه درختان طبقه بندی، تقسیم های چند سطحی را انجام می دهد.

MARS: درختان تصمیم را برای مدیریت بهتر داده های عددی گسترش می دهد.

الگوریتم‌های ساخت درخت‌های تصمیم معمولاً از بالا به پایین کار می‌کنند، با انتخاب متغیری در هر مرحله که به بهترین نحو مجموعه موارد را تقسیم می‌کند. الگوریتم‌های مختلف از معیارهای مختلفی برای اندازه‌گیری "بهترین" استفاده می‌کنند. اینها عموماً همگنی متغیر هدف را در زیر مجموعه‌ها اندازه‌گیری می‌کنند. چند نمونه در زیر آورده شده است. این معیارها برای هر زیر مجموعه نامزد اعمال می‌شود و مقادیر حاصل با هم ترکیب می‌شوند (مثلاً میانگین می‌شوند) تا معیاری از کیفیت تقسیم ارائه شود.

درختان استنتاج شرطی^۸ رویکرد مبتنی بر آمار که از آزمون‌های ناپارامتریک به عنوان معیارهای تقسیم استفاده می‌کند، برای جلوگیری از برازش بیش از حد برای آزمایش‌های متعدد تصحیح شده است. این رویکرد منجر به انتخاب بی‌طرفانه پیش‌بینی‌کننده می‌شود و نیازی به هرس ندارد.

ناخالصی جینی Gini impurity

ناخالصی جینی، شاخص تنوع جینی، یا شاخص جینی سیمپسون در تحقیقات تنوع زیستی، توسط الگوریتم CART (درخت طبقه بندی و رگرسیون) برای طبقه بندی درختان استفاده می‌شود، ناخالصی جینی (به نام ریاضیدان ایتالیایی کورادو جینی) معیاری برای سنجش است. اگر به طور تصادفی بر اساس توزیع برچسب‌ها در زیرمجموعه برچسب‌گذاری شود، هر چند وقت یکبار یک عنصر به طور تصادفی از مجموعه به اشتباه برچسب‌گذاری می‌شود.

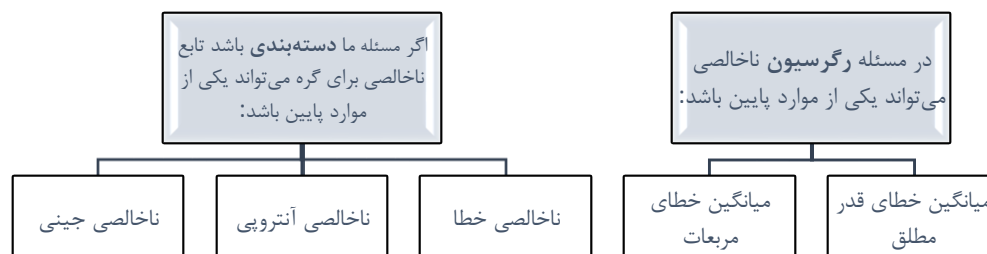
ناخالصی جینی را می‌توان با جمع کردن احتمال p_i یک آیتم با برچسب i اشتباه در برابر احتمال $1 - p_i$ محاسبه کرد و زمانی که همه موارد در گره در یک دسته هدف قرار می‌گیرند به حداقل خود (صفر) می‌رسد.

Information gain

Information gain در درخت تصمیم توسط الگوریتم‌های تولید درخت ID3، C4.5 و C5.0 استفاده شده. کسب اطلاعات مبتنی بر مفهوم "آنتروپی و محتوای اطلاعاتی" از نظریه اطلاعات^۹ است. آنتروپی به صورت زیر تعریف می‌شود:

$$H(T) = I_E(p_1, p_2, \dots, p_J) = - \sum_{i=1}^J p_i \log_2 p_i$$

به طور دقیق تر در هر گره ما سعی می‌کنیم یک بُعد از متغیرهایی وابسته را به همراه یک آستانه انتخاب کنیم و داده‌ها را برحسب این بُعد و آستانه به دو نیم تقسیم کنیم، به قسمی که بطور متوسط در هر دو نیم متغیرهای مستقل γ خیلی به هم نزدیک و همسان شده باشند. این بُعد و آستانه را می‌نامیم. حال سؤال اینجاست که کدام بُعد از متغیرهای وابسته و چه آستانه‌ای را باید انتخاب کرد. به زبان ریاضی باید آن تتایی را انتخاب کرد که ناخالصی داده را کم کند.



⁸ Conditional Inference Trees

⁹ Information theory

سوال ۱۳. هرس درخت تصمیم : Prepruning و Postpruning

در درخت تصمیم مشکلی با نام overfitting مطرح است . ممکن است در ایجاد درخت تعداد زیادی شاخه به وجود آید که دلیل آن وجود آنومالی در داده‌ها است . آنومالی به دلیل وجود نویز و داده‌های پرت به وجود آید از طرفی ممکن است درخت ایجاد شده برای داده‌های جدید ضعیف عمل کند

راه حل: هرس نمودن شاخه‌های زائد است و اینکار باعث کوچک شدن، ساده شدن و به الطبع فهم آسان درخت خواهد شد. از جهتی برای داده‌های تست نیز عملکرد بهتری خواهد داشت . برای هرس دو روش وجود دارد:

۱. Prepruning

۲. Postpruning

Prepruning

- این روش در زمان ساخت درخت اعمال می‌شود
- درواقع در صورتی که تشخیص داده شود که تقسیم بیشتر یک شاخه بهبودی در دقت نخواهد داشت آن شاخه هرس شده و دیگر تقسیم انجام نمی‌گیرد
- کلاس آن شاخه براساس بیشترین تعداد اعضای آن محاسبه می‌شود
- مثلاً در صورتی که بیشتر داده‌های باقی‌مانده در آن شاخه **yes** باشند، شاخه را دیگر تقسیم نکرده و یک برگ با مقدار برچسب **yes** قرار می‌دهیم
- تقسیم شدن یا نشدن براساس تعیین یک آستانه مناسب است
- تعیین این آستانه ساده نیست

Postpruning

- روش متداول تر است
- زمانی که درخت تصمیم ایجاد شد اعمال می‌شود
- در صورتی که زیر درخت یک گره، هرس شود آن گره تبدیل به برگ می‌شود
- کلاس این برگ بیشترین درصد کلاس‌های آن زیر درخت خواهد شد
- این روش از روش قبلی هزینه برتر است

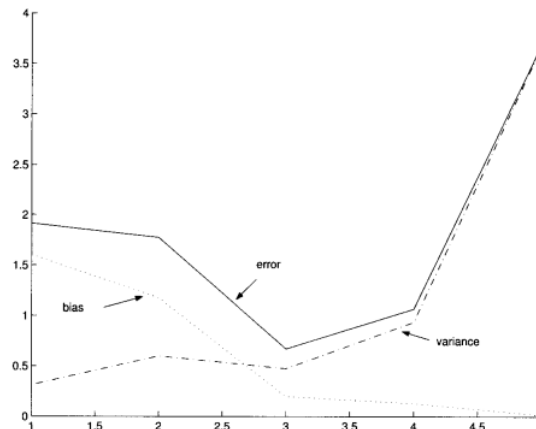
بعنوان نمونه :

- بطور مثال در درخت CART از معیار **cast complexity** برای پس هرس استفاده می‌شود
- این معیار توسط دو عامل تعیین می‌شود
 - تعیین برگ‌های درخت (number of leaves)
 - نسبت خطای درخت (error rate)
- هرس از پایین درخت شروع می‌شود و براساس این معیار هرس شدن یا نشدن یک شاخه تعیین می‌شود
- برای تعیین خطای درخت از مجموعه **pruning set** برای تعیین درصد خطا استفاده میشود

توجه داشته باشید که هرس از پایین به بالا و چپ به راست است.

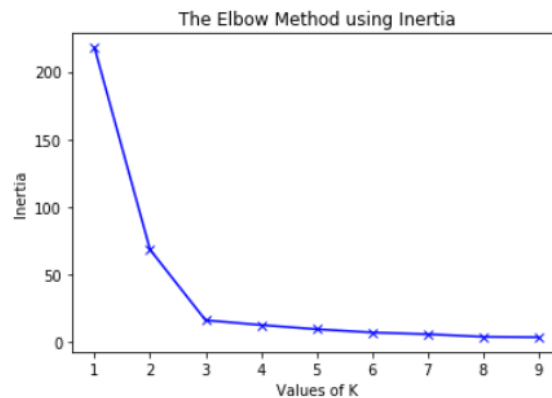
سوال ۱۴. The elbow method

CSS^{۱۰} مجموع مجذور فاصله بین هر نقطه و مرکز در یک خوشه است. هنگامی که $WCSS^{11}$ را با مقدار K رسم می کنیم، نمودار شبیه یک Elbow به نظر می رسد. با افزایش تعداد خوشه ها، مقدار $WCSS$ شروع به کاهش می کند. برای تعیین تعداد بهینه خوشه ها، باید مقدار k را در "زانو (elbow)" انتخاب کنیم، یعنی نقطه ای که پس از آن اعوجاج/انحراف (distortion/inertia) به صورت خطی شروع به کاهش می کند.



شکل ۱. نمودار بایاس و واریانس بنا بر مرتبه های مختلف مدل

بنابراین مثلاً برای داده های داده شده در شکل زیر، نتیجه می گیریم که تعداد بهینه خوشه برای داده ها ۳ است.



خطای بایاس: وجود فرضیه های مختلف روی مدل و الگوریتم یادگیری منجر به ایجاد خطای بایاس (اریبی) می شود. بزرگ بودن اریبی می تواند الگوریتم یا مدل آماری را از کشف روابط بین ویژگی ها (Features) و متغیر پاسخ (Target Variable) باز دارد. اغلب بزرگ بودن خطای اریبی، منجر به کم برآزش (Underfitting) می شود.

خطای واریانس: حساسیت زیاد مدل با تغییرات کوچک روی داده های آموزشی، نشانگر وجود واریانس زیاد است. این امر نشانگر آن است که اگر مدل آموزش داده شده را روی داده های آزمایشی به کارگیریم، نتایج حاصل با داده های واقعی فاصله زیادی خواهند داشت. متأسفانه افزایش واریانس در این حالت منجر به مدل بندی مقادیر Noise شده و به جای پیش بینی صحیح، دچار پیچیدگی و مشکل «بیش برآزش» (Overfitting) می شود.

¹⁰ Cluster-Sum of Squared

¹¹ Within-Cluster-Sum of Squared Errors (WCSS)

مدل‌های با واریانس بزرگ (مثلاً رگرسیون چند جمله‌ای هم‌مرتبه با تعداد مشاهدات)، که معمولاً پیچیده‌تر هستند، این امکان را می‌دهد تا داده‌های train به خوبی برازش شوند. با این وجود، ممکن است مشاهدات برازش شده دارای خطا یا نویز باشند که متأسفانه مدل تحت تاثیر آن‌ها، برآوردها را با دقت انجام داده است. به این ترتیب پیش‌بینی آن‌ها باعث افزودن پیچیدگی در مدل شده است. در حالیکه این امر از طرفی دقت برآوردها را هم برای داده‌های آزمایشی کمتر می‌کند. در مقابل، مدل‌هایی که دارای بایاس بزرگی هستند، نسبتاً ساده بوده (مثل مدل رگرسیون دو جمله‌ای یا حتی خطی) اما ممکن است واریانس کوچکتري را براساس مجموعه داده‌های آزمایشی ایجاد کنند.



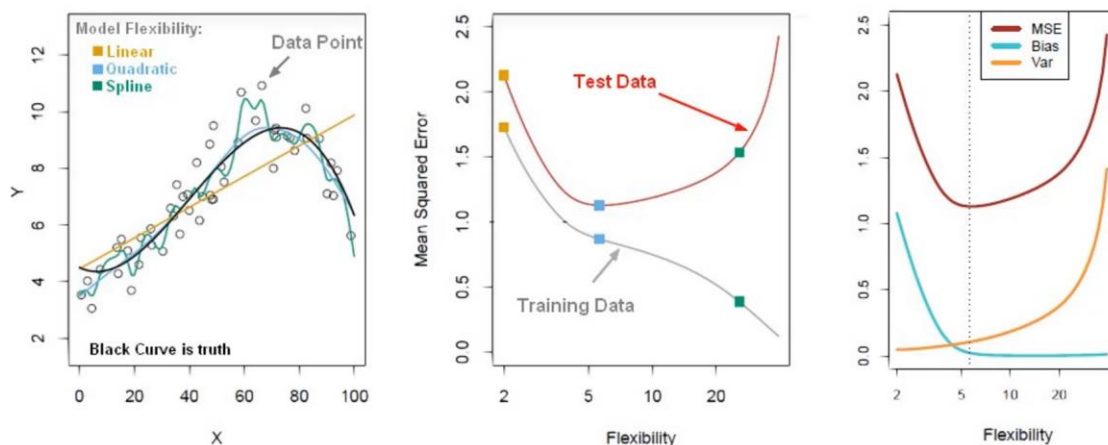
استفاده از روش‌های کاهش بُعد (Dimensionality Reduction) و انتخاب ویژگی (Feature Selection) می‌توانند واریانس را به کمک ساده‌سازی مدل انجام دهند. از طرفی افزایش داده‌های آموزشی باعث کاهش واریانس خواهد شد. افزایش تعداد متغیرهای پیش‌گو (Predictors)، باعث کاهش بایاس (اریبی) می‌شود ولی این امر به قیمت افزایش واریانس خواهد بود. الگوریتم‌های یادگیری، معمولاً از یک پارامتر تنظیم‌کننده برای موازنه واریانس و بایاس استفاده می‌کنند. از آنجایی که همه منابع خطا (مربع اریبی، واریانس مدل یا واریانس عبارت خطا) مثبت یا حداقل نامنفی هستند، می‌توان یک کران پایین برای خطای مدل روی داده‌های آزمایشی در نظر گرفت. هر چه مدل پیچیده‌تر باشد، نقاط بیشتری از داده‌های آموزشی را پوشش می‌دهد و بایاس نیز کم خواهد بود. در حالیکه یک مدل پیچیده، باعث بوجود آمدن خطای زیاد برای برازش داده‌های جدید خواهد شد در نتیجه واریانس مدل را برای چنین داده‌هایی، زیاد می‌کند.

پس برای ساختن یک مدل خوب، باید تعادل خوبی بین واریانس و بایاس پیدا کنیم به طوری که خطای کل را به حداقل برساند. به همین دلیل خطای مدل را به واریانس و بایاس (اریبی) تجزیه می‌کنیم. در نهایت، تابع زیان MSE (میانگین مربعات خطا) یا لگاریتم درستنمایی منفی به کمک محاسبه امید ریاضی و طبق رابطه زیر حاصل می‌شود (خطای کل):

$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

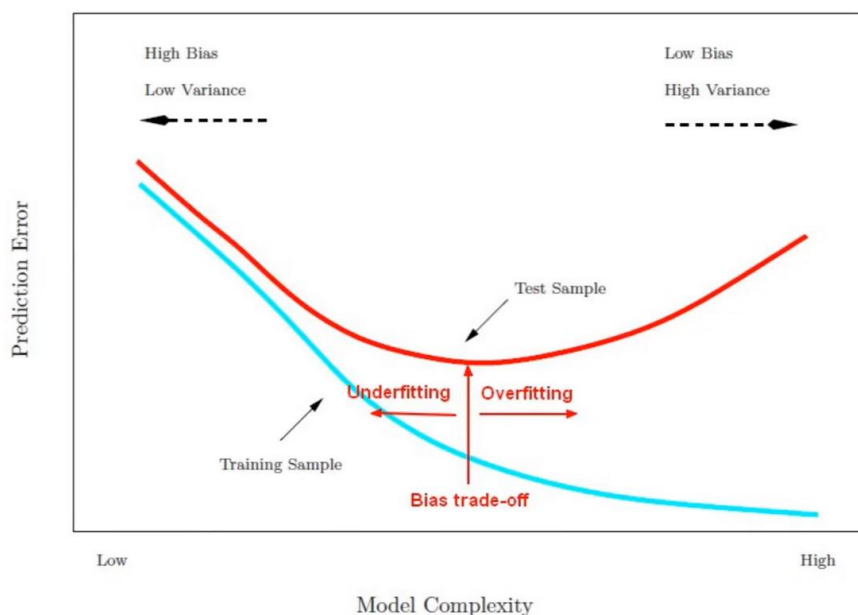
$$\Rightarrow \text{MSE} = E_x \{ \text{Bias}_D [f^{\wedge}(x; D)]^2 + \text{Var}_D [f^{\wedge}(x; D)] \} + \sigma^2$$

تعادل بهینه بین بایاس و واریانس، باعث می‌شود مدل overfitt و underfitt نشود.



در تصویر فوق، نمودار سمت چپ در شکل بالا یک انعطاف پذیری (flexibility) مدل را به صورت خطی، درجه دوم و برازش اسپلاین توصیف می کند. برای مقایسه مدل هایمان با یکدیگر، باید انعطاف پذیری در مقابل میانگین مربعات خطا را همانطور که در نمودار میانی بالا توضیح داده شد، ترسیم می کنیم.

نمودار داده های آموزشی را در مقابل داده های آزمون ترسیم می کند. می بینیم که برای مدل خطی (به رنگ زرد) خطای بالایی داریم. خطا برای مدل درجه دوم (به رنگ آبی) کاهش می یابد، ستاره خطا دوباره برای مدل spline (به رنگ سبز) افزایش می یابد، اما به طور اساسی فقط در مجموعه تست، در عین حال برای مجموعه آموزشی، خطا همچنان کاهش می یابد: این شرط **overfitting** نامیده می شود.



در شکل بالا، خطای پیش بینی را به عنوان تابعی از پیچیدگی مدل (model complexity) نشان می دهد. همانطور که به سمت چپ حرکت می کنیم، بایاس زیاد اما واریانس کم = **underfit** داده ها را داریم. همانطور که به سمت راست حرکت می کنیم، جایی که پیچیدگی بالاتری وجود دارد، یک بایاس کمتر اما واریانس بالا = **overfit** داده ها را دریافت می کنیم.

در اغلب «مدل‌های پیش‌بین (Predictive Model)» وجود بایاس کوچک برای پارامترها منجر به واریانس بزرگ برای مدل خواهد شد. البته برعکس این حالت نیز وجود دارد، به این معنی که با کوچک کردن واریانس مدل، با مشکل بزرگ شدن بایاس یا اریبی پارامترها مواجه خواهیم شد.

مسئله اصلی آن است که در یک مدل مناسب، هم بایاس و هم واریانس باید حداقل ممکن باشند. ولی متأسفانه، کمینه‌سازی (Minimization) هر دو این شاخص‌ها به شکل توأم، امکان‌پذیر نیست. چنین وضعیتی را «تناقض واریانس-اریبی» (Bias-Variance Dilemma) می‌نامند.

تجزیه میزان خطای کل به واریانس و بایاس، همچنین بهره‌گیری از کمینه‌سازی آن‌ها در مدل رگرسیونی، ایده اصلی در روش‌های رگرسیون با قاعده‌سازی (Regularization) نظیر «رگرسیون لاسو» (Lasso Regression) و رگرسیون Ridge Regression است. به این ترتیب با استفاده از تکنیک‌هایی مانند جریمه کردن مدل، مثلاً در «رگرسیون لاسو» (Lasso Regression) یا رگرسیون Ridge Regression موازنه واریانس و بایاس برقرار می‌شود و مدلی بدون underfitting یا overfitting بدست می‌آید.

روش‌های قاعده‌سازی، بایاس را وارد روش حل و برآورد پارامترهای مدل کرده و به این ترتیب واریانس مدل را نسبت به روش‌های رگرسیون عادی (Ordinary Least Squares) یا OLS کاهش می‌دهند. هر چند روش‌های مبتنی بر OLS، برآوردگرهای نارایب برای پارامترهای رگرسیونی ایجاد می‌کنند ولی وجود واریانس کوچکتر در مدل‌های رگرسیونی با قاعده، مفیدتر بوده و مقادیر برآورد شده توسط آن‌ها، دارای خطای کمتری هستند.

سوال ۲ تسک امتیازی: چگونه می‌توان با استفاده از statistical significance tests به مقایسه مدل‌ها پرداخت؟

از آزمون‌های معنی‌داری آماری برای پاسخ به این سؤال استفاده می‌شود: احتمال اینکه آنچه ما فکر می‌کنیم رابطه بین دو متغیر است واقعاً فقط یک اتفاق تصادفی باشد چقدر است؟ اگر نمونه‌های زیادی را از یک جامعه انتخاب کنیم، آیا باز هم رابطه یکسانی بین این دو متغیر در هر نمونه پیدا می‌کنیم؟ اگر بتوانیم یک سرشماری از جامعه انجام دهیم، آیا این رابطه در جمعیتی که نمونه از آن گرفته شده است نیز وجود دارد؟ یا اینکه یافته ما فقط به دلیل شانس تصادفی است؟

آزمون‌های معنی‌داری آماری به ما می‌گویند که احتمال اینکه رابطه‌ای که فکر می‌کنیم پیدا کرده‌ایم فقط به دلیل شانس تصادفی باشد چقدر است. آنها به ما می‌گویند که اگر فرض کنیم که رابطه‌ای وجود دارد، احتمال اشتباه ما چقدر است. ما هرگز نمی‌توانیم به طور کامل ۱۰۰٪ مطمئن باشیم که رابطه‌ای بین دو متغیر وجود دارد. منابع خطا بیش از حد قابل کنترل هستند، به عنوان مثال، خطای نمونه‌گیری، سوگیری محقق، مشکلات مربوط به قابلیت اطمینان و اعتبار، اشتباهات ساده و غیره. اما با استفاده از تئوری احتمال و منحنی نرمال، می‌توانیم احتمال اشتباه بودن را تخمین بزنیم، اگر فرض کنیم که یافتن یک رابطه درست است. اگر احتمال اشتباه اندک باشد، می‌گوییم که مشاهده ما از رابطه یک یافته آماری معنی‌دار است.

معنی‌داری آماری به این معناست که شانس خوبی وجود دارد که در یافتن رابطه‌ای بین دو متغیر درست باشیم. اما اهمیت آماری با اهمیت عملی یکسان نیست. می‌توانیم یافته‌های آماری معنی‌داری داشته باشیم، اما پیامدهای آن یافته ممکن است کاربرد عملی نداشته باشد. محقق باید همیشه هم اهمیت آماری و هم اهمیت عملی هر یافته تحقیق را بررسی کند.

مراحل آزمون اهمیت آماری

- فرضیه تحقیق را بیان کنید.
- فرضیه صفر را بیان کنید.
- انتخاب سطح احتمال خطا (سطح آلفا)
- آزمون را برای اهمیت آماری انتخاب و محاسبه کنید.
- نتایج را تفسیر کنید

➤ آزمون فرضیه‌های آماری برای مقایسه عملکرد دو مدل

رایج‌ترین آزمون فرضیه‌های آماری مورد استفاده برای مقایسه عملکرد مدل‌های ML، آزمون t زوجی (paired Student's t-test) است که از طریق نمونه‌های فرعی تصادفی از مجموعه داده آموزشی (random subsamples of the training dataset) ترکیب شده است. فرضیه صفر در این آزمون این است که بین عملکرد دو مدل کاربردی ML تفاوتی وجود ندارد. آزمون‌های دیگری نیز وجود دارند مثلاً آزمون کای دو، F-Test و غیره.

➤ تست chi-square

برای داده‌های اسمی و ترتیبی از Chi-Square به عنوان آزمون معناداری آماری استفاده می‌شود. برای محاسبه مربع کای، جدولی که توزیع مشترک دو متغیر را نشان می‌دهد مورد نیاز است. سپس در این جدول با مراحل زیر مربع کای را محاسبه می‌کنیم:

مراحل محاسبه chi-square

- نمایش فرکانس‌های مشاهده شده برای هر سلول
- فرکانس‌های مورد انتظار برای هر سلول را محاسبه کنید
- برای هر سلول، فرکانس مورد انتظار منهای مشاهده شده را به مجذور تقسیم بر فرکانس مورد انتظار محاسبه کنید
- تمام نتایج برای تمام سلول‌ها

ما نمی‌توانیم ارزش آمار Chi-Square را به تنهایی تفسیر کنیم. در عوض، ما باید آن را در یک زمینه قرار دهیم. در تئوری، مقدار آمار Chi-Square به طور معمول توزیع می‌شود. یعنی مقدار آمار Chi-Square شبیه یک منحنی معمولی (زنگ شکل) است. بنابراین می‌توانیم از خواص منحنی نرمال برای تفسیر مقدار بدست آمده از محاسبه آمار Chi-Square استفاده کنیم.

اگر مقداری که برای Chi Square به دست می‌آوریم به اندازه کافی بزرگ باشد، می‌توان گفت که سطح معنی‌داری آماری را نشان می‌دهد که در آن می‌توان رابطه بین دو متغیر را فرض کرد. با این حال، اینکه آیا مقدار به اندازه کافی بزرگ باشد به دو چیز بستگی دارد: اندازه جدول احتمالی که آمار Chi-Square از آن محاسبه شده است. و سطح آلفای که انتخاب کرده ایم.

هرچه اندازه جدول بزرگتر باشد، در صورت مساوی بودن سایر موارد، باید مقدار Chi-Square بزرگتر باشد تا به اهمیت آماری برسد. به طور مشابه، هرچه سطح آلفا دقیق‌تر باشد، اگر سایر موارد برابر باشند، باید مقدار Chi-Square بزرگتر باشد تا به معنی‌داری آماری برسد.

اصطلاح "درجات آزادی" برای اشاره به اندازه جدول احتمالی استفاده می‌شود که مقدار آماره Chi-Square بر روی آن محاسبه شده است. درجه آزادی به صورت حاصل ضرب (تعداد سطرهای جدول منهای ۱) برابر (تعداد ستون‌های جدول منهای) محاسبه می‌شود.

هر دو آزمون t و آزمون کای دو آزمون‌های آماری هستند که برای آزمایش و احتمالاً رد یک فرضیه صفر طراحی شده‌اند. فرضیه صفر معمولاً عبارتی است مبنی بر اینکه چیزی صفر است یا چیزی وجود ندارد. به عنوان مثال، می‌توانید این فرضیه را که تفاوت بین دو میانگین صفر است، یا می‌توانید این فرضیه را آزمایش کنید که بین دو متغیر رابطه وجود ندارد. تفاوت بین آزمون کای دو و آزمون t این است که آزمون t به شما این امکان را می‌دهد که بگویید یا "ما می‌توانیم فرضیه صفر میانگین‌های برابر را در سطح ۰,۰۵ رد کنیم" یا "شواهد کافی برای رد صفر میانگین‌های برابر در سطح ۰,۰۵ نداریم." آزمون کای اسکور به شما امکان می‌دهد بگویید یا "ما می‌توانیم فرضیه صفر عدم وجود رابطه را در سطح ۰,۰۵ رد کنیم" یا "شواهد کافی برای رد عدد صفر در سطح ۰,۰۵ نداریم."

➤ آزمایش اهمیت یک مدل با F-TEST :

برای محاسبه F-test اهمیت کلی (overall significance) ، نرم افزار آماری شما فقط باید عبارت‌های مناسب را در دو مدلی که با هم مقایسه می‌کند، قرار دهد. آزمون F-test کلی ، مدلی که شما مشخص کرده اید را با مدلی که هیچ متغیر مستقلی ندارد مقایسه می‌کند. این نوع مدل به عنوان یک مدل فقط رهگیری (intercept-only model) نیز شناخته می‌شود.

هنگام برازش داده‌ها با استفاده از رگرسیون غیرخطی، اغلب مواقعی پیش می‌آید که باید بین دو مدل انتخاب کرد که هر دو به خوبی با داده‌ها مطابقت دارند. پس از ترسیم باقیمانده‌های هر مدل و مشاهده مقادیر r^2 برای هر مدل، هر دو مدل ممکن است متناسب با داده‌ها به نظر برسند.

در این حالت می‌توان آزمون F انجام داد تا ببینیم کدام مدل از نظر آماری بهتر است. این آزمون F پاسخ قطعی می‌دهد و بر تفسیر دلخواه مقدار r^2 یا نمودار باقی مانده تکیه نمی‌کند.

یک آزمون F از یک توزیع F پیروی می‌کند و می‌تواند برای مقایسه مدل‌های آماری استفاده شود. آماره F با استفاده از یکی از دو معادله بسته به تعداد پارامترها در مدل‌ها محاسبه می‌شود:

➤ اگر هر دو مدل تعداد پارامترهای یکسانی داشته باشند، فرمول آماره F به صورت $F = \frac{SS1}{SS2}$ است که در آن SS1 مجموع مربعات باقیمانده برای مدل اول و SS2 مجموع مربعات باقیمانده برای مدل دوم است.

درجات آزادی $N - V$ وجود دارد، که در آن N تعداد نقاط داده و V تعداد پارامترهای تخمین زده شده است (در هر پارامتر تخمین زده شده یک درجه آزادی از دست می‌رود). سپس آماره F حاصل را می‌توان با یک جدول F برای استخراج مقدار p-value مقایسه کرد.

- اگر مقدار p-value بزرگ باشد (بزرگتر از α) مدل اول از نظر آماری بهتر از مدل دوم است.
- اگر مقدار p-value کوچک باشد (کمتر از $1-\alpha$) مدل دوم از نظر آماری بهتر از مدل اول است.

➤ اگر مدل ها تعداد پارامترهای متفاوتی داشته باشند، فرمول به صورت زیر در می آید:

$$F = \frac{(SS_1 - SS_2)/(df_1 - df_2)}{SS_2/df_2}$$

مجموع مربعات هر مدل و درجات آزادی برای هر مدل مانند قبل محاسبه می شود (توجه داشته باشید که مدل ها برای این مورد درجات آزادی متفاوتی خواهند داشت).

علاوه بر این، اولین مدل باید مدلی با پارامترهای کمتر باشد (یعنی مدل ساده تر). یک بار دیگر می توان از آماره F و درجات آزادی برای تعیین مقدار p استفاده کرد. هنگام یافتن مقدار p از درجه آزادی $df_1 - df_2$ و df_2 استفاده کنید. در این مورد، مقدار p کمتر از α نشان می دهد که مدل پیچیده تر (مخرج آماره F) به طور قابل توجهی بهتر از مدل ساده تر با داده ها مطابقت دارد.

🔵 توجه داشته باشید که آزمون F چیزی در مورد اهمیت فیزیکی مدل به شما نمی گوید و موردی ندارد که مدلی را انتخاب کنید که علیرغم عدم همخوانی مناسب با داده ها، معنای فیزیکی بهتری داشته باشد.

هنگام مقایسه میانگین دو گروه از آزمون تی استفاده می شود هنگام مقایسه میانگین بیش از دو گروه از آزمون های ANOVA و MANOVA استفاده می شود.

تحلیل واریانس (ANOVA)

برای مقایسه برآزش های دو مدل، می توانید از تابع `anova()` با اشیاء رگرسیون به عنوان دو آرگومان جداگانه استفاده کنید. تابع `anova()` اشیاء مدل را به عنوان آرگومان می گیرد و یک ANOVA را برمی گرداند و آزمایش می کند که آیا مدل پیچیده تر به طور قابل توجهی در گرفتن داده ها بهتر از مدل ساده تر است یا خیر. همچنین، آزمون ANOVA امکان مقایسه بیش از دو گروه را به طور همزمان دارد برای تعیین اینکه آیا رابطه ای بین آنها وجود دارد یا خیر.

فرمول ANOVA به شکل زیر است:

$$F = \frac{MST}{MSE}$$

که در آن :

F=ANOVA coefficient

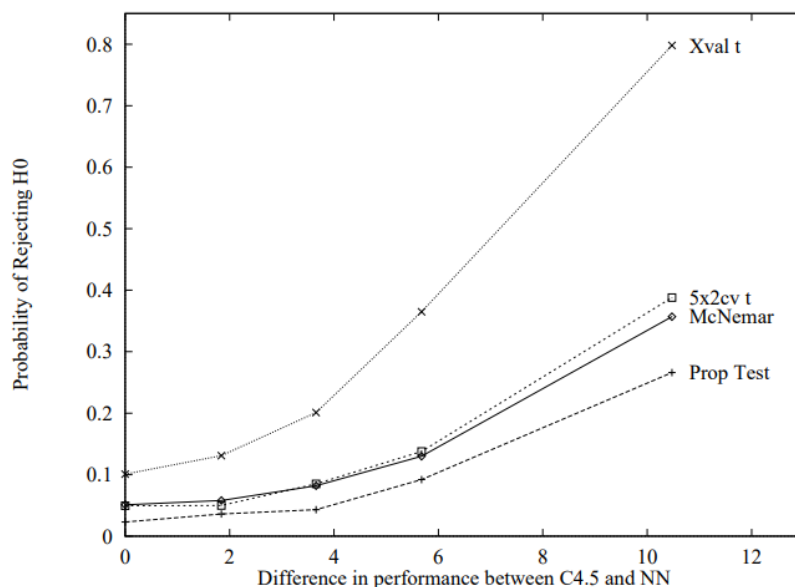
MST=Mean sum of squares due to treatment

MSE=Mean sum of squares due to error

در ANOVA، فرضیه صفر این است که تفاوتی بین میانگین‌های گروهی وجود ندارد. اگر هر گروهی به طور قابل توجهی با میانگین کلی گروه تفاوت داشته باشد، آنگاه ANOVA یک نتیجه آماری معنی دار را گزارش خواهد کرد. تفاوت بین ANOVA و chi-square این است که از chi-square برای بررسی اینکه آیا توزیع طبقات و با یک مدل توزیع سازگار است (اغلب توزیع برابر، اما نه همیشه) استفاده می شود، در حالی که از ANOVA برای بررسی اینکه آیا تفاوت در میانگین بین نمونه ها معنی دار است یا خیر استفاده می شود.

آیا می توانیم از R-Squared برای مقایسه مدل ها استفاده کنیم؟

از R-Squared برای مقایسه مدل ها استفاده نکنید. دلایل متفاوتی برای این وجود دارد مثلاً در بسیاری از موقعیت‌ها، R-Squared در مقایسه با مدل‌ها گمراه‌کننده است. به طور کلی، r-squared بالاتر نشان می دهد که تنوع بیشتر توسط مدل توضیح داده شده است. با این حال، همیشه اینطور نیست که r-squared بالا برای مدل رگرسیون خوب باشد. R-squared نشان نمی دهد که آیا یک مدل رگرسیون کافی است یا خیر. شما می توانید یک مقدار R-squared پایین برای یک مدل خوب، یا یک مقدار R-squared بالا برای مدلی که با داده ها مطابقت ندارد، داشته باشید! R-squared در خروجی شما یک برآورد مغرضانه از جمعیت R-squared است.



تصویر فوق، قدرت چهار آزمون آماری در تسک تشخیص حروف. محور افقی تعداد نقاط درصدی را ترسیم می کند که دو الگوریتم C4.5 و NN در هنگام آموزش روی training set با اندازه ۳۰۰ تفاوت دارند.

به طور خلاصه، آزمون‌های معنی‌داری آماری برای تخمین احتمال اینکه رابطه مشاهده‌شده در داده‌ها تنها به صورت تصادفی رخ داده است، استفاده می‌شود. احتمال اینکه متغیرها واقعاً در جامعه نامرتبط باشند. می توان از آنها برای فیلتر کردن فرضیه های غیرمنتظره استفاده کرد.

آزمون‌های اهمیت آماری به این دلیل استفاده می‌شوند که معیار مشترکی را تشکیل می‌دهند که برای بسیاری از افراد قابل درک است، و اطلاعات ضروری را در مورد یک پروژه تحقیقاتی که می‌تواند با یافته‌های پروژه‌های دیگر مقایسه شود، به اشتراک می‌گذارد.

با این حال، آنها اطمینان نمی دهند که این تحقیق با دقت طراحی و اجرا شده است. در واقع، آزمون‌های اهمیت آماری ممکن است گمراه‌کننده باشند، زیرا اعداد دقیقی هستند. اما هیچ ارتباطی با اهمیت عملی یافته‌های تحقیق ندارند. در نهایت، همیشه باید از معیارهای ارتباط^{۱۲} همراه با آزمون‌ها برای معناداری آماری استفاده کرد. دومی احتمال وجود رابطه را تخمین می زند. در حالی که اولی قدرت (و گاهی اوقات جهت) رابطه را تخمین می زند. هر کدام کاربرد خود را دارند و زمانی که با هم استفاده می شوند بهترین هستند.

نتیجه

- آزمون‌های فرضیه‌های آماری می‌توانند به مقایسه مدل‌های یادگیری ماشین و انتخاب مدل نهایی کمک کنند.
- استفاده ساده از آزمون‌های فرضیه‌های آماری می‌تواند منجر به نتایج گمراه‌کننده شود.
- استفاده صحیح از آزمون‌های آماری چالش برانگیز است، و برای استفاده از آزمون مک‌نمار یا اعتبارسنجی متقابل ۲×۵ با آزمون زوجی اصلاح‌شده توافق نظر وجود دارد.

سوال ۳ تسک امتیازی : معیار **Matthews Correlation Coefficient (MCC)** چیست و کجا استفاده می شود؟

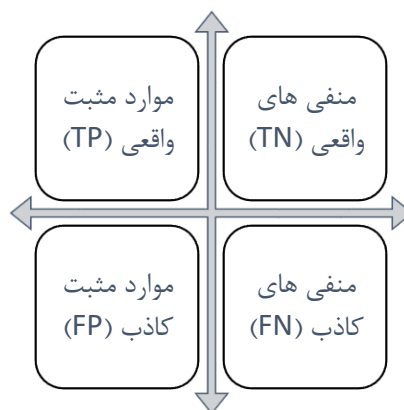
ضریب همبستگی متیوز چیست؟

ضریب همبستگی متیوز که به اختصار MCC نیز نامیده می شود توسط برایان متیوز در سال ۱۹۷۵ اختراع شد. MCC یک ابزار آماری است که برای ارزیابی مدل استفاده می شود. وظیفه آن اندازه‌گیری یا اندازه‌گیری تفاوت بین مقادیر پیش‌بینی شده و مقادیر واقعی است و معادل آمار مربع کای (chi-square statistics) برای یک جدول احتمالی ۲×۲ است که در آن n تعداد کل مشاهدات است.

$$|MCC| = \sqrt{\frac{\chi^2}{n}}$$

فرمول ضریب همبستگی متیوز

MCC بهترین معیار طبقه‌بندی تک‌مقداری است که به خلاصه کردن confusion matrix یا ماتریس خطا کمک می کند. یک ماتریس confusion دارای چهار موجودیت است:



¹² measures of association

و با فرمول زیر محاسبه می شود:

$$MCC = \frac{TN \times TP - FN \times FP}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$$

اگر پیش‌بینی نرخ‌های خوبی را برای هر چهار مورد از این نهادها به دست آورد، گفته می‌شود که معیار قابل اعتمادی است که نمرات بالایی ایجاد می‌کند. و برای مطابقت با اکثر ضرایب همبستگی، MCC نیز بین ۱ و -۱ متغیر است:

- ضریب ۱+ یک پیش‌بینی کامل را نشان می‌دهد
- ۰ بهتر از پیش‌بینی تصادفی نیست
- ۱- نشان‌دهنده اختلاف کامل بین پیش‌بینی و مشاهده است
- اگر MCC برابر با ۱-، ۰ یا ۱+ نباشد، نشانگر قابل اعتمادی نیست که نشان دهد یک پیش‌بینی‌کننده چقدر شبیه حدس‌های تصادفی است، زیرا MCC به مجموعه داده‌ها وابسته است.

دقت و امتیاز F1 محاسبه‌شده بر روی confusion matrix ها یکی از محبوب‌ترین معیارهای پذیرفته‌شده در تسک های طبقه‌بندی باینری بوده است (و هنوز هم هستند). با این حال، این معیارهای آماری می‌توانند به طور خطرناکی نتایج بیش از حد خوش بینانه را به ویژه در مجموعه داده های نامتعادل نشان دهند.

در حالی که هیچ راه کاملی برای توصیف ماتریس اشتباه مثبت و منفی درست و غلط با یک عدد وجود ندارد، ضریب همبستگی متیوز به طور کلی به عنوان یکی از بهترین معیارها در نظر گرفته می‌شود. سایر معیارها، مانند نسبت پیش‌بینی‌های صحیح (که دقت نیز نامیده می‌شود)، زمانی که دو کلاس اندازه‌های بسیار متفاوتی دارند، مفید نیستند. به عنوان مثال، اختصاص دادن هر شی به مجموعه بزرگتر به نسبت بالایی از پیش‌بینی های صحیح دست می‌یابد، اما به طور کلی طبقه بندی مفیدی نیست.

در عوض، ضریب همبستگی متیوز (MCC)، نرخ آماری قابل اعتمادتری است که تنها در صورتی امتیاز بالایی ایجاد می‌کند که پیش‌بینی نتایج خوبی در هر چهار دسته confusion matrix (مثبت واقعی، منفی کاذب، منفی درست و نادرست) متناسب با اندازه عناصر مثبت و اندازه عناصر منفی در مجموعه داده بدست آورد.

MCC در ارزیابی طبقه‌بندی‌های باینری امتیاز آموزنده‌تر و درست‌تری نسبت به دقت و امتیاز F1 ایجاد می‌کند، و توصیه می‌شود ضریب همبستگی متیوز باید به دقت و امتیاز F1 در ارزیابی وظایف طبقه‌بندی باینری توسط همه جوامع علمی ترجیح داده شود.