



تحلیل آماری کاربردی

عنوان:

"پیش بینی قیمت و پست بازیکن فوتبال

با استفاده از رگرسیون و طبقه بندی"

گردآورنده:

امیرعلی خطیب

۹۷۲۲۷۴۳۱۴۹

استاد:

جناب آقای دکتر علیرضا شادمان

تاریخ تحویل:

تیر ۱۴۰۱

فهرست مطالب

مقدمه	۱
تشریح داده‌ها	۱
آماده‌سازی داده‌ها	۶
پیاده‌سازی روش‌ها و ارزیابی	۸
مدل رگرسیون خطی	۹
کا نزدیک ترین همسایه	۱۱
مدل درخت رگرسیون	۱۳
جمع‌بندی نهایی برای مدل‌های رگرسیون	۱۵
رگرسیون لجستیک	۱۶
مدل تحلیل تفکیکی خطی	۱۷
مدل تحلیل تفکیکی کوآدراتیک	۱۸
مدل بیز ساده	۱۹
مدل درخت تصمیم	۱۹
جمع‌بندی نهایی برای مدل‌های طبقه‌بندی	۲۴
جمع‌بندی نهایی	۲۵

مقدمه:

در این گزارش ما قصد داریم به بررسی عملکرد الگوریتم‌های یادگیری ماشین^۱ برای داده‌های بازیکنان فوتبال^۲ بپردازیم. در مجموعه داده اشاره شده ما حدود ۱۹۲۳۹ بازیکن داریم که تمام مشخصات آن‌ها نظیر سن، ملیت، تیم باشگاهی، و توانایی‌های فردی آن‌ها ذخیره شده است. در این پروژه ما قصد داریم تو گروه از الگوریتم‌های رگرسیون^۳ و طبقه‌بندی^۴ را با اهداف متفاوت بر روی داده‌ها پیاده‌سازی کنیم تا پیش‌بینی‌های مورد نظرمان را انجام دهیم. به طور کلی ما ۲ هدف را در این پروژه دنبال می‌کنیم:

- (۱) **هدف مسئله رگرسیون:** هدف از پیاده‌سازی الگوریتم‌های رگرسیون، پیش‌بینی ارزش بازیکن در بازار نقل و انتقالات با در نظر گرفتن پارامترهایی نظیر سن، پتانسیل رشد، شهرت و توانایی‌های فردی می‌باشد. با کمک گرفتن از این الگوریتم‌ها باشگاه‌ها می‌توانند مدیریت بهتری روی خرید و فروش خود داشته باشند؛ زیرا می‌توانند تا حدی قیمت بازیکن موردنظر خود را با وارد کردن اطلاعات شخصی و توانایی‌های فردی بازیکن تخمین بزنند.
- (۲) **هدف مسئله طبقه‌بندی:** هدف از پیاده‌سازی الگوریتم‌های طبقه‌بندی، پیش‌بینی پست بازیکن است. به این صورت که با دیدن توانایی‌های فردی بازیکن و اندازه‌گیری آن‌ها، می‌تواند دید خوبی به مربی جهت به کارگیری آن‌ها در پست مناسب بدهد. این روش بیشتر در زمینه استعدادیابی بازیکنان در رده‌های پایه می‌تواند مفید واقع شود. زیرا بازیکنان را به بازی در پستی که در آن مستعدتر هستند، سوق می‌دهد.

تشریح داده‌ها:

در این قسمت ابتدا تمامی ویژگی‌ها را به اختصار توضیح می‌دهیم (در اسامی ویژگی‌ها به لاتین نوشته می‌شود؛ زیرا معادل دقیقی برای بعضی از ویژگی‌ها موجود نمی‌باشد):

- (۱) **Acceleration:** بیانگر این است که چقدر زمان صرف می‌شود تا پس از شروع به دویدن، بازیکن به حداکثر سرعتی که دارد برسد. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- (۲) **Sprint speed:** بیانگر بالاترین سرعت بازیکن در هنگام دویدن. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- (۳) **Finishing:** قابلیت تمام کنندگی گویند که بیانگر دقت شوت در محوطه پنالتی می‌باشد. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).

¹ Machine learning

² [FIFA 22 dataset](#)

³ Regression

⁴ Classification

- ۴) Long shot: دقت شوت در خارج از محوطه جریمه. (در قالب امتیاز از ۰ تا ۹۹ درجه‌بندی می‌شود).
- ۵) Penalties: دقت در ضربات پنالتی. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۶) Positioning: قابلیت جایگیری بازیکن در فضاهاى خالی برای دریافت توپ و خلق موقعیت خطرناک. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۷) Shot power: بیانگر میزان قدرتی ست که بازیکن هنگام زدن شوت با حفظ دقت شوت به توپ منتقل می‌کند. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۸) Volleys: توانایی در زدن شوت هنگامی توپ در هوا ست. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۹) Crossing: کیفیت سانتر چه در جریان بازی و در هنگام کرنر یا ضربه آزاد. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۰) Curve: توانایی دادن کات به توپ هنگام پاس و شوت. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۱) Free kick: دقت در زدن ضربات آزاد. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۲) Long pass: توانایی و دقت دادن پاس بلند. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۳) Short pass: توانایی و دقت دادن پاس کوتاه. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۴) Vision: میزان آگاهی بازیکن از جایگاه دقیق هم تیمی‌هایش در زمین بازی. بالا بودن این قابلیت احتمال موفقیت پاس‌های بلند رو افزایش می‌دهد. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۵) Agility: بیانگر سرعت بازیکن در کنترل توپ می‌باشد. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۶) Balance: توانایی حفظ کردن تعادل در تقابل‌های فیزیکی در زمین. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۷) Ball control: توانایی کنترل توپ در هنگام دریافت پاس از هم تیمی‌ها. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۸) Composure: توانایی تحمل فشار از جانب بازیکن تیم حریف. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۱۹) Dribbling: توانایی دریبل زدن بازیکن تیم حریف. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).

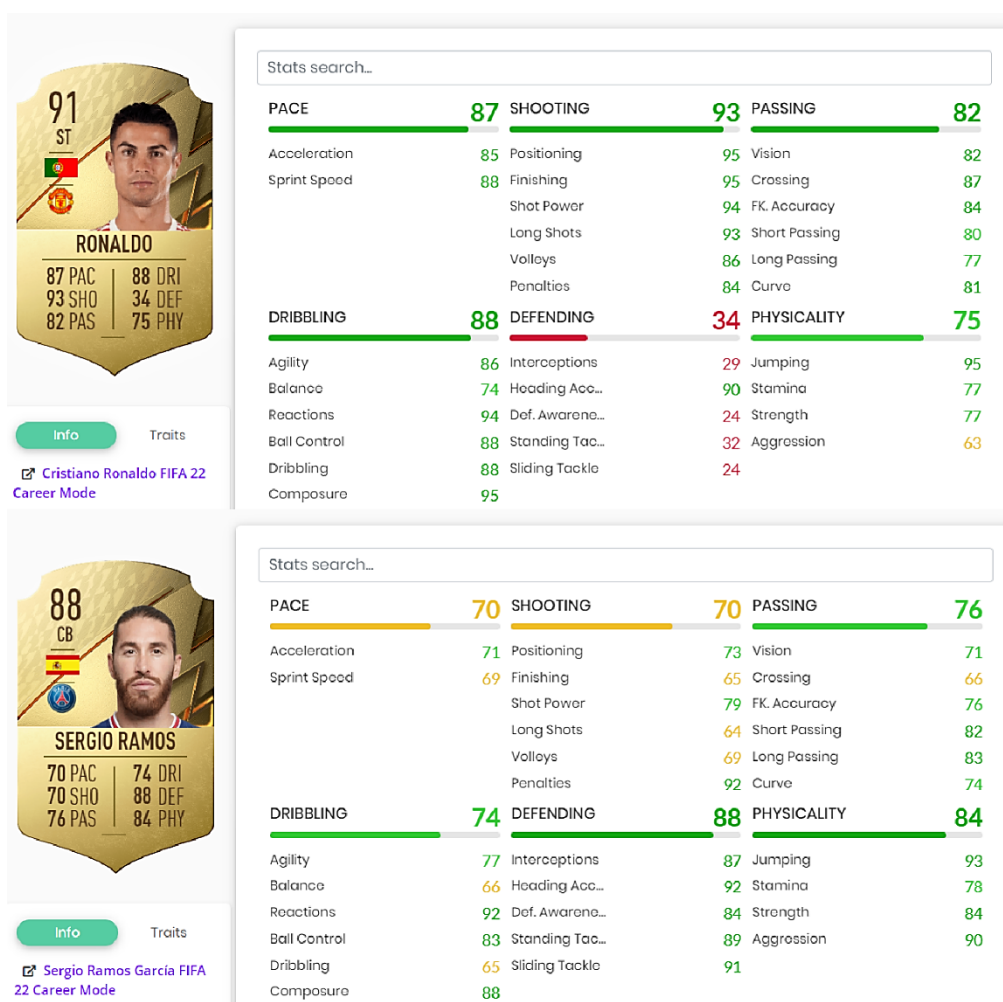
- ۲۰) Aggression: جگندگی. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۱) Jumping: توانایی پرش. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۲) Stamina: بازیکنی که این مقدار را به میزان بالا در قابلیت‌های فردی خود دارد، انتظار می‌رود تا به زودی را بدون خستگی در زمین بازی کند. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۳) Strength: قدرت از لحاظ فیزیکی. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۴) Defensive awareness: شامل قابلیت جا گیری مدافع و ندادن فرصت با حمله‌کننده در ایجاد خلق موقعیت با پاس یا شوت ناگهانی. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۵) Heading accuracy: دقت ضربه سر برای پاس و شوت. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۶) Interception: قابلیت قطع توپ. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۷) Sliding tackle: بیانگر دقت و زمان‌بندی مناسب در زدن تکل. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۸) Gk_diving: قدرت شیرجه دروازه بان. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۲۹) Gk_handling: قدرت گرفت توپ و در دست گرفتن آن. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۳۰) Gk_kicking: بیانگر قدرت و دقت ضربه دروازه. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۳۱) Gk_positioning: جایگیری دروازه بان در چهارچوب دروازه. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۳۲) Gk_speed: سرعت دروازه بان در مصاف‌های تک به تک. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۳۳) Gk_reflexes: انعطاف پذیری دروازه بان رد شیرجه زدن برای دفع توپ. (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۳۴) Age: سن بازیکن
- ۳۵) Potential growth: میزان رشد احتمالی در آینده
- ۳۶) Overall: معدل کل توانایی‌های فردی بازیکن (در قالب امتیاز از ۰ تا ۹۹ درجه بندی می‌شود).
- ۳۷) International reputation: درجه معروف بودن بازیکن

دو ویژگی دیگر نیز وجود دارند که به عنوان متغیر پاسخ^۵ در مسائل رگرسیون و طبقه‌بندی مورد استفاده قرار می‌گیرند:

۳۸) Position: پست بازیکن در زمین. که به ۴ کلاس دروازه‌بان، دفاع، هافبک، و مهاجم دسته‌بندی می‌شدند، ولی به دلیل محدودیت‌های موجود در روش‌هایی که به آن‌ها مسلط هستیم و نلبالانس بودن داده‌ها، کلاس‌دفاع و دروازه‌بان تحت عنوان دفاع- دروازه‌بان، و کلاس مهاجم و هافبک تحت عنوان هافبک-مهاجم در نظر گرفته شد.

۳۹) Value euro: ارزش بازیکنان در بازار نقل انتقالات به یورو.

به عنوان مثال توانایی‌های یک مدافع و یک مهاجم در شکل ۶۱ آمده است.

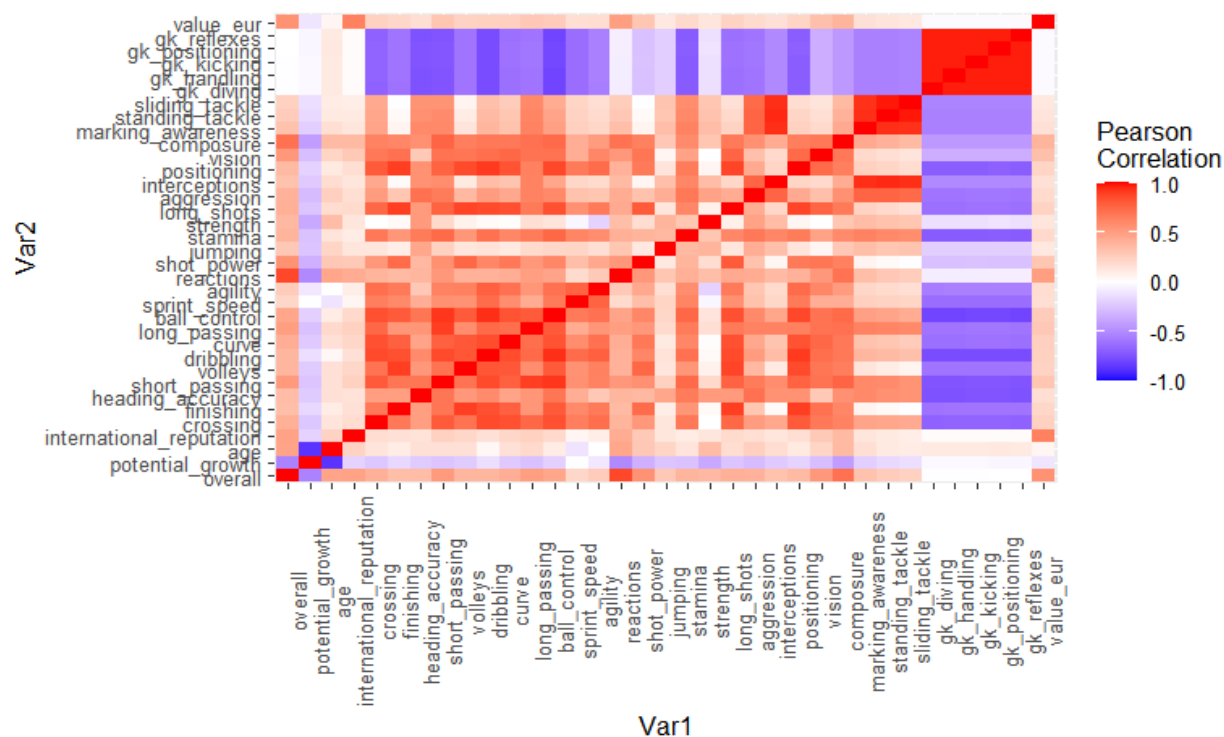


شکل ۱: نمونه‌ای از توانایی‌ها دو بازیکن در پست دفاع و مهاجم

⁵ Output variable

⁶ <https://www.futbin.com/22/player/589/sergio-ramos-garcia>
<https://www.futbin.com/22/player/426/cristiano-ronaldo>

به دلیل تعداد بالای ویژگی‌ها و وابستگی تعداد قابل توجهی از آن‌ها به یکدیگر، نمی‌توان تصویر جامعی با کمک روش‌های مبتنی بر تحلیل اکتشافی^۷ بدست آورد. تنها تصویر جامع ممکن از این مجموعه داده، مصورسازی دو به دو ویژگی‌ها در مقادیر همبستگی^۸ می‌باشد، تا بتوان بدان وسیله در انتخاب متغیرهای مناسب در مدل‌های پیش‌بینی به ما کمک کند.



نمودار ۱: نمودار شدت همبستگی بین ویژگی‌ها

آنچه در نمودار ۱ مشهود است همبستگی تعداد قابل توجهی از ویژگی‌ها با یکدیگر است. برای مثال قابلیت‌های دروازه‌بانی رابطه مستقیمی با یکدیگر دارند این در حالی است که با سایر توانایی‌ها رابطه عکس دارند زیرا دروازه‌بانی تفاوت زیادی با سایر پست‌ها در بازی فوتبال دارد که قابلیت‌های آن نیز تفاوت قابل توجهی با دیگر مهارت‌ها دارد.

برای ذکر مثال دیگر می‌توان به رابطه سن با پتانسیل رشد اشاره کرد زیرا سن مفید در فوتبال به طور متوسط ۳۳ سال تا ۳۴ سال تخمین زده شده و بازیکنان از این سن به بعد دچار افول می‌شود پس معمولاً فرصتی برای رشد مهارت‌های خود نخواهند داشت.

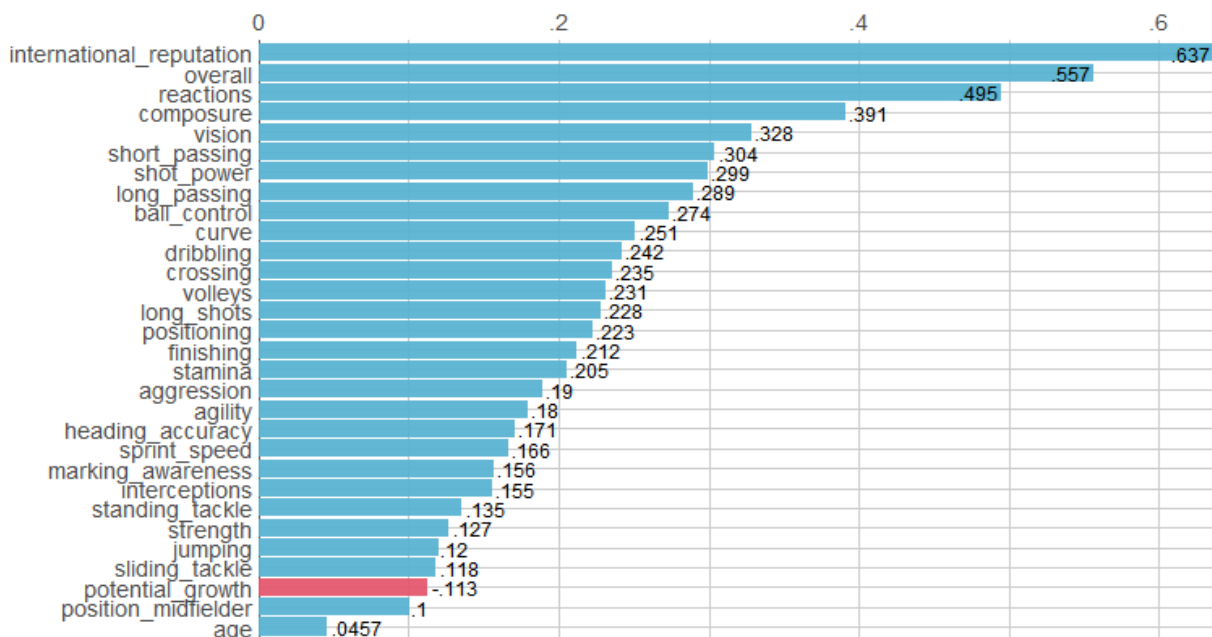
^۷ Exploratory data analysis

^۸ Correlation

مورد دیگری که می‌توان به آن پرداخت به طور خاص رتبه‌بندی میزان همبستگی ویژگی‌ها با قیمت بازیکن می‌باشد که می‌تواند به طور خاص در مسئله رگرسیون راهگشا باشد.

Correlations of value_eur

30 largest correlation variables (original & dummy)



نمودار ۲: رتبه بندی ویژگی‌ها در همبستگی برای انتخاب ویژگی

آماده‌سازی داده‌ها:

داده‌های موجود دارای مقداری داده گمشده می‌باشند. که در این مرحله باید پاکسازی شوند. با پاکسازی صورت گرفته تعداد مشاهدات به ۱۹۱۶۵ بازیکن می‌رسد.

همانطور که در شکل ۲ قابل مشاهده است، مرحله بعدی در تقسیم‌بندی داده‌ها، داده‌ها به سه گروه داده اعتبارسنجی^۹، برآورد^{۱۰} و آزمایشی^{۱۱} تقسیم‌بندی می‌شوند. به این صورت که ۸۰ درصد داده‌ها به داده‌های آموزشی^{۱۲} و ۲۰ درصد باقی‌مانده به داده‌های آزمایشی تخصیص می‌یابد. همچنین ۸۰ درصد داده‌های تخصیص داده شده به داده‌های آموزشی،

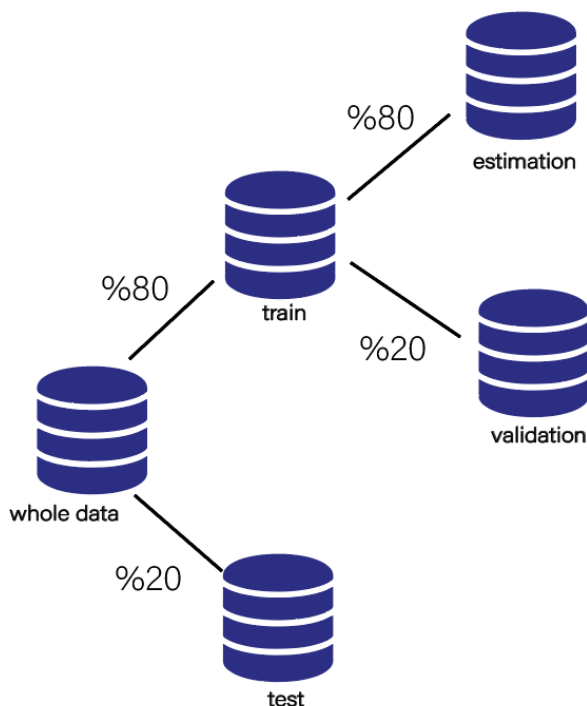
^۹ Validation set

^{۱۰} Estimation set

^{۱۱} Test set

^{۱۲} Train set

خود نیز به دو گروه داده تقسیم‌بندی می‌شود که ۸۰ درصد آن به مجموعه داده برآورد و ۲۰ درصد آن به مجموعه داده اعتبارسنجی تخصیص می‌یابد.



شکل ۲: شماتیک تقسیم‌بندی داده‌ها

همچنین در مراحل پیش رو نیاز به استفاده از داده‌های نرمال سازی شده داریم (در روش k نزدیک ترین همسایه^{۱۳}). به همین جهت داده‌های برآورد را با استفاده از میانگین و انحراف معیار بدست آمده از خود آن داده‌ها نرمال سازی می‌کنیم و داده‌های اعتبارسنجی را با کمک میانگین و انحراف معیار داده‌های برآورد نرمال سازی^{۱۴} می‌کنیم.

برای نرمال سازی داده‌های آزمایشی باید ابتدا داده‌های آموزشی را نرمال سازی کرده سپس داده‌های آزمایشی را با استفاده از میانگین و انحراف معیار داده‌های آموزشی نرمال سازی می‌کنیم.

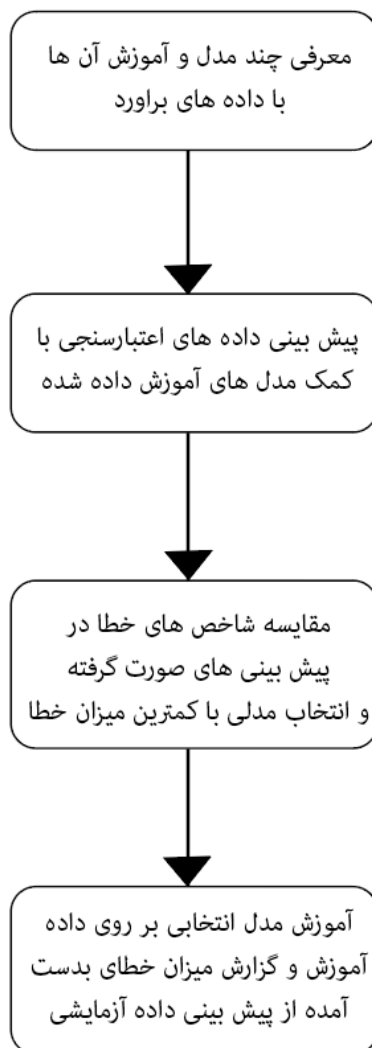
و همانطور که در توضیح ویژگی موقعیت بازیکن گفتیم، مجبور به کاهش تعداد کلاس‌های مان برای مسئله طبقه‌بندی شدیم.

¹³ K-nearest-neighbor

¹⁴ Normalization

پیاده‌سازی روش‌ها و ارزیابی:

به طور کلی در بدست آوردن مدل بهینه در مسائل یادگیری ماشین، گام‌های ذکر شده در شکل ۳ طی می‌شود.



شکل ۳: فلوچارت مراحل گزینش مدل بهینه با کمترین میزان خطا

این گام‌ها به تفصیل در دو پروژه رگرسیون و طبقه‌بندی انجام شده.

معیار ارزیابی در مسئله رگرسیون:

مدل‌های رگرسیون مجذور مربعات خطا^{۱۵} می‌باشد. که به صورت زیر می‌باشد.

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

معیار ارزیابی^{۱۶} در مسئله طبقه‌بندی:

در مسائل طبقه‌بندی مختلف زمانی که با دو کلاس مواجه هستیم شاخص‌های ارزیابی مختلفی برای مدل مطرح می‌گردد. نظیر sensitivity, specificity, precision, accuracy و غیره. موضوعی که در رابطه با این شاخصه‌های ارزیابی مطرح می‌شود این است که این‌ها زمانی مفید واقع می‌شود که اشتباه یادداشت پیش‌بینی کردن یک کلاس نسبت به کلاس دیگر ارجحیت داشته باشد. برای مثال در پیش‌بینی بیماری سرطان به اشتباه سالم تشخیص دادن بیمار به مراتب تأثیرات مخرب‌تری دارد نسبت به حالتی که به اشتباه فرد سالم رو بیمار تشخیص بدهیم و می‌تواند موجب وقوع اتفاقات جبران ناپذیری شود. در نتیجه شاخص ارزیابی ما sensitivity می‌باشد. و سعی می‌شود با پایین بردن آستانه مقدار این شاخص را افزایش دهیم.

اما در اینجا ما دو کلاس داریم یک کلاس مدافع-دروازه بان و کلاس هافبک-مهاجم هر دوی این دو کلاس در پیش‌بینی یکسان می‌باشد و ارجحیتی نسبت به یکدیگر ندارند. لذا ترجیح ما این است که شاخص accuracy را ملاک ارزیابی قرار دهیم.

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

رگرسیون:

(۱) مدل رگرسیون خطی:

جزئیات مدل‌های آموزش داده شده در جدول ۱ آمده است.

¹⁵ Root mean squared error

¹⁶ Evaluation metric

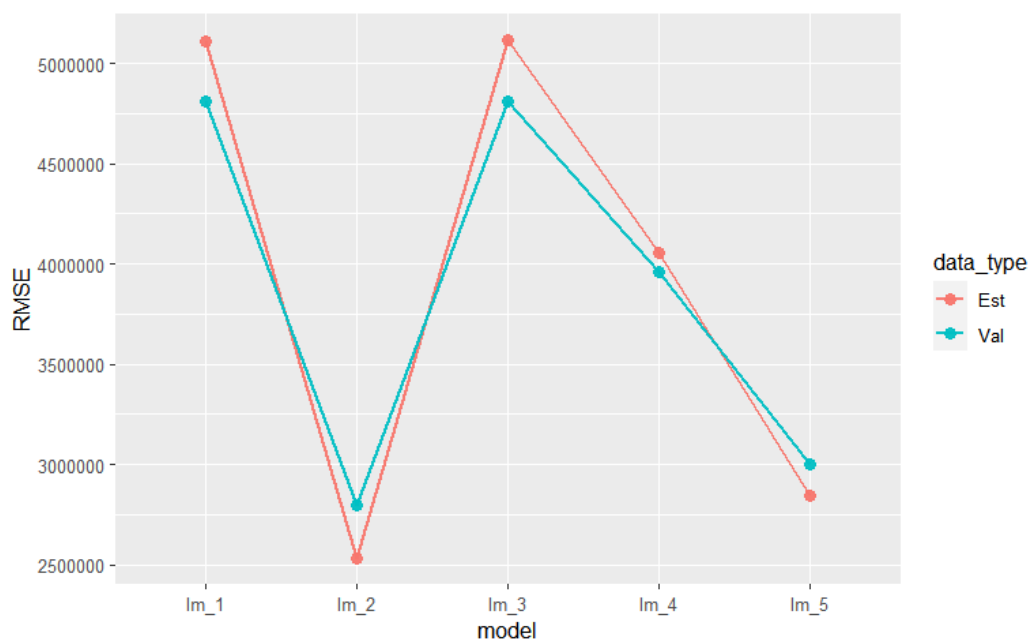
جدول ۱: توضیحات مدل‌های برازش داده شده‌ی رگرسیون خطی روی داده‌های برآورد

شماره مدل	متغیرهای مورد استفاده	توضیحات مدل
۱	همه متغیرها	-
۲	همه متغیرها	اثرات متقابل
۳	همه متغیرها	دستور گام برای انتخاب متغیر با روش دو طرفه
۴	تعدادی از متغیرها	استفاده از دستور $\text{poly}(x,2)$ برای تمامی متغیرهای انتخاب شده
۵	تعدادی از متغیرها	استفاده از دستور $\text{poly}(x,2)$ برای تمامی متغیرهای انتخاب شده

مقادیر مجذور مربعات خطا مدل‌ها بر روی داده‌های اعتبارسنجی مطابق جدول ۲ و نمودار ۳ می‌باشد.

جدول ۲: میزان خطای مجذور مربعات خطای مدل‌های رگرسیون در خطی در بیش‌بینی داده‌های اعتبارسنجی

شماره مدل	RMSE-validation
۱	۴۸۱۳۸۸۵
۲	۲۷۹۷۶۰۶
۳	۴۸۱۱۷۶۳
۴	۳۹۶۰۹۳۹
۵	۳۰۰۰۲۶۶



نمودار ۳: نمودار میزان خطای مدل‌های رگرسیون خطی روی داده‌های اعتبارسنجی و برآورد

۲) مدل K نزدیک‌ترین همسایه:

مدل را به ازای مقادیر مختلف k روی داده‌های برآورد استاندارد شده و غیر استاندارد شده برازش دادیم. که جزئیات مدل‌های آموزش داده شده در جدول ۳ آمده است.

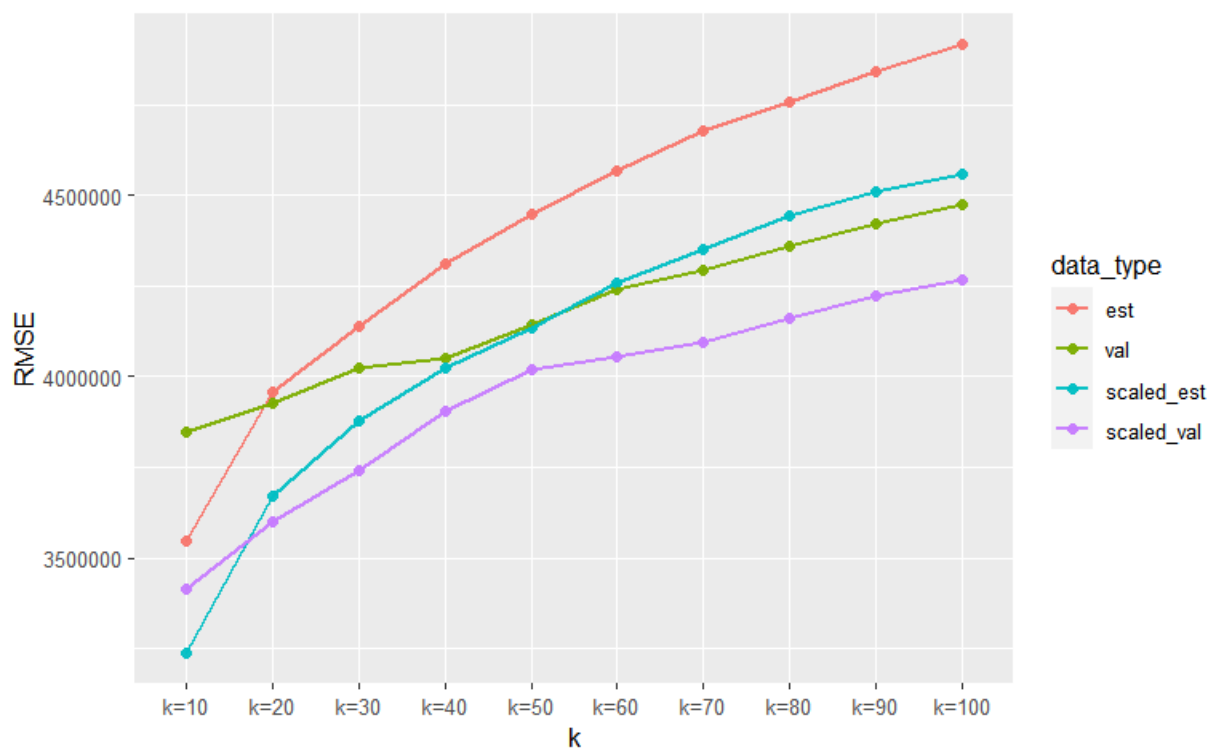
جدول ۳: توضیحات مربوط به مدل‌های برازش داده شده‌ی ک نزدیک‌ترین همسایه روی داده‌های برآورد

K	متغیرهای مورد استفاده	توضیحات مدل
۱۰	همه متغیرها	داده استاندارد شده
۲۰	همه متغیرها	داده استاندارد شده
۳۰	همه متغیرها	داده استاندارد شده
۴۰	همه متغیرها	داده استاندارد شده
۵۰	همه متغیرها	داده استاندارد شده
۶۰	همه متغیرها	داده استاندارد شده
۷۰	همه متغیرها	داده استاندارد شده
۸۰	همه متغیرها	داده استاندارد شده
۹۰	همه متغیرها	داده استاندارد شده
۱۰۰	همه متغیرها	داده استاندارد شده
۱۰	همه متغیرها	داده غیر استاندارد
۲۰	همه متغیرها	داده غیر استاندارد
۳۰	همه متغیرها	داده غیر استاندارد
۴۰	همه متغیرها	داده غیر استاندارد
۵۰	همه متغیرها	داده غیر استاندارد
۶۰	همه متغیرها	داده غیر استاندارد
۷۰	همه متغیرها	داده غیر استاندارد
۸۰	همه متغیرها	داده غیر استاندارد
۹۰	همه متغیرها	داده غیر استاندارد
۱۰۰	همه متغیرها	داده غیر استاندارد

مقادیر مجذور مربعات خطا مدل‌ها بر روی داده‌های اعتبارسنجی استاندارد شده و غیر استاندارد مطابق جدول ۴ و نمودار ۴ می‌باشد.

جدول ۴: مقادیر مجذور مربعات خطای مدل‌های ک نزدیک ترین همسایه روی داده‌های اعتبارسنجی استاندارد شده و غیر استاندارد

RMSE (validation-scaled)	RMSE (validation-non-scaled)	k
۳۴۱۵۴۶۲	۳۸۴۷۷۱۱	۱۰
۳۵۹۹۹۶۲	۳۹۲۸۱۴۳	۲۰
۳۷۴۰۹۷۵	۴۰۲۴۴۰۴	۳۰
۳۹۰۴۸۸۸	۴۰۵۲۳۳۵	۴۰
۴۰۱۹۱۲۸	۴۱۴۳۶۰۱	۵۰
۴۰۵۶۸۵۷	۴۲۳۹۰۵۶	۶۰
۴۰۹۴۸۸۴	۴۲۹۳۶۵۶	۷۰
۴۱۶۲۸۳۶	۴۳۵۹۵۱۶	۸۰
۴۲۲۴۷۶۱	۴۴۲۰۴۰۹	۹۰
۴۲۶۷۸۷۲	۴۴۷۲۴۶۱	۱۰۰



نمودار ۴: مقادیر مجذور مربعات خطا برای مدل مدل‌های ک نزدیک ترین همسایه بر روی داده های برآورد و اعتبارسنجی استاندارد شده و غیر استاندارد

۳) مدل درخت رگرسیون^{۱۷}:

مدل درخت رگرسیون را به ازای مقادیر cp و $minsplit$ روی داده برآورد آموزش دادیم. در مدل درخت رگرسیون انشعاباتی با برش روی ویژگی‌های مختلف ایجاد می‌شود. برای متوقف شدن هر یک از این برش‌ها باید یکی از دو شرط ذیل برقرار گردد. اولاً این که با ایجاد برش آماره R^2 کمتر از مقدار cp افزایش یابد. دوماً اینکه تعداد داده‌هایی که می‌خواهیم روی آن برش بنزیم تعدادشان کوچکتر مساوی $minsplit$ باشد.

ویژگی مدل‌های برازش^{۱۸} داده شده مطابق جدول ۵ می‌باشد.

جدول ۵: توضیحات مربوط به مدل‌های برازش داده شده‌ی درخت رگرسیون روی داده‌های برآورد

متغیرهای استفاده شده	cp	$minsplit$
همه متغیرها	۰.۱	۱۰
همه متغیرها	۰.۱	۴۰
همه متغیرها	۰.۱	۷۰
همه متغیرها	۰.۱	۱۰۰
همه متغیرها	۰.۰۱	۱۰
همه متغیرها	۰.۰۱	۴۰
همه متغیرها	۰.۰۱	۷۰
همه متغیرها	۰.۰۱	۱۰۰
همه متغیرها	۰.۰۰۱	۱۰
همه متغیرها	۰.۰۰۱	۴۰
همه متغیرها	۰.۰۰۱	۷۰
همه متغیرها	۰.۰۰۱	۱۰۰

مقادیر مجذور مربعات خطا مطابق جدول ۶ و نمودار ۵ می‌باشد.

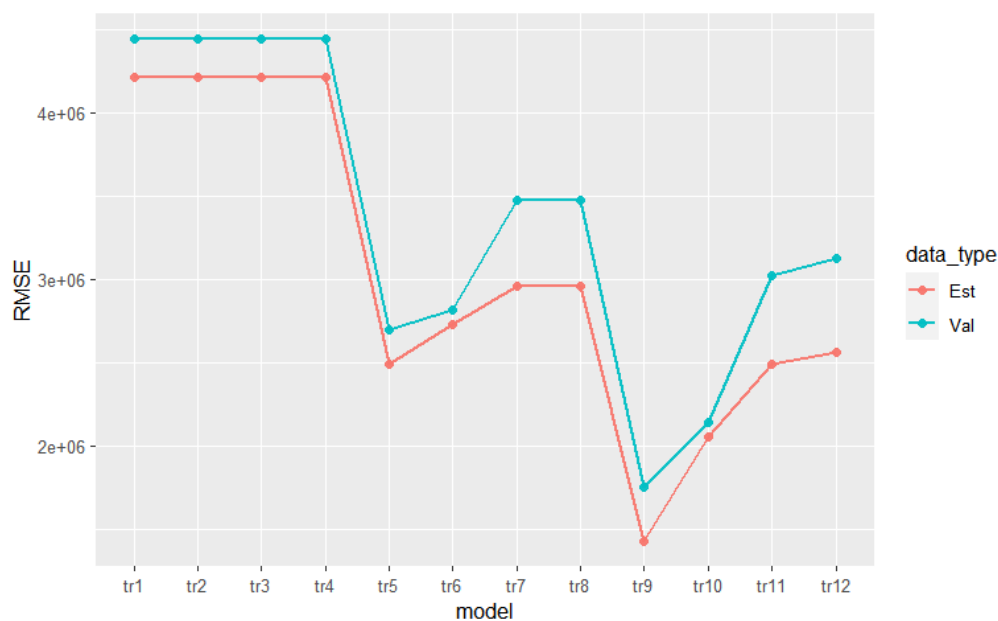
جدول ۶: مقادیر مجذور مربعات خطا برای مدل‌های درخت رگرسیون در بیش‌بینی داده‌های اعتبارسنجی

RMSE-validation	cp	$minsplit$
۴۴۳۸۷۳۵	۰.۱	۱۰
۴۴۳۸۷۳۵	۰.۱	۴۰
۴۴۳۸۷۳۵	۰.۱	۷۰
۴۴۳۸۷۳۵	۰.۱	۱۰۰
۲۶۹۴۴۴۳	۰.۰۱	۱۰

^{۱۷} Regression tree

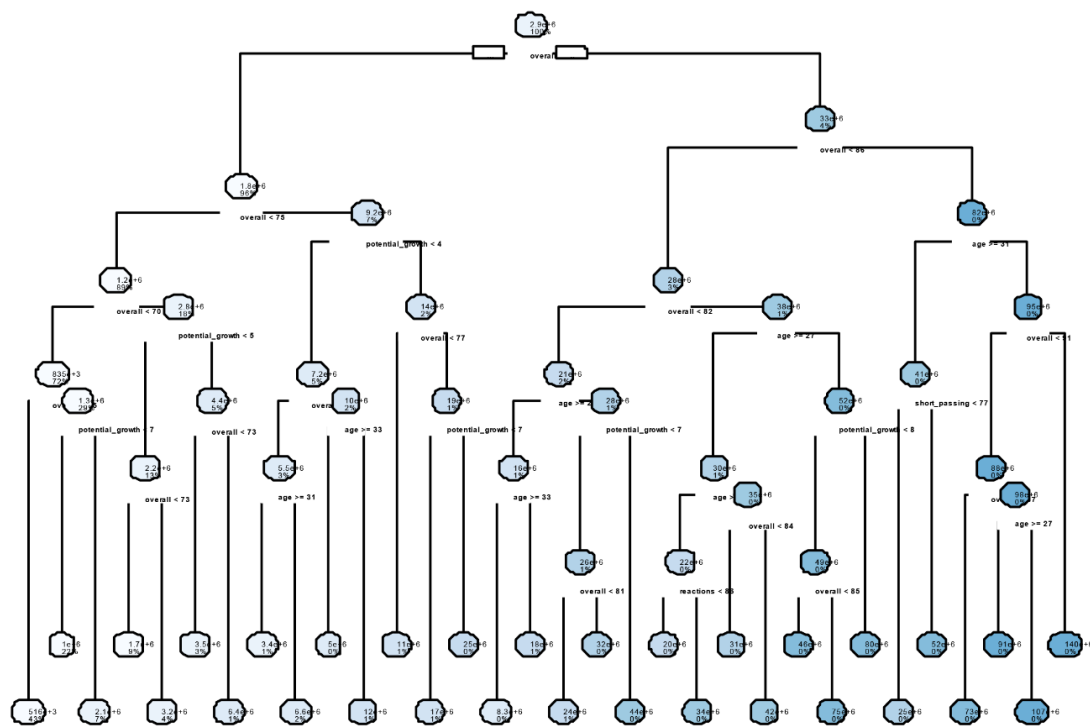
^{۱۸} fit

۲۸۱۳۳۳۵	۰.۰۱	۴۰
۳۴۷۸۸۲۰	۰.۰۱	۷۰
۳۴۷۸۸۲۰	۰.۰۱	۱۰۰
۱۷۵۴۵۴۲	۰.۰۰۱	۱۰
۲۱۴۰۳۷۷	۰.۰۰۱	۴۰
۳۰۲۲۴۲۸	۰.۰۰۱	۷۰
۳۱۲۷۸۷۸	۰.۰۰۱	۱۰۰



نمودار ۵: مقادیر مجذور مربعات خطا برای مدل‌های درخت رگرسیون در بیش‌بینی داده‌های اعتبارسنجی و برآورد

با توجه به نمودار ۵ و جدول ۶، مدل درخت شماره ۹ ($cp = 0.001$, $minsplit = 10$) با کمترین مجذور مربعات خطا، مدل انتخابی ما از بین درخت‌ها می‌باشد. نمودار انشعابات مدل درخت رگرسیون انتخابی مطابق نمودار ۷ می‌باشد.

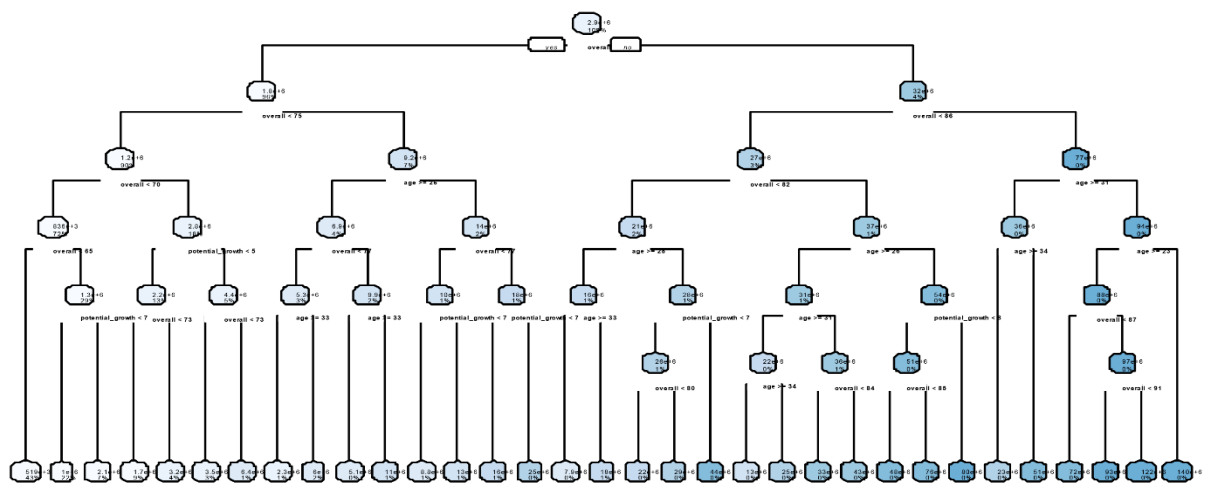


نمودار ۷: نمودار انشعابات درخت برازش داده‌شده (روی داده‌های برآورد) انتخابی از بین مدل‌های درخت رگرسیون

جمع‌بندی نهایی برای مدل‌های رگرسیون:

از بین تمامی مدل‌های انتخابی، مدل درخت رگرسیون ($\text{minsplit} = 10$, $\text{cp} = 0.001$) با میزان مجذور مربعات خطای ۱۷۵۴۵۴۲ به عنوان بهترین مدل انتخاب گردید.

حال مدل انتخابی را روی مجموعه داده آموزش برازش می‌دهیم و مقدار مجذور مربعات خطای بدست آمده برابر با ۲۱۴۳۰۷۹ می‌باشد. درخت برازش داده شده مطابق نمودار ۸ است.



نمودار ۸: نمودار انشعابات درخت رگرسیون (مربوط به مدل منتخب که روی داده‌های آموزش برازش داده شده است).

طبقه‌بندی:

(۱) رگرسیون لجستیک:

جزئیات مدل‌های آموزش داده شده در جدول ۷ آمده است.

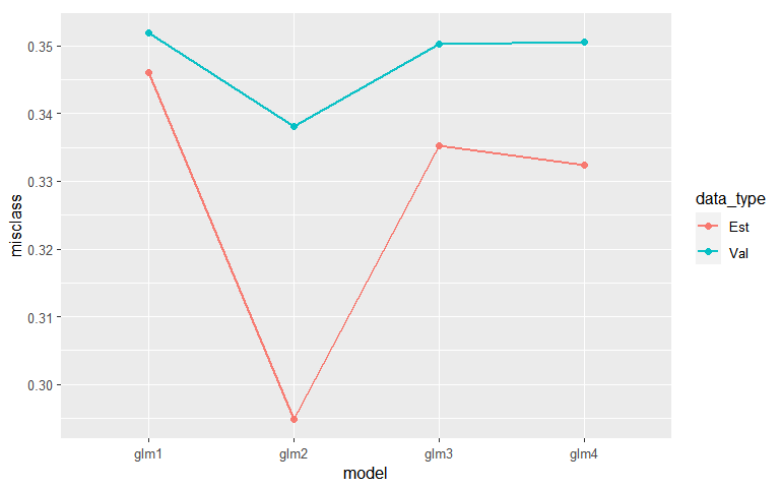
جدول ۷: توضیحات مدل‌های رگرسیون لجستیک برازش داده شده روی داده‌های برآورد

شماره مدل	متغیرهای مورد استفاده	توضیحات مدل
۱	همه متغیرها	-
۲	همه متغیرها	اثرات متقابل
۳	تعدادی از متغیرها	استفاده از دستور $\text{poly}(x,2)$ برای تمامی متغیرهای انتخاب شده
۴	تعدادی از متغیرها	استفاده از دستور $\text{poly}(x,3)$ برای تمامی متغیرهای انتخاب شده

مقادیر خطای طبقه‌بندی^{۱۹} مدل‌های برازش داده شده در پیش‌بینی داده‌های اعتبارسنجی مطابق جدول ۸ و نمودار ۹ آمده است.

جدول ۸: مقادیر خطای طبقه‌بندی مدل‌های رگرسیون لجستیک در پیش‌بینی داده‌های اعتبارسنجی

شماره مدل	خطای طبقه‌بندی (اعتبارسنجی)
۱	۰.۳۵۱۸۰۹۶
۲	۰.۳۳۸۱۱۵۴
۳	۰.۳۵۰۱۷۹۳
۴	۰.۳۵۰۵۰۵۴



نمودار ۹: مقادیر خطای طبقه‌بندی مدل‌های رگرسیون لجستیک در پیش‌بینی داده‌های برآورد و اعتبارسنجی

¹⁹ Misclassification error

۲) مدل تحلیل تفکیکی خطی^{۲۰}:

جزئیات مدل‌های آموزش داده شده به شرح زیر در جدول ۹ بیان شده است.

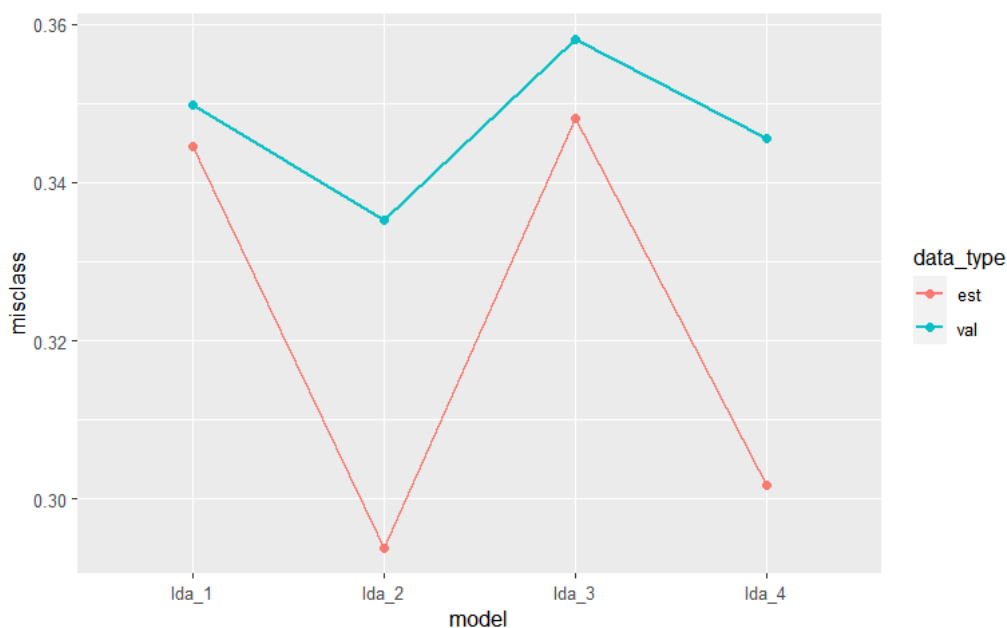
جدول ۹: توضیحات مدل‌های تحلیل تفکیکی خطی برآزش داده شده روی داده‌های برآورد

شماره مدل	متغیرهای مورد استفاده	توضیحات مدل
۱	همه متغیرها	برآورد احتمال پیشین از روی داده برآورد
۲	همه متغیرها	اثرات متقابل، $\text{poly}(x,2)$ و برآورد احتمال پیشین از روی داده برآورد
۳	همه متغیرها	احتمال پیشین یکنواخت
۴	همه متغیرها	اثرات متقابل، $\text{poly}(x,2)$ و احتمال پیشین یکنواخت

مقادیر خطای طبقه‌بندی مدل‌های برآزش داده شده در پیش‌بینی داده‌های اعتبار سنجی مطابق جدول ۱۰ و نمودار ۱۰

است. جدول ۱۰: مقادیر خطای طبقه بندی مدل‌های تحلیل تفکیکی خطی در پیش‌بینی داده‌های اعتبارسنجی

شماره مدل	خطای طبقه بندی (اعتبارسنجی)
۱	۰.۳۴۹۸۵۳۳
۲	۰.۳۳۵۱۸۱۰
۳	۰.۳۵۸۰۰۴۶
۴	۰.۳۴۵۶۱۴۶



نمودار ۱۰: مقادیر خطای طبقه بندی مدل‌های تحلیل تفکیکی خطی در پیش‌بینی داده‌های برآورد و اعتبارسنجی

²⁰ Linear discriminant analysis

۳) مدل تحلیل تفکیکی کوآدراتیک^{۲۱}:

جزئیات مدل‌های آموزش داده شده به شرح جدول ۱۱ است.

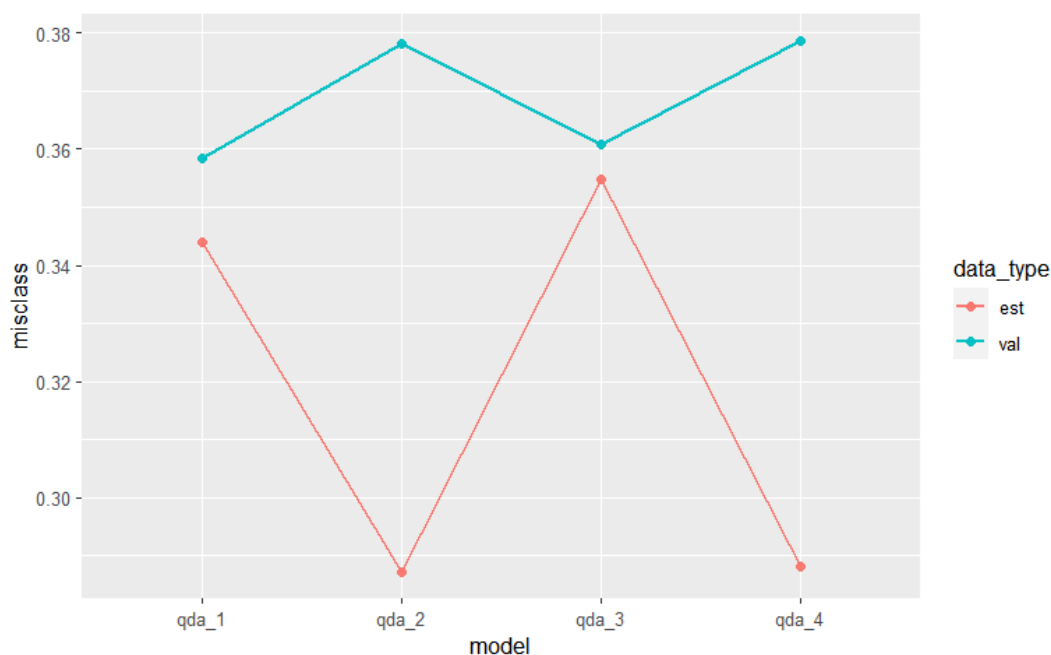
جدول ۱۱: توضیحات مدل‌های تحلیل تفکیکی کوآدراتیک در پیش‌بینی داده‌های اعتبارسنجی

شماره مدل	متغیرهای مورد استفاده	توضیحات مدل
۱	همه متغیرها	برآورد احتمال پیشین از روی داده برآورد
۲	همه متغیرها	اثرات متقابل، $\text{poly}(x,2)$ و برآورد احتمال پیشین از روی داده برآورد
۳	همه متغیرها	احتمال پیشین یکنواخت
۴	همه متغیرها	اثرات متقابل، $\text{poly}(x,2)$ و احتمال پیشین یکنواخت

مقادیر خطای طبقه‌بندی مدل‌های برازش داده شده در پیش‌بینی داده‌های اعتبارسنجی مطابق جدول ۱۲ و نمودار ۱۱ است.

جدول ۱۲: مقادیر خطای طبقه‌بندی مدل‌های تحلیل تفکیکی کوآدراتیک در پیش‌بینی داده‌های اعتبارسنجی

شماره مدل	خطای طبقه‌بندی (اعتبارسنجی)
۱	۰.۳۵۸۳۳۰۶
۲	۰.۳۷۸۲۱۹۸
۳	۰.۳۶۰۹۳۹۰
۴	۰.۳۷۸۵۴۵۸



نمودار ۱۱: مقادیر خطای طبقه‌بندی مدل‌های تحلیل تفکیکی کوآدراتیک در پیش‌بینی داده‌های اعتبارسنجی و برآورد

²¹ Quadratic discriminant analysis

مشاهده می‌شود که مدل تحلیل تفکیکی کوآدراتیک دچار پیش‌برازش شده است. و مدل مناسبی برای طبقه‌بندی داده‌های ما نیست.

(۴) مدل بیز ساده^{۲۲}:

با مستقل در نظر گرفتن ویژگی‌ها نسبت به یکدیگر مدل تحلیل تفکیکی خطی به مدل بیز ساده تغییر فرم می‌دهد در اینجا ما یک مدل بیز ساده را با در نظر گرفتن تمام ویژگی‌ها روی داده برآورد برازش دادیم. که مقدار خطای طبقه‌بندی آن روی مجموعه داده اعتبارسنجی ۰.۴۴۷۳۴۳ شده است.

(۵) مدل درخت تصمیم^{۲۳}:

مدل‌های درخت طبقه‌بندی را با مقادیر cp برابر با ۰.۱ و ۰.۰۱ و ۰.۰۰۱ با $minsplit$ های بین ۱۰ تا ۲۰۰ با گام ۱۰ بر روی داده‌های برآورد برازش داده‌ایم که در مجموع ۶۰ مدل به دست آمده که در جدول ۱۳ قابل مشاهده می‌باشد. همچنین نمودارها به تفکیک cp در نمودار ۱۲ قابل مشاهده است.

جدول ۱۳: مقادیر خطای طبقه‌بندی مدل‌های درخت تصمیم در پیش‌بینی داده‌های اعتبارسنجی

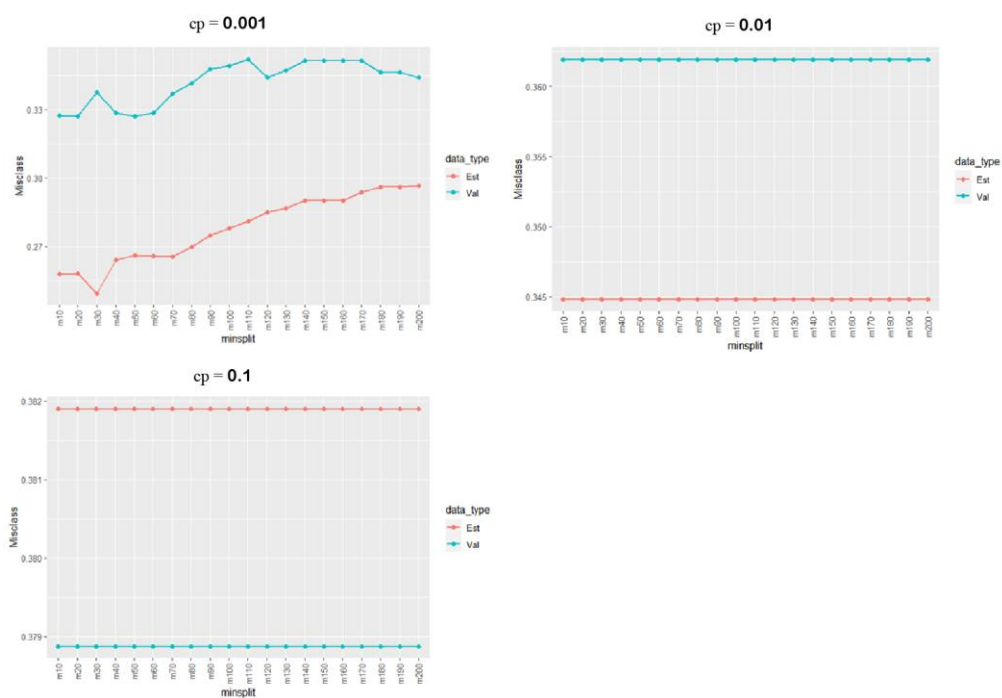
خطای طبقه بندی	cp	minsplit
۰.۳۷۸۸۷۱۹	۰.۱	۱۰
۰.۳۷۸۸۷۱۹	۰.۱	۲۰
۰.۳۷۸۸۷۱۹	۰.۱	۳۰
۰.۳۷۸۸۷۱۹	۰.۱	۴۰
۰.۳۷۸۸۷۱۹	۰.۱	۵۰
۰.۳۷۸۸۷۱۹	۰.۱	۶۰
۰.۳۷۸۸۷۱۹	۰.۱	۷۰
۰.۳۷۸۸۷۱۹	۰.۱	۸۰
۰.۳۷۸۸۷۱۹	۰.۱	۹۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۰۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۱۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۲۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۳۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۴۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۵۰

^{۲۲} Naïve Bayes

^{۲۳} Decision tree

۰.۳۷۸۸۷۱۹	۰.۱	۱۶۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۷۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۸۰
۰.۳۷۸۸۷۱۹	۰.۱	۱۹۰
۰.۳۷۸۸۷۱۹	۰.۱	۲۰۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۲۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۳۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۴۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۵۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۶۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۷۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۸۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۹۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۰۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۱۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۲۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۳۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۴۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۵۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۶۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۷۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۸۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۱۹۰
۰.۳۶۱۹۱۷۲	۰.۰۱	۲۰۰
۰.۳۲۷۳۵۵۷	۰.۰۰۱	۱۰
۰.۳۲۷۰۲۹۷	۰.۰۰۱	۲۰
۰.۳۳۷۴۶۳۳	۰.۰۰۱	۳۰
۰.۳۲۸۶۵۹۹	۰.۰۰۱	۴۰
۰.۳۲۷۰۲۹۷	۰.۰۰۱	۵۰
۰.۳۲۸۶۵۹۹	۰.۰۰۱	۶۰
۰.۳۳۷۱۳۷۳	۰.۰۰۱	۷۰
۰.۳۴۱۷۰۲۰	۰.۰۰۱	۸۰
۰.۳۴۷۸۹۷۰	۰.۰۰۱	۹۰
۰.۳۴۹۲۰۱۲	۰.۰۰۱	۱۰۰
۰.۳۵۲۱۳۵۶	۰.۰۰۱	۱۱۰

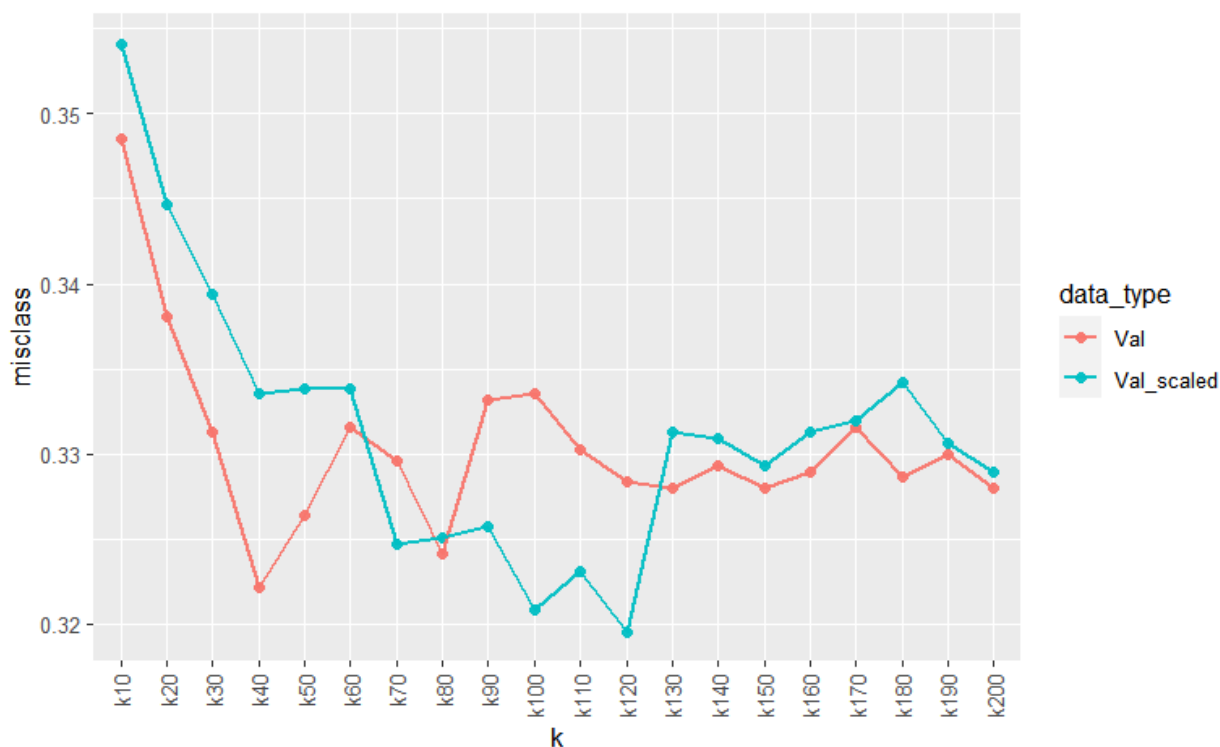
۰.۳۴۳۹۸۴۳	۰.۰۰۱	۱۲۰
۰.۳۴۷۲۴۴۹	۰.۰۰۱	۱۳۰
۰.۳۵۱۴۸۳۵	۰.۰۰۱	۱۴۰
۰.۳۵۱۴۸۳۵	۰.۰۰۱	۱۵۰
۰.۳۵۱۴۸۳۵	۰.۰۰۱	۱۶۰
۰.۳۵۱۴۸۳۵	۰.۰۰۱	۱۷۰
۰.۳۴۶۲۶۶۷	۰.۰۰۱	۱۸۰
۰.۳۴۶۲۶۶۷	۰.۰۰۱	۱۹۰
۰.۳۴۳۹۸۴۳	۰.۰۰۱	۲۰۰



نمودار ۱۲: مقادیر خطای طبقه‌بندی مدل‌های درخت تصمیم در پیش‌بینی داده‌های اعتبارسنجی و برآورد

درخت بدست آمده از برآزش بر روی داده‌های برآورد با ($cp = 0.001$ و $minsplit = 50$) در نمودار ۱۳ آمده است.

۰.۳۳۳۸۷۶۸	۰.۳۳۱۵۹۴۴	۶۰
۰.۳۲۴۷۴۷۳	۰.۳۲۹۶۳۸۱	۷۰
۰.۳۲۵۰۷۳۴	۰.۳۲۴۰۹۵۲	۸۰
۰.۳۲۵۷۲۵۵	۰.۳۳۳۲۲۴۶	۹۰
۰.۳۲۰۸۳۴۷	۰.۳۳۳۵۵۰۷	۱۰۰
۰.۳۲۳۱۱۷۱	۰.۳۳۰۲۹۰۲	۱۱۰
۰.۳۱۹۵۳۰۵	۰.۳۲۸۳۳۳۹	۱۲۰
۰.۳۳۱۲۶۸۳	۰.۳۲۸۰۰۷۸	۱۳۰
۰.۳۳۰۹۴۲۳	۰.۳۲۹۳۱۲۰	۱۴۰
۰.۳۲۹۳۱۲۰	۰.۳۲۸۰۰۷۸	۱۵۰
۰.۳۳۱۲۶۸۳	۰.۳۲۸۹۸۶۰	۱۶۰
۰.۳۳۱۹۲۰۴	۰.۳۳۱۵۹۴۴	۱۷۰
۰.۳۳۴۲۰۲۸	۰.۳۲۸۶۵۹۹	۱۸۰
۰.۳۳۰۶۱۶۲	۰.۳۲۹۹۶۴۱	۱۹۰
۰.۳۲۸۹۸۶۰	۰.۳۲۸۰۰۷۸	۲۰۰



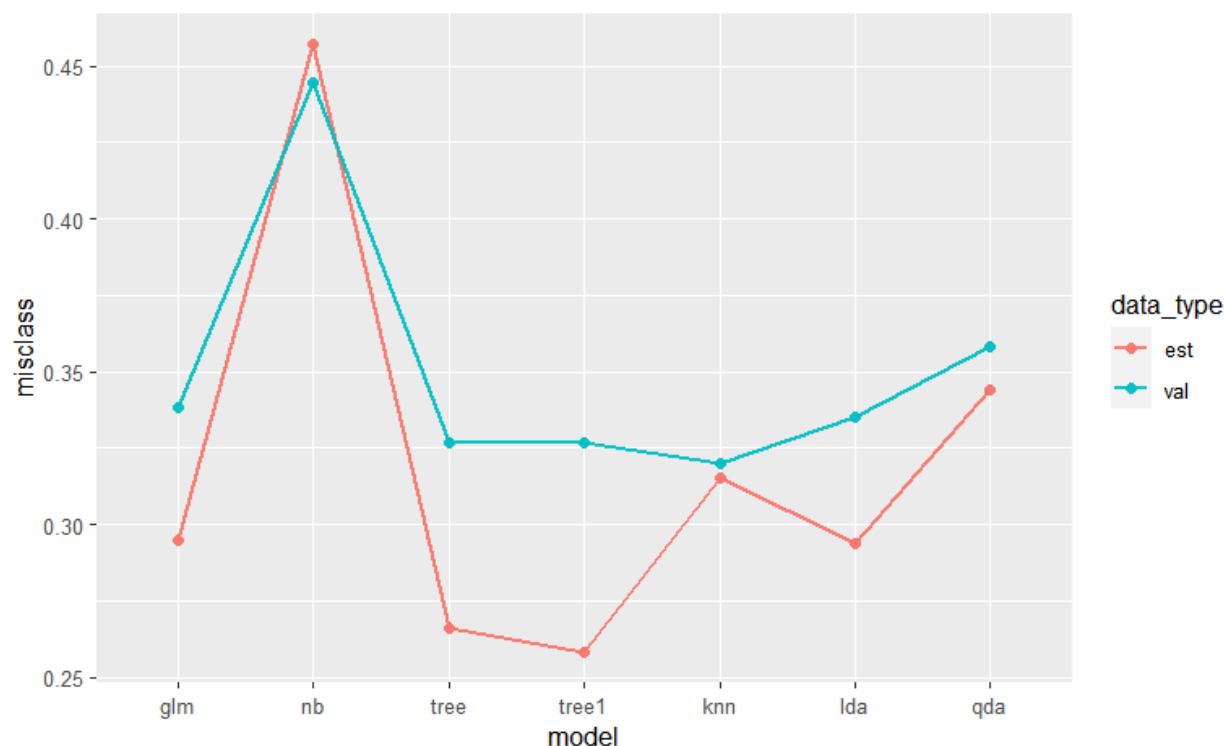
نمودار ۱۴: مقادیر خطای طبقه‌بندی مدل‌های درخت تصمیم در پیش‌بینی داده‌های اعتبارسنجی استاندارد شده و غیر استاندارد

جمع‌بندی نهایی برای مدل‌های طبقه‌بندی:

در این مرحله مدل‌های منتخب از هر یک از مدل‌های طبقه‌بندی در اینجا جمع می‌شوند و کوچک‌ترین خطای طبقه‌بندی (روی داده‌های اعتبارسنجی) انتخاب شده و مدل منتخب روی داده‌های آموزش برازش داده می‌شود و سپس مقدار خطای طبقه‌بندی روی داده‌های آزمایشی محاسبه می‌شود. (خطای طبقه‌بندی مدل k نزدیک‌ترین همسایه روی داده‌های اعتبارسنجی استاندارد شده بدست آمده)

جدول ۱۵: جمع‌بندی خطای طبقه‌بندی برای مدل‌های مختلف انتخابی در پیش‌بینی داده‌های اعتبارسنجی

خطای طبقه‌بندی	مدل
۰.۳۳۸۱۱۵۴	رگرسیون لجستیک
۰.۳۵۸۳۳۰۶	تحلیل تفکیکی کوآدراتیک
۰.۳۳۵۱۸۱۰	تحلیل تفکیکی خطی
۰.۳۱۹۵۳۰۵	کا نزدیک‌ترین همسایه
۰.۳۲۷۰۲۹۷	درخت تصمیم ۱
۰.۳۲۷۰۲۹۷	درخت تصمیم ۲
۰.۴۴۴۷۳۴۳	بیز ساده



نمودار ۱۵: جمع‌بندی خطای طبقه‌بندی برای مدل‌های مختلف انتخابی در پیش‌بینی داده‌های اعتبارسنجی و برآورد

با مقایسه میزان خطای طبقه‌بندی بین مدل‌ها، مدل کا نزدیک‌ترین همسایه که روی داده‌های استاندارد شده‌ی برآورد نیز برازش داده شده بود، به عنوان بهترین مدل شناسایی گردید. حال مقدار این مدل روی داده‌های آموزشی استاندارد شده برازش داده شده است که مقدار خطای پیش‌بینی آن روی داده‌های آزمایشی استاندارد شده برابر با ۰.۳۲۶۶۳۷ می‌باشد.

جمع‌بندی نهایی:

در این پروژه در هر دو قسمت رگرسیون و طبقه‌بندی مدل‌های انتخابی برازش شد و مدل نهایی برای هر دو مسئله انتخاب گردید و رویدادهای آموزش برازش گردید.

اما با این وجود هنوز میزان خطا مقدار قابل توجهی است دلیل این امر این است که تعداد ویژگی‌های دخیل در برازش داده‌ها بسیار زیاد است و باید از روش‌های پیشرفته‌تر و انعطاف‌پذیرتری در حوزه یادگیری ماشین استفاده کرد. که می‌توان به الگوریتم‌های کاهش ابعاد^{۲۴} و شبکه‌های عصبی^{۲۵} اشاره داشت که شاید بتوانیم از این روش‌ها میزان خطای بسیار کمتری به دست آوریم.

²⁴ Principal component analysis and dimensionality reduction

²⁵ Artificial Neural Network