



Semiparametric Multinomial Logistic Regression Model to Classify Shape Data

Moghimbeygi, M. 

Department of Mathematics, Faculty of Mathematics and Computer Science, Kharazmi University, Tehran, Iran.

Corresponding author: M. Moghimbeygi, M.Moghimbeygi@khu.ac.ir

Received: 28/8/2022 **Revised:** 22/11/2022 **Accepted and Published Online:** 6/12/2022.

Introduction

Statistical shape analysis is one of the fields of multivariate statistics, where the main focus is on the geometric structures of objects. This analysis method is widely used in many scientific fields, such as medicine and morphology. One of the tools for diagnosing diseases or determining animal species is the images and the shapes extracted from them. Introducing methods of classifying shapes can be a solution to determine the class of each observation. Usually, in regression modelling, explanatory and dependent variables are quantitative. However, one may want to measure the relationship between an explanatory variable (with continuous values) and a dependent variable with qualitative values. One option is to use the multinomial logistic regression model. Therefore, a semiparametric multinomial logistic regression model to classify shape data is introduced in this paper.

Material and Methods

The power-divergence criterion is a measure for hypothesis testing in multinomial data. This criterion is used to define the kernel function of explanatory variables. The model is a multinomial logistic regression model based on kernel function as a function of explanatory variables and an intercept. Since the shapes' geometric structure and size play a key role in the classification of shapes, the kernel function is determined based on the shape distances. The smoothing parameter was estimated using the least square cross-validation method. Also, the estimation of model parameters was done using the neural network method.

Results and Discussion

The shape space is a manifold, but most of the methods presented in the literature for classifying shapes were done in the shape tangent space or used linear transformations. Since mapping from the manifold to linear space decreases data information, applying tangent space and linear spaces will reduce classification accuracy. Therefore, the shape space is used to classify the shape data. The performance of the model in a simulation study and two real data sets were investigated in the paper. The two real data sets used in this paper are taken from the shape package in R software. The first data set is related to schizophrenia patients and people as control, and the second one is associated with the skull of three species of apes of two sexes. The classification of these data showed an accuracy of 82% and 84%, respectively. Also, a comparison was made with the previous methods based on a real data set, which showed the proper performance of our approach compared to the other two techniques.

Conclusion

Since in the nonparametric kernel function, suitable distances of the shape space were used, the introduced method performs better than those based on Euclidean spaces. Also, the ability to use other shape distances, such as partial, full Procrustes and Riemannian distances, makes the model more flexible in classifying different types of shape data. On the other hand, size-and-shape distance can be used in the kernel function to classify data whose size plays a key role in their geometric structure. Furthermore, since few statistical distributions have been introduced in the shape space, nonparametric methods can be helpful in the analysis of shape data. However, using nonparametric methods in the shape space is time-consuming from the point of view of computer calculations.

Keywords: Logistic regression, Semiparametric regression, Shape data, Classification.

Mathematics Subject Classification (2010): 57N25, 62G08.



©The Author(s). The Publisher is Iranian Statistical Society.

This is an open access article distributed under the terms and conditions of [\(CC BY-NC 4.0\)](https://creativecommons.org/licenses/by-nc/4.0/)



رگرسیون لوژستیک چندجمله‌ای نیمه پارامتری برای رده‌بندی داده‌های شکل

میثم مقیم‌بیگی

گروه ریاضی، دانشکده علوم ریاضی و کامپیوتر، دانشگاه خوارزمی

نویسنده مسئول: میثم مقیم‌بیگی، M.Moghimbeygi@khu.ac.ir

تاریخ دریافت: ۱۴۰۱/۶/۲۷ تاریخ بازنگری: ۱۴۰۱/۹/۷ تاریخ پذیرش و انتشار: ۱۴۰۱/۹/۱۵

چکیده: در این مقاله یک مدل رگرسیون لوژستیک چند جمله‌ای نیمه پارامتری برای رده‌بندی پیکربندی‌های برجسب‌دار معرفی شده است. در مدل رگرسیونی متغیر تبیینی تابع هسته‌ای است که با استفاده از معیار توان-واگرایی به دست آمده است. همچنین متغیر پاسخ به صورت رسته‌ای بوده و رده هر پیکربندی را نشان می‌دهد. این مدل رگرسیونی نیمه پارامتری بر اساس فواصل تعریف شده در فضای شکل معرفی شده و به همین دلیل میزان رده‌بندی درست اشکال با استفاده از این روش در مقایسه با روش‌های پیشین بهبود یافته است. عملکرد این مدل در قالب یک مطالعه شبیه‌سازی مورد بررسی قرار گرفته است. در انتها نیز کاربردی از این روش در رده‌بندی دو مجموعه داده واقعی به نمایش گذاشته شد. همچنین روش ارائه شده در این مقاله با روش‌های معرفی شده در نوشتگان مقایسه گردید که نشان از عملکرد مناسب این روش در رده‌بندی پیکربندی‌ها دارد.

واژه‌های کلیدی: رگرسیون لوژستیک، رگرسیون نیمه پارامتری، داده شکل، رده‌بندی.
کد موضوع‌بندی ریاضی (۲۰۱۰): 62G08, 57N25

۱ مقدمه

تحلیل آماری شکل یکی از حوزه‌های آمار چند متغیره است که در آن تمرکز اصلی بر ساختارهای هندسی اجسام قرار دارد. "شکل یک شیء تمام اطلاعات هندسی باقی مانده از یک شیء پس از حذف اثرات مکان، مقیاس و دوران است." این تعریف رسمی برای اولین بار توسط **کنندال** (۱۹۷۷) در مورد شکل ارائه شد. پس از تعریف



شکل سؤالی که مطرح می‌شود این است که “چگونه می‌توان اطلاعات هندسی یک شیء را به صورت کمی بیان کرد؟” برای جواب به این سؤال باید گفت که یکی از روش‌های ثبت شیء استفاده از ایده ریخت‌شناسان است. آن‌ها یک شیء را با کمک نقاطی که بر سطح خارجی آن مشخص کرده‌اند معرفی می‌کنند. انتخاب این نقاط باید به گونه‌ای باشد که تا حد امکان اطلاعات کلی موجود در شیء را بیان کند. واضح است که هرچه تعداد این نقاط بیشتر باشد، اطلاعات بیشتری از شیء در دسترس خواهد بود. در تحلیل آماری شکل به این نقاط کلیدی نقاط شاخص^۱ می‌گویند. به مجموعه تمام نقاط شاخصی که برای یک شیء در نظر گرفته می‌شود، پیکربندی^۲ آن شیء گویند (درآیدن و ماردیا، ۲۰۱۶). چنانچه به نقاط شاخص نام یا برچسب اختصاص یابد در این صورت به آن‌ها پیکربندی‌های برچسب‌دار گفته می‌شود. به منظور تحلیل آماری شکل، ابتدا لازم است پیکربندی‌ها ثبت شوند. دو رویکرد اصلی در ثبت پیکربندی‌ها استفاده از نقاط شاخص و منحنی‌ها است. در این مقاله رویکرد استفاده از نقاط شاخص دنبال می‌شود. روش‌های مختلفی برای ثبت اشکال براساس نقاط شاخص وجود دارد. برای مثال ثبت نقاط با استفاده از زوایا یا اعداد مختلط را می‌توان نام برد (درآیدن و ماردیا، ۲۰۱۶). دو سیستم مختصات شناخته شده برای تحلیل آماری اشکال توسط کندال (۱۹۸۴) و بوکشتین (۱۹۸۶) معرفی شده‌اند. در حالی که مختصات اول عمدتاً برای تحلیل‌های نظری مورد استفاده قرار می‌گیرد، مختصات دوم ابزار مهمی در تحلیل اکتشافی است. منابع برای آشنایی با تحلیل آماری شکل درآیدن و ماردیا (۲۰۱۶) و بسته shape در نرم افزار R است (درآیدن، ۲۰۱۲).

برای خلاصه کردن داده‌ها رویکردهای متفاوتی در آمار وجود دارد. یکی از این روش‌ها، رده‌بندی داده‌ها بر اساس یک یا چند ویژگی آن‌ها است. برای رده‌بندی اشکال نیز روش‌های مختلفی معرفی شده است. برای مثال، دباولار و همکاران (۲۰۲۰) خوشه‌بندی داده‌های شکل که در طول زمان تغییر می‌کردند را در رده‌های مختلف بدون برچسب با استفاده از یک مدل ترکیبی جدید مطالعه کردند. مدل آن‌ها یک مسیر با استفاده از میانگین اشکال برای هر خوشه ایجاد می‌کند که از این منحنی‌ها برای خوشه‌بندی استفاده می‌شود. از نقطه نظر کندال (۱۹۷۷) شکل یک شیء معادل نقطه‌ای روی یک کره است. لازم به ذکر است که منقید فضای توپولوژی است که در هر نقطه به صورت موضعی شبیه فضای اقلیدسی است. با توجه به این‌که کره یک منیفلد است، یکی دیگر از رویکردها، خوشه‌بندی بر روی منیفلد است. این رویکرد برای خوشه‌بندی اشکال، مشابه استفاده از الگوریتم‌های تصویر غیرخطی (یانکف و کتوق، ۲۰۰۶). نمونه دیگری از این رویکرد، تصویر اشکال در فضای خطی و استفاده از مولفه‌های اصلی برای رده‌بندی اشکال با استفاده از مدل رگرسیون لوژستیک است (مقیم بیگی و نودهی، ۲۰۲۲).

رگرسیون لوژستیک رابطه بین متغیرهای توضیحی و یک متغیر پاسخ دو وجهی را توصیف می‌کند. این مدل‌بندی ابزار مناسبی برای رده‌بندی داده‌ها با پاسخ‌های دو وجهی است. تعمیم این روش را می‌توان در قالب پاسخ‌های چند وجهی به صورت رگرسیون لوژستیک چند جمله‌ای تعریف کرد. با توجه به محدودیت‌های موجود در مسائل کاربردی، محققان مدل رگرسیون لوژستیک را برای غلبه بر آن محدودیت‌ها تعمیم داده‌اند. برای مثال، لین و چن (۲۰۰۸) مدل‌های رگرسیون لوژستیک با پاسخ‌های ترتیبی را با استفاده از مدل‌های لجیت جمعی تعمیم دادند. آن‌ها همچنین

¹Landmarks²Configuration

یک روش هموارسازی خطی موضعی ناپارامتری را برای برازش به پاسخهای ترتیبی با استفاده از متغیرهای کمکی و طبقه‌ای معرفی کردند. رده‌بندی بیماران دیابتی استفاده کردند. به‌عنوان مثال دیگری از کاربردهای این روش **ریفادا و همکاران (۲۰۲۱)** مدل‌های پارامتری را به مدل‌های ناپارامتری برای استفاده در رگرسیون لوژستیک تعمیم دادند. **سیمو و همکاران (۲۰۲۰)** مدل‌های نیمه خطی تعمیم‌یافته با متغیرهای کمکی روی منیفلد ریمانی را معرفی کردند که این مدل خود شامل یک مدل خطی از متغیرهای تبیینی، یک تابع از اندازه‌ها روی منیفلد ریمانی و خط‌هایی از توزیع نرمال بود. **اکومورا و نایتو (۲۰۰۶)** روش دیگری برای رده‌بندی داده‌ها معرفی کردند. آن‌ها یک روش هموارسازی هسته برای رگرسیون چند جمله‌ای ارائه کردند. برآوردگر توابع رگرسیون در مدل آن‌ها با به حداقل رساندن معیار توان-واگرایی^۱ موضعی به دست می‌آید. این معیار کلاسی از آزمون‌های نیکویی برازش را ارائه می‌دهد.

بیشتر روش‌های ارائه شده برای تحلیل اشکال در فضای مماسی شکل انجام گرفته یا از ابزارهای خطی بهره برده‌اند. هرچند مطالعاتی نیز در فضای شکل انجام گرفته است. برای مثال، **مقیم بیگی و گل‌علی‌زاده (۱۳۹۸)** به مدل‌بندی رگرسیونی داده‌های شکل بر روی کره پرداخته‌اند و **فتوحی و گل‌علی‌زاده (۱۳۹۱)** عملکرد تحلیل ژئودزیک در فضای شکل را بهبود داده‌اند. همانگونه که اشاره شد فضای شکل یک منیفلد است. لذا بنظر می‌رسد استفاده از فضای مماسی و فضاهای خطی از دقت رده‌بندی بکاهد. از این رو در این مقاله قصد داریم از ابزارهای مناسب فضای شکل برای رده‌بندی اشکال استفاده کنیم. در این مقاله یک مدل رگرسیون لوژستیک چند جمله‌ای^۲ (MLR) بر اساس یک برآوردگر ناپارامتری از پیکربندی‌ها ارائه می‌کنیم. در این مدل، پیکربندی‌ها به صورت متغیرهای کمکی ظاهر می‌شوند که در آن معیار توان-واگرایی، بخش ناپارامتری مدل را نتیجه می‌دهد. از این رو در بخش ۲، مقدمه‌ای از تحلیل آماری شکل ارائه خواهد شد. سپس، MLR ناپارامتری برای داده‌های شکل با استفاده از معیار توان-واگرایی در بخش ۳ ارائه می‌شود. عملکرد مدل در یک مطالعه شبیه‌سازی و دو مجموعه داده واقعی در بخش ۴ مورد بررسی قرار خواهد گرفت. همچنین مقایسه‌ای نیز با روش‌های پیشین بر اساس یک مجموعه داده واقعی انجام خواهد گرفت. در انتها نیز مقاله با یک بخش نتیجه‌گیری به پایان می‌رسد.

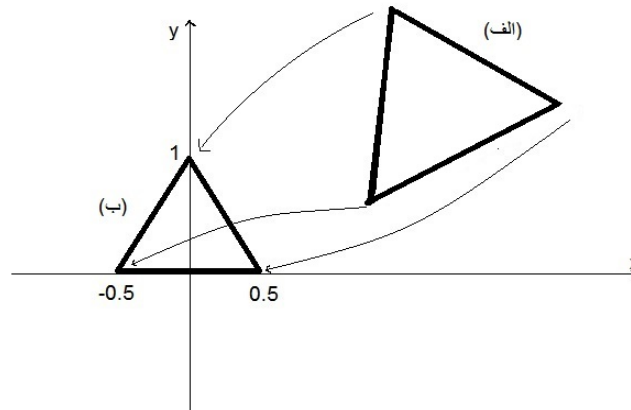
۲ مقدمه‌ای از تحلیل آماری شکل

ر تحلیل آماری شکل به مجموعه نقاط محدودی که هندسه یک شیء را تعیین می‌کند، نقاط شاخص و مجموعه تمام نقاط شاخص در نظر گرفته شده برای یک شیء را پیکربندی می‌نامند. علاوه بر این، از یک ماتریس $k \times m$ بعدی مانند X برای تعریف پیکربندی استفاده می‌شود که در آن k تعداد نقاط شاخص و m بعد پیکربندی است. طبق تعریف **کندل (۱۹۷۷)**، "شکل"، تمام اطلاعات هندسی باقیمانده از یک شیء پس از حذف اثرات مکان، مقیاس و دوران است. مثالی از این تعریف در شکل ۱ قابل مشاهده است. بنابراین، برای به دست آوردن شکل یک پیکربندی، لازم است اثرات مکان، مقیاس و دوران حذف شود. برای حذف اثر مکان می‌توان از ماتریس‌های مرکزی‌کننده مانند

¹Power-divergence

²Multinomial logistic regression

زیرماتریس هلمرت^۱، نمایش داده شده با نماد H ، استفاده کرد. (درآیدن و ماردیا، ۲۰۱۶). از این رو، ماتریس HX پیکربندی مرکزی یا هلمرتی شده نامیده می‌شود. در ادامه چند معیار فاصله که می‌توانند در تحلیل آماری داده‌های شکل مورد استفاده قرارگیرند معرفی می‌شوند.



شکل ۱. حذف اثر مکان، مقیاس و دوران مثلث متساوی الاضلاع در مختصات بوکشتین با انتقال دو نقطه شاخص مثلث (الف) به $(-0.5, 0)$ و $(0.5, 0)$. مثلث (ب) شکل و مثلث (الف) پیکربندی است.

تعریف ۱. فرض کنید X_1 و X_2 دو ماتریس پیکربندی $k \times m$ باشند. برای $j = 1, 2$ ، با فرض $Z_j = HX_j / \|HX_j\|$ که به‌عنوان قبل-شکل^۲ نامیده می‌شود، فاصله ریمانی کوتاه‌ترین فاصله بین Z_1 و Z_2 بر روی کره قبل-شکل است و به‌صورت $\rho(X_1, X_2) = \arccos(\sum_{i=1}^m \lambda_i)$ محاسبه می‌شود، که در آن λ_i ها مقادیر ویژه حاصل از تجزیه مقدار تکین عبارت زیر است

$$Z_1^T Z_2 = V \text{diag}(\lambda_1, \dots, \lambda_m) U^T. \quad (1)$$

تعریف ۲. فاصله پروکراسستس^۳ جزئی با انطباق بهینه ماتریس‌های قبل-شکل Z_1 و Z_2 بر اساس دوران حاصل می‌شود و از رابطه $d_p(X_1, X_2) = \sqrt{2(1 - \sum_{i=1}^m \lambda_i)^{1/2}}$ به‌دست می‌آید.

تعریف ۳. فاصله پروکراسستس^۴ تام که با نماد $d_F(X_1, X_2)$ نمایش داده می‌شود با انطباق بهینه ماتریس‌های قبل-شکل Z_1 و Z_2 بر اساس دو عامل مقیاس و دوران به‌دست می‌آید و به‌صورت $d_F(X_1, X_2) = [1 - \sum_{i=1}^m \lambda_i]^{1/2}$ محاسبه می‌شود.

¹ Helmert
² pre-shape
³ Procrustes
⁴ Procrustes

$(\sum_{i=1}^m \lambda_i)^2)^{1/2}$ قابل محاسبه است.

تعریف ۴. فاصله ریمانی در فضای اندازه-و-شکل^۱ که با نماد $d_S(X_1, X_2)$ نمایش داده می‌شود به صورت $d_S(X_1, X_2) = \sqrt{S_1^2 + S_2^2 - 2S_1S_2 \cos \rho(X_1, X_2)}$ است، که در آن S_1 و S_2 اندازه‌های مرکزی X_1 و X_2 هستند.

بنابه **درآیدن و ماردیا (۲۰۱۶)**، با استفاده از فاصله ریمانی، می‌توان توزیع اندازه-و-شکل را مستقیماً در فضای اندازه-و-شکل مشخص کرد. یک توزیع اندازه-و-شکل متقارن دورانی^۲، متناسب با عبارت

$$\exp\{-h\phi(d_S(X, \mu))\}, \quad (2)$$

است، که در آن $\phi(x)$ یک تابع مثبت است. یک انتخاب مناسب برای این تابع $\phi(x) = x^2$ است که در آن ثابت نرمال‌ساز تابع چگالی، وابسته به پارامتر μ نیست. ملاک RMS معیاری است که در تحلیل آماری شکل برای محاسبه همگنی استفاده می‌شود. این معیار تعدیلی از اندازه d_S است که به صورت $RMS = d_S/\sqrt{k}$ تعریف می‌شود که در آن k تعداد نقاط شاخص است. مقدار کوچک RMS به معنای همگن بودن مجموعه داده شکل است.

تعریف ۵. تجزیه و تحلیل جزئی پروکروستس معمولی^۳ (POPA) با استفاده از تبدیلات اقلیدسی جابجایی و دوران دو پیکربندی با حداقل رساندن عبارت

$$\|X_2 - X_1\Gamma - \gamma_k\gamma^T\|^2, \quad (3)$$

تعریف می‌شود. می‌توان نشان داد مقادیر بهینه برای Γ و γ با استفاده از تجزیه مقدار تکین (۱) به صورت $\hat{\gamma} = 0$ و $\hat{\Gamma} = UV^T$ به دست می‌آید، که در آن U و V ماتریس‌های متعامد ویژه ($SO(m)$) تعریف شده در معادله (۱) هستند (**درآیدن و ماردیا، ۲۰۱۶**). بنابراین انطباق جزئی X_1 بر روی X_2 به صورت $X_2^P = X_1\hat{\Gamma} + \gamma_k\hat{\gamma}^T$ است.

۳ رگرسیون لوژستیک چندجمله‌ای نیمه پارامتری

رگرسیون نیمه پارامتری در واقع ترکیبی از مدل‌های رگرسیونی پارامتری و ناپارامتری است. این مدل‌ها اغلب زمانی استفاده می‌شوند که مدل ناپارامتری عملکرد خوبی نداشته باشد یا زمانی که محقق می‌خواهد از یک مدل پارامتری استفاده کند اما ساختار داده‌ها به گونه‌ای است که مدل‌بندی آن‌ها با استفاده از مدل‌های شناخته شده پارامتری امکان‌پذیر نیست. فرض کنید X_i یک ماتریس پیکربندی و یک متغیر (تیبینی) با متغیر پاسخ دو وجهی $Y_i \in \{0, 1\}$ برای

¹Size-and-shape

²Rotationally symmetric

³Partial Ordinary Procrustes Analysis

$i = 1, 2, \dots, n$ باشد. در زیر بخش بعد قصد داریم یک مدل MLR برای رده‌بندی پیکربندی‌ها بر اساس عضویت آن در گروه‌های متناظرشان ارائه کنیم. برای این منظور، ابتدا تبدیل ناپارامتری پیکربندی‌ها را بر اساس معیار توان-واگرایی ارائه می‌کنیم. سپس یک مدل لوژستیک نیمه پارامتری در فضای اندازه-و-شکل معرفی می‌کنیم. مدل معرفی شده توانایی تفکیک پیکربندی‌ها را براساس معیارهای آماری مناسب دارد.

۳.۱ تبدیل ناپارامتری پیکربندی‌ها

برای $i = 1, \dots, K$ ، فرض کنید x_i ها متغیرهای تبیینی Y_i ها متغیرهای پاسخ و $Y_i = (Y_{i1}, \dots, Y_{ir})^T$ دارای توزیع چندجمله‌ای $MN(p_1(x_i), \dots, p_r(x_i); N(x_i))$ باشد که در آن $N_i = N(x_i) = \sum_{j=1}^r Y_{ij}$ و $Y_{ij} = Y_j(x_i)$. به منظور ارزیابی برآوردگرها، روش‌های متفاوتی وجود دارد. یکی از آن‌ها را که می‌توان برای آزمایش‌هایی که از توزیع چندجمله‌ای تبعیت می‌کنند استفاده کرد، آزمون نیکویی برازش است. این آزمون بر اساس معیار توان-واگرایی به دست آمده است که توسط کرسی و رید (۱۹۸۴) برای ارزیابی برآوردگر به صورت

$$I_\lambda(\mathbf{p}_i : \mathbf{q}) = \frac{1}{\lambda(1+\lambda)} \sum_{j=1}^r p_{ij} \left\{ \left(\frac{p_{ij}}{q_j} \right)^\lambda - 1 \right\},$$

به دست آورده شده است، که در آن $\lambda \in \mathbb{R}$ ، $\mathbf{p}_i = (p_{i1}, \dots, p_{ir})$ و $\mathbf{q} = (q_1, \dots, q_r)$ پارامترهای توزیع احتمال هستند، معرفی کردند. مشابه روش معرفی شده توسط اکومورا و نایتو (۲۰۰۶)، ما یک نسخه موضعی از توان-واگرایی را با استفاده از تابع هسته ناپارامتری $K_h(\cdot)$ به صورت

$$L_\lambda(\mathbf{q}, \gamma) = \frac{1}{\lambda(1+\lambda)} \sum_{i=1}^K K_h(x_i - x) \left[\sum_{j=1}^r p_{ij} \left\{ \left(\frac{p_{ij}}{q_j} \right)^\lambda - 1 \right\} + \gamma \left(1 - \sum_{j=1}^r q_j \right) \right],$$

در نظر می‌گیریم که γ ضریب لاگرانژ است. پس برآوردگر $q_j(x)$ با به حداقل رساندن $L_\lambda(\mathbf{q}, \gamma)$ به صورت

$$\hat{q}_{j\lambda}(x; h) = \frac{\left\{ \sum_{i=1}^K K_h(x_i - x) p_{ij}^{\lambda+1} \right\}^{1/(\lambda+1)}}{\sum_{j=1}^r \left\{ \sum_{i=1}^K K_h(x_i - x) p_{ij}^{\lambda+1} \right\}^{1/(\lambda+1)}}, \quad (4)$$

به دست می‌آید. به وضوح $0 \leq \hat{q}_{j\lambda}(x; h) \leq 1$ و به ازای هر x تساوی $\sum_{j=1}^r \hat{q}_{j\lambda}(x; h) = 1$ برقرار است. توجه شود که اگر $\lambda = 0$ ، آنگاه برآوردگر حاصل همان برآوردگر نادارایا-واتسن^۱ است که به صورت

$$\hat{q}_{j,0}(x; h) = \frac{\sum_{i=1}^K K_h(x_i - x) p_{ij}}{\sum_{i=1}^K K_h(x_i - x)},$$

^۱Nadaraya-Watson

نوشته می‌شود. در این مقاله از برآوردگر (۴) به‌ازای $\lambda = 0$ برای رده‌بندی پیکربندی‌ها استفاده می‌شود. برای استفاده از مدل معرفی شده در داده‌های شکل، ما r دسته از پیکربندی‌های X_{ij} را برای $1 \leq j \leq r$ و $1 \leq i \leq K$ در نظر می‌گیریم. همچنین برای تعریف معیار خود از آماره آزمون گودالز^۱ استفاده می‌کنیم. این آماره به‌منظور آزمون برابری میانگین اشکال در دو جامعه مورد استفاده قرار می‌گیرد. از بخش ۹ در **درآیدن و ماردیا (۲۰۱۶)**، با فرض X_{ij} به‌عنوان یک متغیر تصادفی $k \times m$ از $N_{km}(\mu_{ij}, \sigma^2 I_{km})$ ، فاصله پروکروسس تام X_{ij} از هر پیکربندی ثابت μ_j برابر $d_F^*(X_{ij}, \mu_j)$ است که تقریباً مشابه χ^2 با درجه آزادی $km - m(m+1)/2 - 1$ توزیع شده است. برای r دسته از پیکربندی‌ها،

$$B_{ij}(X_{ij}, \mu_j) = \frac{d_F^*(X_{ij}, \mu_j)}{\sum_{j=1}^r d_F^*(X_{ij}, \mu_j)}$$

یک متغیر تصادفی از توزیع بتا است. چون فاصله هر شکل با اعضای هر رده می‌تواند نقش مهمی در تعیین رده آن شکل داشته باشد لذا در این مقاله از فواصل معرفی شده در بخش ۲ برای رده‌بندی اشکال استفاده شده است. به‌عبارت دیگر $B_{ij}(X_{ij}, X)$ ‌های بدست آمده از پیکربندی X و سایر پیکربندی‌ها نقش اصلی را در رده‌بندی بازی می‌کنند. برآوردگر q_j برای داده‌های شکل را می‌توان با جایگذاری $B_{ij}(X_{ij}, X)$ به‌جای p_{ij} در (۴) به‌صورت

$$\hat{q}_{j\lambda}(X; h) = \frac{\{\sum_{i=1}^K K_h(X_{ij} - X) d_F^{\lambda+1}(X_{ij}, X)\}^{1/(\lambda+1)}}{\sum_{j=1}^r \{\sum_{i=1}^K K_h(X_{ij} - X) d_F^{\lambda+1}(X_{ij}, X)\}^{1/(\lambda+1)}}, \quad (5)$$

نوشت. از آنجایی که مقدار کوچک $d_F^*(X_{ij}, X)$ به معنای این است که دو پیکربندی X_{ij} و X در یک رده یکسان هستند، بنابراین $\hat{q}_{j\lambda}(X; h)$ احتمال عدم تعلق X به رده j تعریف می‌شود. برای اشکالی که به هم شبیه هستند (فاصله ریمانی کمتری دارند)، داریم $d_p = d_F + O(d_F^*)$ و $\rho = d_F + O(d_F^*)$ (**درآیدن و ماردیا، ۲۰۱۶**).

در برخی از مسائل واقعی اندازه پیکربندی‌ها نقش مهمی در رده‌بندی آن‌ها دارد. برای مثال حجمه یک انسان بالغ از نظر شکلی مشابه حجمه یک کودک است و تنها تفاوت می‌تواند اندازه حجمه‌ها باشد. لذا در این مسائل استفاده از d_p که اندازه را نیز لحاظ می‌کند، می‌تواند معیار مناسب‌تری در مدل معرفی شده باشد. از این رو، d_p می‌تواند به‌جای d_F در معادله (۵) استفاده شود. برای انتخاب تابع هسته، از توزیع (۲) با تابع $\phi(x) = x^2$ استفاده می‌کنیم. انتخاب پهنای باند در تحلیل ناپارامتری از اهمیت بالایی برخوردار است. راهکارهای زیادی برای انتخاب پهنای باند وجود دارد. برای مثال می‌توان به روش اعتبارسنجی متقابل حداقل مربعات و اعتبارسنجی متقابل حداکثر درست‌نمایی اشاره کرد (**دوین، ۱۹۷۶**). از آنجایی که هدف اصلی این مقاله ارائه روشی برای رده‌بندی پیکربندی‌ها است، پهنای باند برآورد شده باید بتواند حداکثر رده‌بندی صحیح را انجام دهد. بنابراین، برای برآورد پارامتر پهنای باند از روش اعتبارسنجی متقابل حداقل مربعات استفاده می‌شود تا فاصله بین احتمال برآورد شده و متغیر پاسخ

¹Goodalls

کمترین مقدار را داشته باشد. اصلاحی از این معیار به صورت

$$CV(h) = \sum_{j=1}^r \sum_{i=1}^K (y_{ij} - q_{j,\lambda}(-X_{ij}, h))^2,$$

معرفی می‌شود که $q_{j,\lambda}(-X_{ij}, h)$ برآوردگری است که در آن مشاهده X_{ij} در هر مرحله کنار گذاشته شده است (واسرمن، ۲۰۰۶؛ ریفادا و همکاران، ۲۰۲۱). برای انتخاب پهنای باند بهینه براساس این معیار لازم است هدف محقق مبتنی بر به حداقل رساندن $CV(h)$ باشد. از آنجایی که در تحلیل‌های ناپارامتری معمولاً برآورد پهنای باند فرم بسته‌ای بر اساس حداقل کردن $CV(h)$ ندارد، لذا برآورد h معمولاً با استفاده از الگوریتم‌های بهینه‌سازی و به صورت عددی انجام می‌گیرد. از آنجا که B_{ij} ها با فرض این‌که پیکربندی از توزیع چند متغیره نرمال به دست می‌آید، انتظار می‌رود که برآوردگر معرفی شده (۵) یک برآوردگر اریب باشد. این مقاله از رگرسیون لوژیستیک خطی ساده برای حذف این سوگیری‌ها استفاده می‌کند.

۳.۲ مدل لوژیستیک

در مدل‌بندی رگرسیون، مدلی برای پیش‌بینی متغیر وابسته بر اساس متغیر تبیینی ایجاد می‌شود. معمولاً در مدل ایجاد شده، هم متغیرهای تبیینی و هم متغیرهای وابسته کمی هستند. با این حال، ممکن است بخواهیم رابطه بین یک متغیر تبیینی (با مقادیر پیوسته) و یک متغیر وابسته با مقادیر کیفی را اندازه‌گیری کنیم. روش استاندارد رگرسیون خطی در این مورد کارا نخواهد بود که در این موارد می‌توان از MLR استفاده کرد. MLR یک روش برای رده‌بندی داده‌ها به دو یا چند رده مختلف است. در رگرسیون خطی استاندارد، مقدار مورد انتظار متغیر پاسخ Y ترکیبی خطی از متغیرهای تبیینی است. یک نسخه نیمه پارامتری از مدل خطی ساده را می‌توان به صورت $\mathbb{E}(Y|X = x) = \alpha + \beta q(x)$ تعریف کرد که α پارامتر عرض از مبدا و $q(x)$ یک تابع ریاضی است. هنگامی که برآمد Y دو وجهی است، یک انتخاب استفاده از رگرسیون لوژیستیک است. گزینه‌های دیگر برای تحلیل این مدل‌ها را می‌توان در [کاکس و اسنیل \(۱۹۸۹\)](#) یافت. با استفاده از توزیع لوژیستیک، می‌توان میانگین Y به شرط x را به صورت

$$\pi(x) = \mathbb{E}(Y|x) = \frac{e^{\alpha + \beta q(x)}}{1 + e^{\alpha + \beta q(x)}},$$

تعریف کرد که در آن $\pi(x)$ به عنوان احتمال موفقیت نیز شناخته می‌شود. یک تبدیل ساده از $\pi(x)$ که نقش کلیدی در این مقاله ایفا می‌کند تبدیل لوچیت $\ln\left[\frac{\pi(x)}{1-\pi(x)}\right] = \alpha + \beta q(x)$ است. برای محاسبه MLR، ابتدا p متغیر کمکی با r رده متفاوت (با پاسخ‌های اسمی Y) و بردار x از متغیرهای تبیینی به طول $p + 1$ را در نظر بگیرید. با

انتخاب r -امین رده به عنوان یک محور داریم

$$g_m(\mathbf{x}) = \ln \frac{Pr(Y = m|\mathbf{x})}{Pr(Y = r|\mathbf{x})} = \mathbf{x}^T \beta_m, \quad m = 1, \dots, r-1,$$

که در آن برداری به طول $p+1$ و $\mathbf{x} = (1, x_1, \dots, x_p)$ و $\beta_m^T = (\beta_{m0}, \beta_{m1}, \dots, \beta_{mp})$ بردار پارامترها است. تعمیمی از احتمال شرطی در r رده ممکن به صورت

$$\pi_j(\mathbf{x}) = Pr(Y = j|\mathbf{x}) = \frac{e^{g_j(\mathbf{x})}}{\sum_{k=0}^r e^{g_k(\mathbf{x})}}, \quad j = 0, \dots, r-1,$$

است. برای تابع درستنمایی، از r متغیر دو وجهی U_i که مقادیر 0 یا 1 را می‌گیرند استفاده می‌شود. بدین صورت که اگر $Y = i$ ، آن‌گاه $U_i = 1$ و برای $i \neq j$ ، $U_j = 0$ باشد. در این صورت بدون در نظر گرفتن مقدار U_i داریم $\sum_{k=1}^r U_k = 1$. بنابراین تابع درستنمایی برای n نمونه مستقل از مشاهدات به صورت

$$L(\beta) = \prod_{i=1}^n \prod_{k=1}^r \pi_k(\mathbf{x}_i)^{u_{ki}},$$

است. با لگاریتم‌گیری از تابع درستنمایی و در نظر گرفتن شرط $\sum_{k=1}^r U_k = 1$ ، داریم:

$$\ell(\beta) = \sum_{i=1}^n \sum_{k=1}^r [u_{ki} \pi_k(\mathbf{x}_i)] - \sum_{i=1}^n \ln \left[\sum_{k=1}^r e^{g_k(\mathbf{x}_i)} \right].$$

با مشتق‌گیری از تابع $\ell(\beta)$ نسبت به پارامترهای مدل، برای $k = 1, \dots, r$ و $l = 0, \dots, p$ داریم

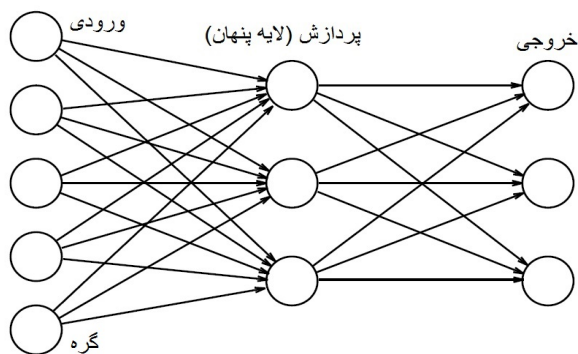
$$\frac{\partial \ell(\beta)}{\partial \beta_{kl}} = \sum_{i=1}^n x_{li} (u_{ki} - \pi_k(\mathbf{x}_i)),$$

که در آن $x_{0i} = 1$. چون حل تحلیلی معادله $\partial \ell(\beta) / \partial \beta_{kl} = 0$ میسر نیست یا بسیار دشوار است، باید به روش‌های عددی متوصل شد. شبکه عصبی مصنوعی در واقع یک مدل محاسباتی است که از روش‌های محاسباتی جدید برای یادگیری ماشینی، نمایش دانش و در انتها اعمال دانش به دست آمده در جهت پیش‌بینی پاسخ‌های خروجی از سامانه‌های پیچیده استفاده می‌کند. هر شبکه عصبی از سه لایه تشکیل شده است: ورودی، خروجی و پردازش. هر لایه شامل گروهی از سلول‌های عصبی است که در حالت کلی این سلول‌ها با تمام نوروهای لایه‌های دیگر در ارتباط هستند. یکی از شبکه‌های عصبی که در مدل‌های رگرسیونی مورد استفاده قرار می‌گیرد شبکه عصبی پیش‌خور¹

¹Feedforward

است. شبکه عصبی پیشخور اولین و ساده‌ترین نوع شبکه عصبی است. در این شبکه اطلاعات تنها از در یک مسیر حرکت می‌کنند. به عبارت دیگر اطلاعات با شروع از گره‌های ورودی و گذر از لایه‌های پنهان به سمت گره‌های خروجی می‌روند. نمایی از یک شبکه عصبی پیش‌خور در شکل ۲ قابل مشاهده است.

در شبکه‌های عصبی از یک تابع هزینه برای برآورد پارامترها استفاده می‌شود. در مدل‌های رگرسیونی معمولاً منفی لگاریتم تابع درستنمایی را به عنوان تابع هزینه در نظر می‌گیرند. سپس با استفاده از روش گرادینت کاهشی به برآورد پارامترها می‌پردازند. به این معنی که گرادینت تابع را حساب کرده، کمی در خلاف جهت آن حرکت کرده و این کار را آنقدر ادامه داده تا تابع هزینه خیلی کوچک شود.



شکل ۲. شبکه عصبی پیش‌خور همراه با سه لایه ورودی، پنهان و خروجی به ترتیب دارای ۵، ۳ و ۳ گره.

برای $i = 1, \dots, K$ و $j = 1, \dots, r$ ، فرض کنید ماتریس‌های پیکربندی X_{ij} به عنوان متغیرهای تبیینی متناظر با پاسخ‌های اسمی y_{ij} در r رده مختلف باشند. برای یک مدل لگ-خطی چند جمله‌ای با r رده، یک شبکه عصبی با r خروجی، منفی لگاریتم تابع درستنمایی شرطی به صورت

$$E = - \sum_{i=1}^K \sum_{j=1}^r t_{ij} \log \pi_{ij}, \quad \pi_{ij} = \frac{e^{y_{ij}}}{\sum_{j=1}^r e^{y_{ij}}}$$

است، که در آن t_{ij} هدف و y_{ij} خروجی برای ورودی i -ام است. یک منبع مناسب برای آشنایی با شبکه‌های عصبی ونابلز و ریپلی (۲۰۰۲) است. همچنین نحوه برآورد پارامترها در مدل MLR با استفاده از شبکه عصبی در نرم افزار R، در این منبع موجود است. در این مقاله از بسته `nnet` برای برازش MLR به منظور رده‌بندی پیکربندی‌ها استفاده می‌شود. پس از مدل‌بندی رگرسیونی داده‌ها با استفاده از مدل MLR لازم است عملکرد این مدل رگرسیونی ارزیابی شود. لذا در ادامه چند معیار شناخته شده که در ارزیابی مدل‌های MLR به کار می‌رود، معرفی می‌شود. اولین

ملاک، معیار اطلاع آکائیک^۱ (AIC) است. این معیار تخمین گر خطای پیش بینی و در نتیجه کیفیت نسبی مدل های آماری برای مجموعه ای از داده ها است (آکائیک، ۱۹۷۳). با توجه به مجموعه ای از مدل های کاندید برای داده ها، AIC کیفیت هر مدل را نسبت به هر یک از مدل های دیگر تخمین می زند. این معیار برای انتخاب مدل به صورت

$$AIC = -2\ell(\hat{\beta}) + 2(p + 1)$$

تعریف می شود. مدل ترجیحی مدلی است که دارای حداقل مقدار AIC باشد. توسعه ای از معیار AIC به نام معیار اطلاع بیزی^۲ (BIC) توسط شوارز (۱۹۷۸) به صورت

$$BIC = -2\ell(\hat{\beta}) + \ln(n)(p + 1)$$

ارائه شد. معیار دیگری که برای ارزیابی مدل ها مورد استفاده قرار می گیرد، معیار ضریب تعیین است. ضریب تعیین در واقع نسبت تغییرات در متغیر وابسته است که قابل پیش بینی از طریق متغیر(های) مستقل است. این ضریب با R^2 نشان داده می شود و معمولاً در رگرسیون خطی از آن استفاده می شود. در رگرسیون لوژستیک ضرایب تعیین های متنوعی معرفی شده است که به عنوان مثال می توان به ضرایب تعیین مک فادن^۳ (MF)، کاکس اسنیل^۴ (CS) و ناگلکرک^۵ (Na) اشاره کرد، که به ترتیب به صورت

$$R_{MF}^2 = 1 - \frac{\ell(\hat{\beta})_k}{\ell(\hat{\beta})_{null}}, \quad R_{CS}^2 = 1 - \left[\frac{\ell(\hat{\beta})_{null}}{\ell(\hat{\beta})_k} \right]^{2/n}, \quad R_{Na}^2 = \frac{1 - \left[\frac{\ell(\hat{\beta})_{null}}{\ell(\hat{\beta})_k} \right]^{2/n}}{1 - (\ell(\hat{\beta})_{null})^{2/n}}$$

تعریف می شوند. در این ضرایب $\ell(\hat{\beta})_{null}$ مقدار تابع درستنمایی تحت فرض $\beta_{ij} = 0$: برای H_0 و $j = 1, \dots, p$ و $i = 1, \dots, r - 1$ مقدار تابع درستنمایی تحت فرض تمامی فضای پارامتری است.

۴ مطالعه شبیه سازی و تحلیل داده های واقعی

در این بخش قصد داریم داده های شبیه سازی شده و واقعی را رده بندی کنیم. لذا در زیربخش اول، داده ها را شبیه سازی کرده، بر اساس مدل MLR رده بندی و عملکرد روش معرفی شده را ارزیابی می کنیم. در بخش دوم نیز دو مجموعه داده واقعی که مرتبط با ماهیت پیکربندی (در مفهوم شکل) هستند را رده بندی می کنیم.

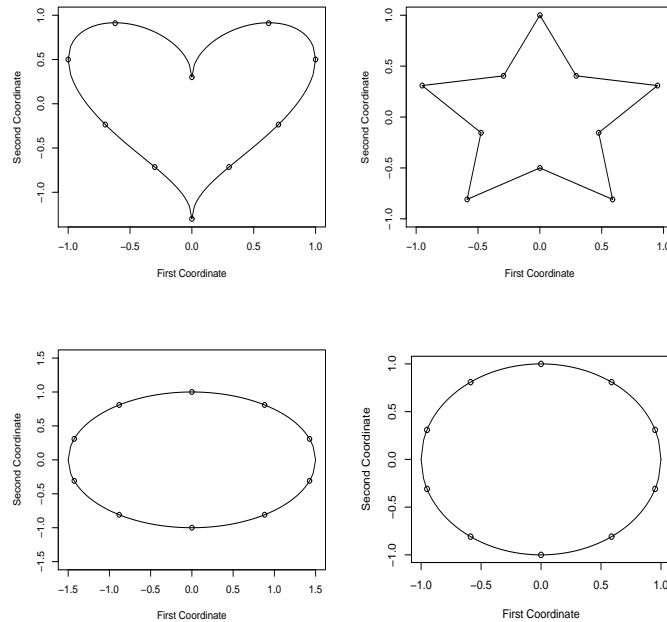
¹Akaike information criterion

²Bayesian information criterion

³McFadden

⁴CoxSnell

⁵Nagelkerke



شکل ۳. نمایی از شکل‌های مرجع ستاره، قلب، دایره و بیضی استفاده شده در مطالعه شبیه‌سازی

۴.۱ مطالعه شبیه‌سازی

به منظور مطالعه شبیه‌سازی، از روشی که در نیبل و گل‌علی‌زاده (۲۰۱۶) توضیح داده شده است برای شبیه‌سازی پیکربندی‌ها استفاده می‌کنیم. مدل به صورت

$$X_{ij} = \beta_{ij}\mu_j\Gamma_{ij} + \gamma_{ij}^T + \varepsilon_{ij} \quad i = 1, \dots, K, \quad j = 1, \dots, r, \quad (6)$$

تعریف می‌شود که در آن μ_j ، j -امین پیکربندی مرجع و ε_{ij} یک ماتریس تصادفی $k \times m$ به عنوان منبع خطا است که از توزیع نرمال چند متغیره تولید می‌شود. ما از چهار پیکربندی مرجع، ستاره، قلب، دایره و بیضی با ده نقطه شاخص در دو بعد استفاده کردیم. نمایشی از اشکال مرجع در شکل ۳ نشان داده شده است. ابتدا ۱۰، ۳۰ یا ۵۰ پیکربندی از مدل (۶) برای هر یک از شکل‌های مرجع شبیه‌سازی می‌کنیم. به عنوان مثال در حالت ۱۰ پیکربندی، مجموعاً ۴۰ پیکربندی داریم که از مدل (۶) با ۴ میانگین مختلف شبیه‌سازی شده‌اند. سپس با استفاده از روش معرفی شده با استفاده از پهنای باند برآورد شده به رده‌بندی داده‌ها می‌پردازیم. این عمل را ۱۰۰ بار تکرار می‌کنیم.

جدول ۱. میانگین و انحراف معیار ضریب تعیین نما (مکفادن، کاکس اسنیل و ناگلکرک) و معیارهای برازش مدل (AIC، BIC، OCP) بر اساس ۱۰۰ بار تکرار از مدل (۶) برای مقادیر مختلف σ .

| تعداد پیکربندی | σ | برآورد | ضریب تعیین نما | | | معیارهای برازش مدل | | |
|----------------|----------|--------------|----------------|------------|---------|--------------------|--------|--------|
| | | | مکفادن | کاکس اسنیل | ناگلکرک | AIC | BIC | OCP |
| ۱۰ | ۰/۱ | میانگین | ۰/۹۹۳۹ | ۰/۸۳۶۴ | ۰/۹۹۸۸ | ۲۵/۳۵۹ | ۶۳/۴۷۱ | ۰/۹۹۶۲ |
| | ۰/۳ | انحراف معیار | ۰/۸۶۷۲ | ۰/۹۰۷۵ | ۰/۹۶۸۰ | ۳۸/۷۲۸ | ۲۸/۹۰۳ | ۰/۰۸۴ |
| | ۰/۵ | میانگین | ۰/۷۰۶۰ | ۰/۸۵۷۱ | ۰/۹۱۴۲ | ۸۹/۲۲۲ | ۱۲۷/۳۳ | ۰/۸۰۶۲ |
| | ۰/۳ | انحراف معیار | ۰/۸۴۱ | ۰/۹۰۲۰۱ | ۰/۹۲۱۵ | ۹۳/۲۵۷ | ۹۳/۲۵۷ | ۰/۵۱۴ |
| | ۰/۱ | انحراف معیار | ۰/۵۹۲ | ۰/۷۳۳ | ۰/۷۳۳ | ۱۳/۱۲۱ | ۱۳/۱۲۱ | ۰/۵۲۵ |
| ۳۰ | ۰/۱ | میانگین | ۰/۹۹۷۶ | ۰/۸۳۷۱ | ۰/۹۹۹۵ | ۲۴/۷۹۳ | ۶۹/۳۹۳ | ۰/۹۹۸۳ |
| | ۰/۳ | انحراف معیار | ۰/۹۱۷۱ | ۰/۹۲۱۰ | ۰/۹۸۲۴ | ۵۱/۵۸۳ | ۲۵/۰۲۱ | ۰/۰۵۳ |
| | ۰/۵ | میانگین | ۰/۷۱۵۵ | ۰/۸۶۲۱ | ۰/۹۱۹۶ | ۱۱۸/۶۶ | ۱۲۳/۱۰ | ۰/۲۰۸ |
| | ۰/۳ | انحراف معیار | ۰/۸۶۲۱ | ۰/۹۱۹۶ | ۰/۹۸۲۴ | ۱۱۸/۶۶ | ۱۲۳/۱۰ | ۰/۲۰۸ |
| | ۰/۱ | انحراف معیار | ۰/۲۷۶ | ۰/۱۰۴ | ۰/۱۱۰ | ۹/۱۶۶۵ | ۹/۱۶۶۵ | ۰/۳۵۰ |
| ۵۰ | ۰/۱ | میانگین | ۰/۹۹۸۹ | ۰/۸۳۷۳ | ۰/۹۹۹۸ | ۲۴/۵۹۹ | ۷۷/۳۷۲ | ۰/۹۹۹۵ |
| | ۰/۳ | انحراف معیار | ۰/۹۱۴۱ | ۰/۹۲۰۶ | ۰/۹۸۲۰ | ۷۱/۶۵۷ | ۱۸۹/۳۷ | ۰/۰۱۶ |
| | ۰/۵ | میانگین | ۰/۷۳۵۱ | ۰/۸۶۹۳ | ۰/۹۲۷۲ | ۱۷۰/۹۲ | ۱۲۴/۴۳ | ۰/۹۵۱۵ |
| | ۰/۳ | انحراف معیار | ۰/۸۶۹۳ | ۰/۹۲۷۲ | ۰/۹۸۲۰ | ۱۷۰/۹۲ | ۱۲۴/۴۳ | ۰/۹۵۱۵ |
| | ۰/۱ | انحراف معیار | ۰/۳۱۳ | ۰/۱۱۴ | ۰/۱۲۱ | ۱۷/۳۵۴ | ۱۷/۳۵۴ | ۰/۲۰۵ |

نتایج حاصل از رده بندی در قالب معیارهای ضریب تعیین-نمای^۱ مکفادن، کاکس اسنیل و ناگلکرک و همچنین معیارهای AIC، BIC و درصد پیش بینی صحیح کلی (OCP) به دست آمد. میانگین و انحراف معیار این نتایج در جدول ۱ برای مقادیر مختلف از σ محاسبه شده است. همانگونه که قابل مشاهده است با افزایش میزان σ میانگین در تمامی ضرایب تعیین نما کاهش می یابد در صورتی که معیارهای AIC و BIC افزایش چشم گیری دارند. همچنین مقادیر انحراف معیار نیز افزایش می یابد. خلاصه ای از مقادیر جدول ۱ در قالب نمودارها در شکل ۴ به نمایش گذاشته شده است. اصلی ترین معیار برای بررسی عملکرد مدل را می توان معیار OCP در نظر گرفت که میزان رده بندی درست را نشان می دهد. همانگونه که انتظار می رود با افزایش σ مقدار OCP کاهش می یابد. با افزایش تعداد پیکربندی ها در هر رده میزان دقت در رده بندی بر اساس معیار OCP بهبود می یابد.

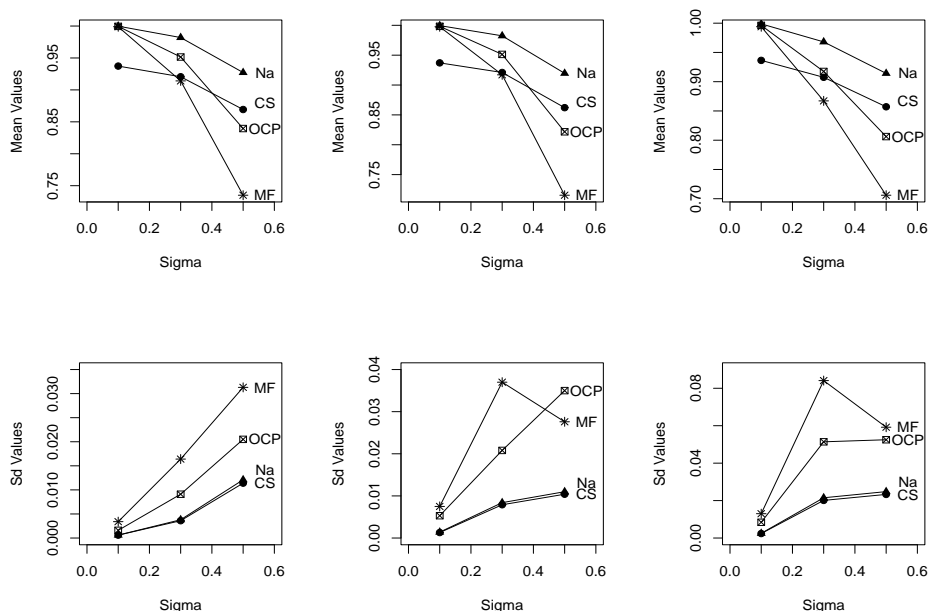
۴.۲ کاربرد

دو مجموعه داده واقعی مورد استفاده در این مقاله از بسته shapes در نرم افزار R گرفته شده اند:

اسکیزوفرنی (Schi): برای ثبت این داده ها، ابتدا تصویربرداری رزونانس مغناطیسی (MRI) در نزدیکی میدساژیتال^۲ از ۱۴ بیمار اسکیزوفرنی و ۱۴ فرد کنترل انجام گرفت. سپس برش های دو بعدی از تصاویر در نظر گرفته و ۱۳ نقطه

¹Pseudo-R-Square

²Midsagittal



شکل ۴. از راست به چپ: نمودارهای میانگین (بالا) و انحراف معیار (پایین) ضرایب تعیین‌نمای مک‌فادن، کاکس‌اسنیل، ناگلکرک، و معیار OCP برای ۱۰، ۳۰ و ۵۰ پیکربندی براساس ۱۰۰ بار تکرار شبیه‌سازی.

شاخص آن‌ها ثابت شد. این داده‌ها در دو بعد هستند (بوکشتین، ۱۹۹۶).

میمون‌ها: داده‌ها مربوط به پیکربندی جمجمه ۱۶۷ میمون بالغ، شامل ۳۰ گوریل ماده، ۲۹ گوریل نر (gorf)، ۲۶ شامپانزه ماده، ۲۸ شامپانزه نر (panf)، ۲۴ اورانگوتان ماده (pongof) و ۳۰ اورانگوتان نر (pongom) است (اهیگینز، ۱۹۸۹؛ اهیگینز و درایدن، ۱۹۹۳).

برای رده‌بندی این دو مجموعه داده، ابتدا پارامتر پهنای باند h را با استفاده از معیار CV برآورد می‌کنیم. با استفاده از \hat{h} ، مدل MLR را برازش داده و پیکربندی‌ها را رده‌بندی می‌کنیم. نتایج در جدول ۲ نشان داده شده است. توجه داشته باشید که برای محاسبه معیارهای گزارش شده در جدول ۲ هر مشاهده را کنار گذاشته و سپس رده آن مشاهده را با استفاده از سایر داده‌ها پیش‌بینی کردیم. برای ۲۸ داده‌ی اسکیزوفرنی، رده مربوط به ۱۲ فرد در رده کنترل و ۱۱ فرد در رده بیمار به‌درستی پیش‌بینی شد. همچنین رده ۲ فرد سالم و ۳ فرد بیمار اشتباه پیش‌بینی شد. برای مجموعه داده میمون‌ها نیز ماتریس درهم‌ریختگی^۱ آن‌ها ایجاد شد. این ماتریس در جدول ۳ قابل مشاهده است. مقادیر داخل جدول تعداد پیش‌بینی‌ها بر اساس رده واقعی هر میمون و رده پیش‌بینی شده را نشان می‌دهد.

¹Confusion matrix

جدول ۲. معیارهای ضریب تعیین نما مک فادن، کاکس اسنیل و ناگلکرک و خلاصه ای از برازش رگرسیون لوژیستیک شامل AIC، BIC و OCP برای داده های اسکیزوفرنی و میمون ها.

| \hat{h} | خلاصه مدل | | | ضریب تعیین نما | | | مجموعه داده اسکیزوفرنی میمون ها |
|-----------|-----------|-------|------|----------------|------------|---------|---------------------------------|
| | OCP | BIC | AIC | ناگلکرک | کاکس اسنیل | مک فادن | |
| ۲۵۷ | ۸۲ | ۳۱۱۶ | ۴۰۴۹ | ۰٫۳۲ | ۰٫۲۴ | ۰٫۲۰ | |
| ۰٫۵۸ | ۸۴ | ۲۱۱٫۴ | ۳۵۱۶ | ۰٫۹۷ | ۰٫۹۵ | ۰٫۸۲ | |

جدول ۳. ماتریس درهم ریختگی گروه های میمون ها با رگرسیون لوژیستیک چند جمله ای نیمه پارامتری.

| رده بندی درست (درصد) | پیش بینی | | | | | | مشاهده |
|----------------------|----------|--------|------|------|------|------|--------------|
| | pongom | pongof | panm | panf | gorm | gorf | |
| ۱۰۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۳۰ | gorf |
| ۱۰۰ | ۰ | ۰ | ۰ | ۰ | ۲۹ | ۰ | gorm |
| ۶۱٫۵ | ۰ | ۰ | ۱۰ | ۱۶ | ۰ | ۰ | panf |
| ۶۳٫۰ | ۰ | ۲ | ۱۷ | ۱۰ | ۰ | ۰ | panm |
| ۹۰٫۸ | ۲ | ۲۰ | ۱ | ۰ | ۰ | ۰ | pongof |
| ۹۳٫۳ | ۲۸ | ۲ | ۰ | ۰ | ۰ | ۰ | pongom |
| ۸۴ | ۱۸ | ۱۳٫۲ | ۱۶٫۲ | ۱۵٫۶ | ۱۷٫۴ | ۱۸ | مجموع (درصد) |

جدول ۴. معیار RMS برای هر دو گروه مجسمه گوریل های نر و ماده قبل و بعد از رده بندی طبق روش های (الف) و (ب) و ارائه شده در این مقاله (ج).

| قبل از رده بندی | ماده | نر | بعد از رده بندی (الف) | بعد از رده بندی (ب) | بعد از رده بندی (ج) |
|-----------------|--------|--------|-----------------------|---------------------|---------------------|
| ۰٫۵۴۳۷ | ۰٫۵۴۳۷ | ۰٫۵۴۳۷ | ۰٫۵۴۳۷ | ۰٫۵۴۳۷ | ۰٫۵۴۳۷ |
| ۰٫۵۴۹۹ | ۰٫۵۴۹۹ | ۰٫۵۴۹۹ | ۰٫۵۴۹۹ | ۰٫۵۴۹۹ | ۰٫۵۴۹۹ |
| ۰٫۵۴۶۸ | ۰٫۵۴۶۸ | ۰٫۵۴۶۸ | ۰٫۵۴۶۸ | ۰٫۵۴۶۸ | ۰٫۵۴۶۸ |
| ۰٫۵۵۸۶ | ۰٫۵۵۸۶ | ۰٫۵۵۸۶ | ۰٫۵۵۹۴ | ۰٫۵۵۸۶ | ۰٫۵۵۸۶ |

به منظور مقایسه روش معرفی شده در این مقاله با روش های دیگر رده بندی پیکربندی ها، ما روش خود را با روش های ارائه شده در الف- نیل و گل علی زاده (۲۰۱۶) و ب- مقیم بیگی و نودهی (۲۰۲۲) بر اساس اندازه گیری RMS و داده های مجسمه گوریل (gorm و gorf) مقایسه می کنیم. این نتایج در جدول ۴ نشان داده شده است. بر اساس معیار RMS اگر مقدار RMS در هر رده، قبل و بعد از رده بندی نزدیک باشد، رده بندی به طور موثر انجام شده است. در واقع این مقدار نشان می دهد که میزان همگنی قبل و بعد از رده بندی به چه اندازه به یکدیگر نزدیک بوده است. همان طور که در این جدول مشاهده می شود، مقادیر RMS در مدل ما قبل و بعد از رده بندی بسیار نزدیک به هم هستند که نشان می دهد مدل ارائه شده در این مقاله (ج) از روش های (الف) و (ب) عملکرد بهتری در رده بندی این مجموعه از داده ها دارد. در واقع یکسان بودن مقادیر جدول قبل و بعد از رده بندی با استفاده از روش ارائه شده در این مقاله نشان دهنده این است که رده بندی به طور کامل و با OCP برابر ۱۰۰ درصد انجام گرفته است.

بحث و نتیجه‌گیری

در این مقاله رده‌بندی پیکربندی‌ها با استفاده از یک مدل نیمه پارامتری انجام گرفت. بخش ناپارامتری مدل با به حداقل رساندن معیار توان-واگرایی موضعی، که در آن تابع وزن با استفاده از توزیع آماری فضای شکل به دست آمده بود، معرفی شد. پارامتر هموارسازی نیز با استفاده از روش اعتبارسنجی متقابل حداقل مربعات برآورد شد. با توجه به این‌که برآورد ناپارامتری یک برآورد اریب است برای رده‌بندی پیکربندی‌ها یک مدل MLR معرفی شد. در این مدل، برآورد پارامترهای MLR با استفاده از روش شبکه عصبی انجام گرفت. سپس پیکربندی‌های شبیه‌سازی شده با استفاده از مدل معرفی شده رده‌بندی شدند. همچنین دو مجموعه داده واقعی نیز با استفاده از مدل رده‌بندی شدند که نشان از عملکرد مناسب مدل داشت. بر اساس تحلیل ارائه شده در زیربخش ۴.۲، روش معرفی شده در این مقاله در مقایسه با روش‌های ارائه شده توسط مقیم بیگی و نودهی (۲۰۲۲) و نییل و گل‌علی‌زاده (۲۰۱۶) از عملکرد مناسب‌تری برخوردار است. اما از نقطه نظر محاسباتی باید به این نکته اشاره کرد که روش نیمه پارامتری ارائه شده در این مقاله زمان‌بر است که هدف تحقیق آتی نویسندگان این مقاله رویارویی با این نقیصه است.

تقدیر و تشکر

نویسنده مقاله از رهنمودها و نظرات ارزنده داوران، سردبیر و ویراستار محترم مجله علوم آماری که سبب ارتقای کیفی مقاله شد، کمال تشکر و قدردانی را دارد.

مراجع

مقیم بیگی، م. و گل‌علی‌زاده م. (۱۳۹۸)، مدل‌بندی رگرسیونی شکل از طریق مثلثی کردن، مجله علوم آماری، ۱۳(۱)، ۱۸۵-۱۹۶.

فتوحی، ح. و گل‌علی‌زاده، م. (۱۳۹۱)، بهبود عملکرد تحلیل ژنودزیک اصلی در تحلیل آمار شکل، مجله علوم آماری، ۲(۶)، ۲۱۹-۲۳۶.

Akaike, H. (1973), Information Theory and an Extension of the Maximum Likelihood Principle. In *B. N. Petrov, & F. Csaki (Eds.), Proceedings of the 2nd International Symposium on Information Theory.*

Bookstein, F. L. (1996), Biometrics, Biomathematics and the Morphometric Synthesis, *Bulletin of Mathematical Biology*, **58**, 313-365.

- Bookstein, F. L. (1986), Size and Shape Spaces for Landmark Data in Two Dimensions (With Discussion), *Statist. Sci.*, **1**, 181–242.
- Cressie, N. and Read, T. R. (1984), Multinomial Goodness-of-Fit Tests, *Journal of the Royal Statistical Society: Series B (Methodological)*, **46**, 440–464.
- Cox, D. R. and Snell, E. J. (1989), *Analysis of Binary Data*, Chapman and Hall/CRC, Boca Raton.
- Debavelaere, V. Durrleman, S. and Allasonnière, S. (2020), Learning the Clustering of Longitudinal Shape Data Sets into a Mixture of Independent or Branching Trajectories, *International Journal of Computer Vision*, **128**, 2794–2809.
- Dryden, I. L. (2012), *Shapes package*, R Foundation for Statistical Computing, Vienna, Austria, Contributed package.
- Dryden, I. L. Hirst, J. D. and Melville, J. L. (2007), Statistical Analysis of Unlabeled Point Sets: Comparing Molecules in Chemoinformatics, *Biometrics*, **63**, 237–251.
- Dryden, I. L. and Mardia, K. V. (2016), *Statistical Shape Analysis: With Applications in R*, John Wiley and Sons, Chichester.
- Duin R. P. W. (1976), On the Choice of Smoothing Parameters for Parzen Estimators of Probability Density Functions, *IEEE Transactions on Computers C-25*, 1175–1179.
- Kendall, D. G. (1977), The Diffusions of Shape, *Advances in Applied Probability*, **9**, 428–430.
- Kendall, D. G. (1984), Shape Manifolds, Procrustean Metrics, and Complex Projective Spaces, *Bulletin of the London Mathematical Society*, **16**, 81–121.
- Lin, K. C. and Chen, Y. J. (2008), Assessing Ordinal Logistic Regression Models via Nonparametric Smoothing, *Communications in Statistics-Theory and Methods*, **37**, 917–930.

- Moghimbeygi, M. and Nodehi, A. (2022), Multinomial Principal Component Logistic Regression on Shape Data, *Journal of Classification*, 1–22.
- Nabil, M. and Golalizadeh, M. (2016), On Clustering Shape Data. *Journal of Statistical Computation and Simulation*, **86**, 2995–3008.
- O’Higgins, P. (1989), *A Morphometric Study of Cranial Shape in the Hominoidea*, PhD thesis, University of Leeds .
- O’Higgins, P. and Dryden, I. L. (1993), Sexual Dimorphism in Hominoids: Further Studies of Craniofacial Shape Differences in Pan, Gorilla, Pongo, *Journal of Human Evolution*, **24**, 183–205.
- Okumura, H. and Naito, K. (2006), Nonparametric Kernel Regression for Multinomial Data, *Journal of multivariate analysis*, **97**, 2009–2022.
- Rifada, M. Chamidah, N. and Ratnasari, V. (2021), Estimation of Nonparametric Ordinal Logistic Regression Model Using Local Maximum Likelihood Estimation, *Commun. Math. Biol. Neurosci.*, 2021, Article-ID.
- Schwarz, G. (1978), Estimating the Dimension of a Model, *The annals of statistics*, **6**, 461–464.
- Simó, A. Victoria Ibáñez, M. Epifanio, I. and Gimeno, V. (2020), Generalized Partially Linear Models on Riemannian Manifolds, *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **69**, 641–661.
- Venables, W. N. and Ripley, B. D. (2002), *Modern Applied Statistics with S. 4th ed.* Springer.
- Wasserman, L. (2006), *All of Nonparametric Statistics*, Springer Science & Business Media.
- Yankov, D. and Keogh, E. (2006), Manifold Clustering of Shapes, *In Proc. IEEE Int. Conf. Data Mining (ICDM)*.