

مدل رگرسیون کلاس پنهان با متغیرهای کمکی مؤثر بر پاسخ پنهان و پاسخهای مشاهده شده

آوات فیضی، جمشید جمالی، سید محسن حسینی
گروه آمار زیستی و اپیدمیولوژی، دانشگاه علوم پزشکی اصفهان

در بسیاری از حوزه های پژوهشی مفاهیم نظری وجود دارند که مستقیم قابل مشاهده یا اندازه گیری نیستند که در ادبیات آماری موسوم به متغیرهای پنهان هستند. این مفاهیم براساس متغیرهای قابل مشاهده بعنوان نشانگرهایی از آنها استخراج می گردند. مدلها می تغیر پنهان ابزار توانمندی در مدلسازی رابطه بین متغیرهای مشاهده شده از طریق استخراج متغیرهای پنهان می باشند. در حوزه متغیرهای پنهان، تحلیل کلاس پنهان با فرض وجود رابطه بین مجموعه ای از متغیرهای مشاهده شده گستته اقدام به مدلبندی رابطه بین آنها در قالب استخراج یک متغیر پنهان گستته می نماید. در سالهای اخیر، مدلها کلاس پنهان کاربرد فراوانی در تجزیه و تحلیل ارتباط بین متغیرهای گستته و متغیرهای کمکی یافته اند که در قالب مدلها رگرسیون کلاس پنهان مورد استفاده قرار گرفته اند. در این مقاله فرم گسترش یافته ای از رگرسیون کلاس پنهان معرفی می شود که امکان ارزیابی اثرگذاری متغیرهای کمکی بر متغیر پاسخ پنهان و نیز متغیرهای پاسخ مشاهده شده را همزمان فراهم می کند. مبانی نظری مدل بیان می شود، روشهای برآورد پارامترها و شناسایی پذیری مدل توضیح داده می شود. روش های نظری بسط مدل که موجب گسترش کارکرد عملی آن می شود مورد اشاره قرار می گیرد. کارکرد عملی مدل با ذکر مثالی در حوزه پزشکی به منظور سنجش میزان آگاهی افراد از علائم هشدار دهنده سرطان و عوامل مؤثر بر آن شرح داده می شود.

واژه های کلیدی: متغیر پنهان، رگرسیون کلاس پنهان، شناسایی پذیری، علائم هشدار دهنده سرطان.

۱ مقدمه

بسیاری از مفاهیم و معیارهای موجود در علوم اجتماعی، رفتاری، پزشکی و اپیدمیولوژی بطور واضح و مستقیم قابل ارزیابی نیستند. نگرشها، استعداد، خصوصیات اخلاقی، رفتارهای اجتماعی، ویژگیهای شخصیتی، میزان آگاهی و بعضی از بیماریهای روان پزشگی مانند

استرس ، افسردگی ، اسکیزوفرنی و ... از این مقوله هستند که در ادبیات آماری به متغیرهای پنهان موسوم هستند. رهیافتی که به استخراج و مطالعه ماهیت واقعی مقاهم پنهان کمک می کند مدلهای متغیر پنهان می باشد ، زیرمجموعه ای از این دسته کلی شامل آنالیز کلاس پنهان و رگرسیون کلاس پنهان می باشد. ساده ترین مدل در این زیر مجموعه ، آنالیز کلاس پنهان می باشد. آنالیز کلاس پنهان یک روش آماری برای استخراج زیر جامعه های همگن از یک جامعه براساس الگوی پاسخ آزمودنی ها به داده های چندمتغیره گستته مشاهده شده می باشد. ایده اصلی آنالیز کلاس این است که آزمودنی های مورد مطالعه با هم همبسته هستند به گونه ای که جوامع موردنظر را می توان ترکیبی از چندین زیر جامعه همگن در نظر گرفت. هدف آنالیز کلاس پنهان یافتن این زیر جوامع می باشد که نقش کلاس های پنهان را ایفا می نمایند. در آنالیز کلاس پنهان از میان آزمودنی ها آنهایی که از لحاظ میزان پاسخ دهی به متغیرهای مشاهده شده شباهت ملموسى دارند در یک کلاس قرار می گیرند. آزمودنی های واقع در هر کلاس کاملا مشابه و آزمودنی های واقع در کلاس های مختلف متفاوت از همیگر هستند و همین خاصیت ، امکان انجام مقایسه در بین کلاس ها را از لحاظ الگوی پاسخ به متغیرهای مشاهده شده فراهم می کند. براین اساس تحلیل کلاس پنهان در داده های گستته را می توان متناظر با تحلیل عاملی در نظر گرفت. تحلیل عاملی با ساختار متغیرها مرتبط است یعنی متغیرها بر مبنای ملاک های همچون ضربه همبستگی در گروهها دسته بندی می شوند اما آنالیز کلاس پنهان با ساختار آزمودنی ها مرتبط است.

رگرسیون کلاس پنهان حالت تعیین یافته آنالیز کلاس پنهان می باشد که در آن متغیرهای کمکی ، برای بهینه کردن تشکیل کلاس های پنهان و ارزیابی نحوه اثرگذاری بر قرار گرفتن آزمودنی ها در کلاس های تشکیل شده (نقش متغیر پاسخ را در این مدلها ایفا می کنند) وارد مدل می شوند. براین مبنای این مدل رگرسیونی کاملاً شبیه مدل رگرسیون لجستیک می باشد با این تفاوت که در رگرسیون لجستیک متغیر پاسخ یک متغیر گستته است که به طور مستقیم قابل مشاهده و اندازه گیری است ، در حالیکه در رگرسیون کلاس پنهان متغیر پاسخ یک متغیر گستته (کلاس های پنهان) است که به طور مستقیم قابل مشاهده یا اندازه گیری نیست. در رگرسیون لجستیک امکان مدل بندی رابطه فقط یک متغیر پاسخ با متغیرهای مستقل وجود دارد در حالیکه مدل رگرسیون کلاس پنهان امکان مدل بندی رابطه چندین متغیر پاسخ را با متغیرهای مستقل فراهم می سازد.

آغاز استفاده از آنالیز کلاس پنهان به اوایل دهه ۱۹۵۰ میلادی توسط پژوهشگران حوزه روانشناسی از جمله لازارسفلد و هنری برمنی گردد، که در مطالعه خویش اقدام به ساختن الگوهایی براساس متغیرهای دو حالتی کردند (لازارسفلد ، ۱۹۵۰ و لازارسفلد و هنری ، ۱۹۶۸). در سال ۱۹۷۴ گودمن با معرفی یک روش نسبتاً ساده برای بدست آوردن

برآوردگرهای ماکریم درستنماهی برای پارامترهای مدل کلاس پنهان ، کاربرد این مدلها را در رشته های مختلف تحقیقی عمومی تر کرد(گودمن ، ۱۹۷۴). فورمن در سال ۱۹۸۲ با استفاده از الگوریتم EM روش پیشنهادی گودمن را توسعه داد(فورمن ، ۱۹۸۲). در سال ۱۹۸۸ دایتون و مک ریدی متغیرهای کمکی را وارد آنالیز کلاس پنهان کردند و تأثیر این متغیرها را بر کلاسهای تشکیل شده بررسی نمودند و در واقع رگرسیون کلاس پنهان را بینانگذاری کردند(دایتون و مک ریدی ، ۱۹۸۸). فورمن در سال ۱۹۹۲ چارچوب کلی از مدلها رگرسیون کلاس پنهان را پیشنهاد کرد که انواع مختلفی از این مدلها را پوشش می داد(فورمن ، ۱۹۹۲). بندهن-روچ در سال ۱۹۹۷ مدل رگرسیون کلاس پنهان را برای متغیرهای گسته چند رده ای بسط داد(بندهن-روچ و همکاران ، ۱۹۹۷). هوانگ در سال ۲۰۰۴ فرم گسترش یافته ای از مدلها رگرسیون کلاس پنهان را که یکی از اجزای مدل کلی پیشنهادی توسط فورمن بود را با جزئیات نظری کامل ترا رائه کرد ، در مدل پیشنهادی هوانگ تأثیر متغیرهای کمکی علاوه بر متغیرهای پنهان بر متغیرهای مشاهده شده نیز سنجیده می شد(هوانگ ، ۲۰۰۴). مدل یاد شده بطور همزمان در قالب یک مدل واحد کارکرد مدل رگرسیون کلاس پنهان معمولی و نیز نتایجی را که از اجرای همزمان چندین رگرسیون لجستیک به آن خواهیم رسید در خود دارد. بنابراین چنین مدلی می تواند به لحاظ کاربردی در بسیاری از حوزه های پژوهشی که با ساختارهای نظری پنهان مواجه هستند بویژه در حوزه پزشکی و بهداشت در مواردی که پژوهشگران با سنجش آگاهی و بویژه با تشخیص های پزشکی سروکار دارند ابزار توانمندی در استخراج ساختار پنهان مورد مطالعه (برای مثال آگاهی ، تشخیص) و ارزیابی عوامل اثرگذار بر آن و بر اجزای سازنده آن (نشانگرهایی که ساختار پنهان براساس آنها ساخته می شود) باشد. در مقاله حاضر مبانی نظری روش پیشنهادی هوانگ معرفی و با توجه به اهمیت کارکرد عملی که در حوزه پزشکی و بهداشت می تواند داشته باشد مثالی در این حوزه شرح داده می شود و چگونگی گسترش نظری این مدل که موجب بسط کارکرد عملی آن خواهد شد مورد اشاره قرار خواهد گرفت.

۲ مبانی نظری مدل

در راستای معرفی ساختار نظری مدلها رگرسیون کلاس پنهان ، فرض کنید' ($Y_i = (Y_{i1}, \dots, Y_{iM})$) نشان دهنده یک بردار با M متغیر مشاهده شده برای آزمودنی نام از یک نمونه n عضوی باشد. فرض می شود که Y_i ها با هم وابسته آماری هستند. ایده اصلی

مدل کلاس پنهان این است که آزمودنی های مورد مطالعه با هم همبسته هستند و جوامع مورد مطالعه ترکیبی از J زیرجامعه می باشد که یکی از اهداف رگرسیون کلاس پنهان یافتن این زیر جوامع و برآورد حجم هر زیر جامعه می باشد.

فرض کنید هر متغیر مشاهده شده مورد بررسی K رسته داشته باشد. $(1 < K < S_i)$ نشان دهنده نامین زیر جامعه (کلاس پنهان) باشد و هر متغیر مشاهده شده مورد بررسی K رسته داشته باشد؛ K_m تعداد رسته های متغیر m می باشد و $S_i = 1, \dots, J$ ، $m = 1, \dots, M$ می باشد . $Y_{im} = 1, \dots, K_m$

حداکثر تعداد کلاسهای ممکن برابر است با $K_1 \times K_2 \times \dots \times K_m \times \dots \times K_M$ که در آن K_M تعداد رسته های متغیر m می باشد.
تابع چگالی Y_i در آنالیز کلاس پنهان بصورت زیر تعریف می شود.

$$p[Y_{i1} = y_1, \dots, Y_{iM} = y_M] = \sum_{j=1}^J \eta_i(x'_i \beta) \prod_{m=1}^M \prod_{k=1}^{k_m} p_{mkj}^{y_{mk}}$$

که در آن $\eta_j = p(s_i = j)$ احتمال عضویت هر آزمودنی در کلاس j و p_{mkj} احتمال انتخاب k مین رسته متغیر m به شرط حضور در کلاس j توسط آزمودنی می باشد. $y_{mk} = 1$ است اگر $y_m = k$ و در غیر اینصورت مقدار آن صفر می باشد.

فورمن و هوانگ تعمیمی از رگرسیون کلاس پنهان را ارائه نمودند که متغیرهای کمکی علاوه بر اثربخشی بر کلاسها ، متغیرهای پاسخ مشاهده شده را نیز تحت تأثیر می گذارند. (فورمن ، ۱۹۹۲ و هوانگ ، ۲۰۰۴)

معادله این نوع رگرسیون کلاس پنهان که همزمان کلاس بندی آزمودنی ها و متغیرهای مشاهده شده اثرپذیر از متغیرهای کمکی می باشد ، بصورت زیر می باشد:

$$p[Y_{i1} = y_1, \dots, Y_{iM} = y_M | x_i, z_i] = \sum_{j=1}^J \eta_i(x'_i \beta) \prod_{m=1}^M \prod_{k=1}^{k_m} p_{mkj}^{y_{mk}} (\gamma_m + z'_{im} \alpha_m)$$

که در آن X بردار متغیرهای کمکی مؤثر بر کلاس بندی (متغیر پاسخ پنهان) و Z بردار متغیرهای کمکی مؤثر بر متغیرهای پاسخ مشاهده شده می باشد. همچنین $\eta_i(x'_i \beta)$ و $p_{mkj}^{y_{mk}} (\gamma_m + z'_{im} \alpha_m)$ توابع پیوندی هستند که در چارچوب مدلهای خطی تعمیم یافته تعریف می شود (مک کولاک و نلدر ، ۱۹۸۹). توابع پیوند گوناگونی می توانند به آسانی برای این منظور مورد استفاده قرار گیرند. عمده ترین پیشنهاد استفاده از تابع لجیت تعمیم یافته می

باشد (اگرستی، ۲۰۰۲).

$$\log\left[\frac{\eta_j(x'_i\beta)}{\eta_J(x'_i\beta)}\right] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{iP} \quad \text{for } i = 1, \dots, n; \quad j = 1, \dots, J-1$$

و

$$\log\left[\frac{p_{mkj}(\gamma_{mj} + z'_{im}\alpha_m)}{p_{mKj}(\gamma_{mj} + z'_{im}\alpha_m)}\right] = \gamma_{mkj} + \alpha_{1mk}z_{im1} + \dots + \alpha_{Lmk}z_{imL}$$

برای

$$i = 1, \dots, n; \quad m = 1, \dots, M; \quad k = 1, \dots, (K_m - 1); \quad j = 1, \dots, J$$

برای انجام رگرسیون کلاس پنهان سه پذیره زیربنایی وجود دارد (بندن-روچ و همکاران، ۱۹۹۷ و هوانگ، ۲۰۰۵).

- ۱- احتمال عضویت آزمودنی ها در کلاس پنهان تنها به متغیر کمکی x مربوط است و متغیر z تأثیری بر کلاس بندی آزمودنی ها ندارد. $p(S_i = j|x_i, z_i) = p(s_i = j|x_i)$
- ۲- در داخل کلاسها، متغیرهای مشاهده شده مورد بررسی از متغیرهای کمکی مستقلند.

$$p[Y_{i1} = y_1, \dots, Y_{iM} = y_M | S_i, x_i, z_i] = p[Y_{i1} = y_1, \dots, Y_{iM} = y_M | S_i, z_i]$$

۳- در داخل کلاسها، متغیرهای پاسخ مشاهده شده مستقل از یکدیگرند. (استقلال موضعی)

$$p[Y_{i1} = y_1, \dots, Y_{iM} = y_M | S_i, z_i] = \prod_{m=1}^M p[Y_{im} = y_m | S_i, z_{im}]$$

۱-۲ شناسایی پذیری

برای بررسی یکتا بودن پارامترها، شناسایی پذیری (باید ارزیابی گردد).
با به تعریف، توزیع F_Y را برای پارامتر ϕ شناسایی پذیر گوئیم هرگاه یک X در همسایگی $F_Y(y; \phi_0) = F_Y(y; \phi)$ $y \in U_Y \leftrightarrow \phi = \phi_0$. $\phi \in \Phi$ وجود داشته باشد بطوریکه:

¹Identifiability

$$X \subset \phi$$

که در آن ϕ فضای پارامتر و U_Y تکیه گاه y می باشد.

به زبان ساده ، شناسایی پذیری با این موضوع سروکار دارد که آیا برای هر یک از پارامترهای آزاد مدل یک مقداری گانه از روی داده های مشاهده شده به دست می آید یا خیر؟

در رگرسیون کلاس پنهان و تعمیم های ارائه شده بر آن ، بیشتر بحث شناسایی پذیری بر روی شناسایی پذیری موضعی متصرکز شده است. شرایط کافی برای شناسایی پذیری مدل رگرسیون کلاس پنهان بصورت زیر می باشد.

قضیه: مدل رگرسیون کلاس پنهان شناسایی پذیر است هرگاه:

- ۱) تعداد پارامترهای مورد برآورد از تعداد پارامترهای مدل اشباع کمتر باشد.
- ۲) مقادیر پارامترهای مدل $(\gamma_{mkj}, \alpha_{qmk}, \beta_{pj})$ و متغیرهای کمکی متناهی باشند.
- ۳) ماتریس اطلاع فیشر مدل ، پرتبه ستونی باشد یعنی مقادیر ویژه ماتریس اطلاع فیشر بزرگتر از صفر باشد.
- ۴) ماتریس طرح متغیرهای کمکی و پیشگو پرتبه ستونی باشد.

در عمل برای بررسی شناسایی پذیری مدل از چندین مجموعه مقادیر اولیه برای برآورد پارامترها استفاده می کنیم ، چنانچه مقادیر متفاوت اولیه در تابع درستنمایی ماکریتم ، منجر به برآورد مشابهی برای پارامترها شود ، مدل شناسایی پذیر خواهد بود (هوانگ ، ۲۰۰۴).

۲-۲ برآورد پارامترها

برآورد پارامترها در مدلهای کلاس پنهان معمولاً از طرق روش حداکثر درستنمایی انجام می شود. از آنجاییکه در مدلهای رگرسیون کلاس پنهان با متغیرهای گستره و حجم وسیع داده ها روبرو هستیم برای برآورد پارامترهای معمولا از روشهای پیشرفته ریاضی و الگوریتم های خاص از جمله الگوریتم EM استفاده می شود.

الگوریتم EM تابع درستنمایی داده های ناکامل را در یک فرآیند تکراری بین مرحله ماکریتم سازی تابع درستنمایی داده های کامل (مرحله M) و جانهی داده های ناکامل از طریق مدلی که پارامترهای آن در آخرین تکرار برآورد شده بودند (مرحله E) ، ماکریتم می کند. عضویت در کلاس های پنهان S_i به عنوان یک متغیر غیر قابل مشاهده ، ساختار داده های ناکامل را در مدل رگرسیونی مدنظر شکل می دهد؛ در زیر به اختصار عملکرد این الگوریتم را در برآورد پارامترهای مدل شرح می دهیم.

با فرض قابل مشاهده بودن S_i فرم لگاریتم تابع درستنمایی برای داده های کامل بصورت زیر خواهد بود.

$$\begin{aligned} \log L_c(\phi; Y, S) &= \sum_{i=1}^n \sum_{j=1}^J S_{ij} [\log \eta_j(x_i' \beta)] \\ &+ \sum_{i=1}^n \sum_{j=1}^J \sum_{m=1}^M \sum_{k=1}^{K_m} S_{ij} Y_{imk} [\log p_{mkj}(\gamma_{mj} + z_{im}' \alpha_m)] \end{aligned}$$

که در آن S_{ij} نشان دهنده عضویت نامین آزمودنی در زمین کلاس و Y_{imk} نشان دهنده k نامین سطح متغیر Im که توسط نامین آزمودنی پاسخ داده شده است. همچنین $\phi = (\gamma_{mj}, \alpha_m, \beta)$ پارامترهای مدل می باشد.
تابع Q را به صورت زیر تعریف می کنیم.

$$Q(\phi|\phi') = E[\log L_c(\phi, Y, S|Y=y, \phi', x, z)]$$

این تابع، امید لگاریتم تابع درستنمایی به شرط داده های مشاهده شده y, x, z و برآوردهای بدست آمده پارامترهادر مراحل قبل یعنی (ϕ') می باشد. الگوریتم EM از ϕ^p به ϕ^{p+1} بصورت زیر است.

مرحله E: محاسبه $Q(\phi|\phi^p)$

$$\begin{aligned} Q(\phi|\phi^p) &= \sum_{i=1}^n \sum_{j=1}^J \theta_{ij}(\phi^p) [\log \eta_j(x_i' \beta)] \\ &+ \sum_{i=1}^n \sum_{j=1}^J \sum_{m=1}^M \sum_{k=1}^{K_m} \theta_{ij}(\phi^p) y_{imk} [\log p_{mkj}(\gamma_{mj} + z_{im}' \alpha_m)] \end{aligned}$$

که در آن

$$\theta_{ij}(\phi^p) = E(S_{ij}|Y_i, \phi^p, x_i, z_i)$$

$$= \frac{\eta_j(x_i' \beta) \prod_{m=1}^M \prod_{k=1}^{K_m} p_{mkj}^{y_{mk}} (\gamma_{mj}^p + z_{im}' \alpha_m^p)}{\sum_{l=1}^J \eta_l(x_i' \beta) \prod_{m=1}^M \prod_{k=1}^{K_m} p_{mlk}^{y_{mk}} (\gamma_{ml}^p + z_{im}' \alpha_m^p)}$$

احتمال پسین عضویت کلاسی به ازای ϕ^p می باشد.
مرحله M: یافتن مقادیری از پارامترهای مدل (ϕ) که $Q(\phi|\phi^p)$ را ماکزیمم می نمایند.

۳-۲ تعیین تعداد کلاسهای پنهان (آزمونهای برازنده‌گی مدل)

تعیین تعداد کلاسهای پنهان امری چالش برانگیز بین محققان می‌باشد. از همین رو ملاکهای گوناگونی برای انتخاب مدل مناسب مطرح شده‌اند. رایجترین روش این است که با نظر محقق تعداد کلاسهای تعیین می‌شود سپس شاخصهای نیکویی برازش مانند $BIC = -2\ln L + 3p$ و $AIC = -2\ln L + 2p$ و $AIC^3 = -\ln L + 3p$ تعداد کلاسهای معین محاسبه می‌کنند و مدلی که مقادیر کمتر این شاخصها را داشته باشد مناسبتر می‌باشد.

تازه‌ترین مطالعات نشان می‌دهد که معیار AIC^3 معیار بهتری نسبت به AIC ، BIC برای تعیین تعداد کلاسهای پنهان در مدل رگرسیون کلاس پنهان می‌باشد (اندروز و کوریم، ۲۰۰۴ و دیاس، ۲۰۰۳).

۳ مثال

برای توضیح رگرسیون کلاس پنهان از داده‌های که به منظور بررسی سطح آگاهی نسبت به علائم هشدار دهنده سرطان و عوامل مؤثر بر آن در ساکین شهر تهران و حومه جمع آوری شده اند استفاده می‌کنیم. این مطالعه در سال ۱۳۸۶ و درین ۲۵۰۰ شهروند ۱۸ ساله و بالاتر تهرانی صورت گرفته است. در این پژوهش، آگاهی یا شناخت به عنوان یک مفهوم انتزاعی (پنهان) در نظر گرفته می‌شود که مستقیم قابل استنتاج نبوده و ارزیابی آن بر مبنای پاسخ به پرسش‌های دوگزینه ای در خصوص شناسایی ۹ علامت هشدار دهنده سرطان (به عنوان متغیرهای پاسخ مشاهده شده) بعنوان یک متغیر پنهان گسترش استخراج می‌شود. این متغیر نقش متغیر پاسخ پنهان را در مدل رگرسیون ایفا می‌کند. در این مطالعه علائم تغییر در عادات دفع ادرار و مدفوع، رخمی که بیش از سه هفته بهبود و ترمیم نیاید، خونریزی یا ترشح غیر معمول، سفتی وجود توده ای در پستان و یا سایر ارگانها، اشکال در بلع (اشکال در قورت دادن غذا)، اشکال در هضم و جذب غذا (سوء هاضمه)، تغییرات قابل ملاحظه در خالها و زگیل‌ها، سرفه‌های مکرر و یا خشونت و تغییر در صدا، کاهش وزن ناگهانی عنوان علائم احتمالی سرطان در نظر گرفته شده اند و تأثیر متغیرهای سن، جنسیت، وضعیت تاہل، سطح تحصیلات، سابقه خانوادگی ابتلا به سرطان، مصرف دخانیات و مصرف الکل بر میزان آگاهی کلی و هر یک از علائم هشدار دهنده سنجیده می‌شود.

جدول ۱ نسبت افراد در کلاسهایی که بر مبنای جواب سوالات مطرح شده در مورد

علائم هشدار دهنده سرطان توسط مدل رگرسیون کلاس پنهان ساخته شده است را نشان می دهد. تفسیر و نامگذاری هر کلاس بر مبنای درصدهای شناخت صحیح علائم صورت می گیرد. بر این اساس بررسی ساختار کلاس های تشکیل شده بر حسب سطح آگاهی افراد نشان می دهد کلاس اول که رده افراد آگاه را تشکیل می دهد ، شامل بخش کمی از پرسش شوندگان (۱۸٪) می باشد و کلاسهای دوم و سوم در بردارنده قسمت اعظم افراد (در مجموع ۱۸٪) هستند که دارای اطلاع متوسط (۵۴٪) و ضعیف (۶۹٪) در مورد علائم هشدار دهنده سرطان هستند، بنابراین می توان نتیجه گیری کرد که در مجموع سطح آگاهی و شناخت نسبت به علائم هشدار دهنده سرطان در جامعه مورد مطالعه پایین است.

جدول ۱ - درصد شناخت صحیح علائم هشدار دهنده سرطان و نسبت افراد واقع در هر کلاس

علائم هشدار دهنده	پاسخ	کلاس ۱	کلاس ۲	کلاس ۳
تغییر در عادات دفع ادرار و مدفوع	بله	۸۵/۸۳	۲۸/۰۶	۴/۷۳
	خیر	۱۴/۱۷	۷۱/۹۴	۹۵/۲۷
زخمی که بیش از سه هفته بهبود و ترمیم نیاید	بله	۹۵/۴۸	۴۴/۸۴	۲/۸۲
	خیر	۴/۵۲	۵۵/۱۶	۹۷/۱۸
خونریزی یا ترشح غیرمعمول	بله	۹۸/۲۲	۵۲/۱۸	۱/۹۵
	خیر	۱/۷۸	۴۷/۸۲	۹۸/۰۵
سفتی و وجود توده ای در پستان و یا سایر ارگانها	بله	۹۹/۸۶	۸۱/۷۴	۱۵/۳۶
	خیر	۰/۱۴	۱۸/۲۶	۸۴/۶۴
اشکال در بلع	بله	۹۵/۰۰	۲۸/۳۵	۰/۸۶
	خیر	۵/۰۰	۷۱/۶۵	۹۹/۱۴
اشکال در هضم و جذب غذا (سوء هاضمه)	بله	۹۴/۴۹	۳۴/۷۹	۳/۹۶
	خیر	۵/۵۱	۶۵/۲۱	۹۶/۰۴
تغییرات قابل ملاحظه در خال ها و زگیل ها	بله	۹۸/۲۰	۵۶/۰۹	۲/۹۸
	خیر	۱/۸۰	۴۳/۹۱	۹۷/۰۲
سرفه های مکرر و یا خشونت و تغییر در صدا	بله	۹۳/۰۰	۴۴/۰۸	۳/۲۷
	خیر	۷/۰۰	۵۵/۹۲	۹۶/۷۳
کاهش وزن ناگهانی	بله	۹۴/۶۶	۷۰/۶۹	۱۰/۶۱
	خیر	۵/۳۴	۲۹/۳۱	۸۹/۳۹
حجم هر کلاس		۱۸/۸۲	۵۴/۴۹	۲۶/۶۹

جدول ۲ تأثیر متغیرهای مستقل بر قرار گرفتن افراد در کلاس های تشکیل شده (سطوح مختلف آگاهی) در قالب ضرایب رگرسیونی را نشان می دهد. بررسی ضرایب متغیرها (نسبت بخت) در جدول نشان می دهد که مؤثرترین عامل بر سطح آگاهی میزان تحصیلات می باشد. به عنوان مثال افزایش یک واحدی در سطح تحصیلات ، شناس قرار گرفتن در رده سه نسبت به رده یک را 70% کاهش می دهد. متغیر مستقل مهم دیگر که از نظر آماری معنادار گردیده است جنسیت می باشد. شناس قرار گرفتن مردان در کلاس سوم نسبت به کلاس اول 147% بیشتر از زنان می باشد. از دیگر متغیرهای مستقلی که از نظر آماری معنادار شده اند وضعیت تأهل می باشد . شناس قرار گرفتن افراد متاهل در کلاس یک نسبت به کلاس سوم ، 2 برابر افراد مجرد است. علاوه بر این با آنکه اثر سن معنادار گردیده است اما افزایش سن تأثیر چندانی در قرار گرفتن افراد در رده های مختلف ندارد.

بطور کلی افراد متاهل ، دارای تحصیلات بالاتر و زنان احتمال بیشتری دارند که در رده افراد دارای آگاهی بالا قرار گیرند.

جدول ۲-برآورده از اثاث متغیرهای مستقل بر متغیر پاسخ کلاس پنهان

متغیرها	ضرایب	برای ضرایب	بازه اطمینان (مقدار)	نسبت بخت	بازه اطمینان برای نسبت بخت
کلاس ۲ (کلاس ۱ بعنوان مرجم)					
عرض از مبدأ	$4,72$	$(3,93-5,50)$	$11,84(0,000)$	$11,63$	$(51,13-243,72)$
سن	$-0,03$	$(-0,04-0,02)$	$-4,74(0,000)$	$0,97$	$(0,96-0,98)$
جنس (زن)	$0,31$	$(0,04-0,57)$	$2,74(0,025)$	$1,36$	$(1,04-1,77)$
وضعیت تأهل (مجرد)	$-0,62$	$(-0,94-0,29)$	$-2,72(0,000)$	$0,54$	$(0,39-0,75)$
سطح تحصیلات	$-0,76$	$(-0,94-0,59)$	$-8,56(0,000)$	$0,47$	$(0,39-0,55)$
سابقه خانوادگی	$-0,28$	$(-0,59-0,04)$	$-1,70(0,089)$	$0,76$	$(0,55-1,04)$
ابتلا (نداشتن)	$0,16$	$(0,16-0,49)$	$0,98(0,330)$	$1,18$	$(0,85-1,63)$
استعمال دخانیات (نکردن)	$0,01$	$(0,05-0,960)$	$-0,51(-0,54)$	$1,01$	$(0,60-1,21)$
کلاس ۳ (کلاس ۱ بعنوان مرجم)					
عرض از مبدأ	$4,73$	$(3,83-5,62)$	$10,34(0,000)$	$11,302$	$(46,12-276,97)$
سن	$-0,02$	$(-0,03-0,00)$	$-2,50(0,012)$	$0,98$	$(0,97-1,00)$
جنس (زن)	$0,91$	$(0,09-1,22)$	$5,71(0,000)$	$2,47$	$(1,8-3,37)$
وضعیت تأهل (مجرد)	$-0,69$	$(-1,08-0,29)$	$-2,42(0,001)$	$0,50$	$(0,34-0,75)$
سطح تحصیلات	$-1,21$	$(-1,41-1,00)$	$-11,62(0,000)$	$0,30$	$(0,24-0,37)$
سابقه خانوادگی	$-0,35$	$(-0,81-0,11)$	$-1,49(0,140)$	$0,70$	$(0,44-1,12)$
ابتلا (نداشتن)	$-0,14$	$(-0,52-0,24)$	$-0,72(0,470)$	$0,87$	$(0,60-1,27)$
استعمال دخانیات (نکردن)	$0,49$	$(0,10-1,07)$	$1,64(0,100)$	$1,63$	$(0,91-2,93)$
صرف الكل (نداشتن)					

مدل رگرسیون کلاس پنهان مورد استفاده در این پژوهش امکان ارزیابی نحوه اثرگذاری متغیرهای مستقل بر تک تک علائم هشدار دهنده سرطان را نیز فراهم می سازد. جدول ۳ ضرایب متغیرهای مستقلی که بر آگاهی از هر یک از علائم هشدار دهنده سرطان مؤثرند را نشان می دهد. در این جدول ، فقط متغیرهای مستقل مؤثر بر شناخت هر یک از علائم هشدار دهنده ذکر شده است و دیگر متغیرها در سطح ۵٪ معنی دار نشده اند.

جدول ۳- برآورد اثر مستقیم متغیرهای مستقل بر تک تک علائم هشدار دهنده سرطان

علائم هشدار دهنده سرطان	متغیرهای کمکی	آماره آزمون(p-مقدار)	نسبت بخت
تغییر در عادات دفع ادرار و مدفوع	سن	۲/۰۴(۰/۰۰۲)	۱/۰۱
خوبی که پیش از سه هفته بهبود و ترمیم نیاید	جنسیت	۱/۹۷(۰/۰۴۹)	۱/۱۳
خوبیزی یا ترشح غیرمعمول	تحصیلات	۲/۲۳(۰/۰۲۰)	۱/۰۹
سفتی و وجود نوده ای در پستان و یا سایر ارگانها	...		
اشکال در بلع	تحصیلات	۲/۴۵(۰/۰۱۴)	۱/۱۰
اشکال در هضم و جذب غذا(سوء هاضمه)	جنسیت	۶/۴۱(۰/۰۰۰)	۱/۶۰
تغییرات قابل ملاحظه در خال ها و زکیل ها	وضعیت تا هل	۲/۸۰(۰/۰۰۵)	۱/۲۹
سرمه های مکرر و یا خشونت و تغییر در صدا	تحصیلات	۴/۸۷(۰/۰۰۰)	۱/۲۵
کاهش وزن ناگهانی	...		
کاهش وزن ناگهانی	جنسیت	۱/۹۷(۰/۰۴۹)	۱/۱۳
...	وضعیت تا هل	۲/۰۴(۰/۰۴۱)	۱/۱۶
...	صرف الكل	۲/۲۳(۰/۰۲۰)	۱/۳۰

۴ بحث و نتیجه گیری

در این پژوهش فرم تعییم یافته ای از مدل رگرسیون کلاس پنهان معرفی گردید و بیان شد که با توجه به ماهیت نظری آن که امکان ارزیابی متغیر پنهان ، نحوه اثرگذاری متغیرهای کمکی بر متغیر پنهان و نیز بر هر یک از متغیرهای پاسخ مشاهده شده که متغیرهای پنهان بر مبنای آنها تشکیل می شود را فراهم می کند می تواند در حوزه پزشکی و بهداشت در مقوله هایی چون بیماریهای روانی ، سنجش نگرشها و از همه مهمتر در تشخیص های پزشکی بسیار مفید واقع شود. در مثالی که برای توضیح کارکرد این روش بیان شد میزان آگاهی یک جمعیت عمومی از علائم هشدار دهنده بیماری سرطان و عوامل مؤثر بر سطح کلی آگاهی و شناخت هر یک از علائم هشدار دهنده ارزیابی گردید.

رگرسیون کلاس پنهان یک روش ناپارامتری است و تنها فرض عمدۀ آن شناسایی پذیر بودن مدل می باشد که در این مطالعه با بیان یک قضیه شرایط کافی برای شناسایی پذیر بودن مدل بیان گردید. برای برآورد پارامترها در رگرسیون کلاس پنهان دو روش عمدۀ ماکزیمم درستنمایی (ML) و روش ماکزیمم پسین^۲ (MAP) وجود دارد. برآورد MAP بیزی از ماکزیمم کردن توزیع لگاریتم پسین بدست می آید. بطوریکه توزیع لگاریتم پسین از مجموع تابع لگاریتم درستنمایی و لگاریتم های پیشین برای پارامترها حاصل می شود. گرچه بطور کلی تفاوت بسیار زیادی بین ML و MAP وجود ندارد. مهمترین مزیت روش MAP نسبت به ML این است که از وقوع جوابهای مرزی جلوگیری می کند بعبارت دیگر احتمالات واریانس صفر نخواهد شد. با کمی اطلاع از توزیع پیشین ، برآوردهای پارامترها در داخل فضای پارامتر قرار خواهند داشت. امروزه ، بیشتر نرم افزارها از الگوریتم EM یا تعدیل یافته آن برای برآورد پارامترهای مدل استفاده می کنند. استفاده از الگوریتم EM باعث یافتن جوابهای بهینه برای پارامترها حتی زمانیکه مقادیر اولیه از جواب بهینه دور باشند می گردد. مهمترین ویژگی ارزشمند این روش پایابی برآوردها می باشد.

مدل بیان شده در این مقاله از طریق وارد کردن متغیرهای پیش بین پنهان در جزء پیش بینی کننده مدل در حالت معمولی و نیز در حالت چندسطحی قابل تعمیم می باشند (فیضی و ورمانت ، ۲۰۰۹). با توجه به آنکه در عمل متغیرهای پیش بینی که امکان ارزیابی مستقیم آنها وجود ندارد (از طرق مدلهای عاملی استخراج می گرددن) به وفور با آن مواجه می شویم و از طرفی با توجه به حوزه کاربری وسیع داده های چندسطحی ، بسط مدل معرفی شده در مقاله حاضر در چنین ساختاری موجب گسترش چشمگیر حوزه کاربری آن می شود.

مراجع

- [1] Agresti A. (2002). *Analysis of Categorical Data*, 2nd Ed., New York:John Wiley and Sons.
- [2] Andrews RL, Currim IS. (2003). *A Comparison of segment retention criteria for finite mixture logit models* , journal of Marketing Research, 40, 235-43.

²MAximum Posterior method

- [3] Banden-Roche K, Miglioretti D, Zeger S, Rathouz P. (1997). *Latent variable regression for multiple discrete outcomes*, journal of the American Statistical Association, 92, 1375-86.
- [4] Dayton CM, Macready GB. (1988). *Concomitant-variable latent class models*, journal of the American Statistical Association, 83, 173-78.
- [5] Dias JG. (2004). *Finite Mixture Models: Review, Application, and computer intensive Methods*, The Netherlands: Groningen of University.
- [6] Feizi A, Vermunt JK. (2009). *Mixed effects latent class regression with latent factor predictors. Multivariate behavioral research*,(under review).
- [7] Formann AK. (1982). *Linear logistic latent class analysis*, Biometrikam, 24(2), 171-90.
- [8] Formann AK. (1992). *Linear logistic latent class analysis for polytomous data*, journal of the American Statistical Association, 87, 476-86.
- [9] Goodman LA. (1974). *Exploratory latent structure analysis using both identifiable and unidentifiable models*, Biometrikam, 61, 215-31.
- [10] Huang GH, Bandeen-Roche K. (2004). *building identifiable latent variable model with covariate effects on underlying and measured variables*, Psychometrika, 69, 5-32.
- [11] Huang GH. (2005). *Selection the number of classes under latent class regression: a factor analytic analogue*, Psychometrika, 70, 325-45.
- [12] Lazarsfeld PF. (1950). *The logical and mathematical foundations of latent structure analysis*. In: Stouffer SS and et al, editors. *Measurement and prediction*, Princeton:Princeton University Press, 362-412.
- [13] Lazarsfeld PF, Henry NW. (1968). *Latent Structure Analysis*, Boston:Houghton Mill.
- [14] McCullagh P, Nelder JA. (1989). *Generalized Linear Models*, 2nd Ed., London:Chapman and Hall.