

## مقدمه‌ای بر استنتاج و یادگیری در شبکه‌های بیزی

فهیمه مرادی<sup>۱</sup>، علی کریم‌نژاد<sup>۲</sup>، سودابه شمه‌سوار<sup>۳</sup>

چکیده:

شبکه‌های بیزی ابزار جدیدی در مدل‌بندی پدیده‌ها و سیستم‌های ایستا و پویا هستند و در زمینه‌های مختلفی از جمله تشخیص بیماری‌ها، پیش‌بینی آب و هوا، تصمیم‌گیری و دسته‌بندی کاربرد دارند. یک شبکه بیزی یک مدل گرافی-احتمالی است که ارتباط‌های علت و معلولی بین متغیرهای تصادفی را نشان می‌دهد و از یک گراف بدون دور جهت‌دار و یک مجموعه از احتمال‌های شرطی تشکیل شده است. دو موضوع مهم در مدل‌بندی یک مجموعه داده با شبکه بیزی یادگیری ساختاری و یادگیری پارامتری شبکه است. در این مقاله یک شبکه بیزی با ساختار معلوم را در نظر می‌گیریم و با شبیه‌سازی تلاش می‌کنیم ساختار شبکه را با استفاده از دو الگوریتم متداول PC و K<sub>2</sub> یادگیریم. سپس، به یادگیری پارامترهای شبکه می‌پردازیم و برآوردهای ماکریم درستنمایی، ماکریم احتمال پسین و میانگین پسین پارامترهای مورد علاقه را به دست می‌آوریم. در ادامه، عملکرد برآوردها را با استفاده از معیار واگرایی کوبلک-لایبلر مقایسه می‌کنیم و در نهایت، با استفاده از یک مجموعه داده واقعی، به یادگیری ساختاری و پارامتری شبکه می‌پردازیم تا امکان پیاده‌سازی روش‌های پیشنهادی بر روی داده‌های واقعی را نشان دهیم.

**واژه‌های کلیدی:** توزیع دیریکله، شبکه بیزی، یادگیری پارامتری، یادگیری ساختاری.

### ۱ مقدمه

دارد. همچنین از شبکه بیزی در مدل‌بندی شبکه تنظیمی بیان ژن،

مدل‌بندی ترافیک بزرگراه و تشخیص صدا نیز استفاده شده است.

مدل‌بندی داده‌ها با شبکه‌ی بیزی از دو مرحله تشکیل شده است.

شبکه بیزی به یک گراف بدون دور جهت‌دار اطلاق می‌شود که پارامترهایی در قالب احتمال شرطی این ساختار را از حالت کیفی

به حالت کمی تبدیل می‌کنند. رأس‌های این گراف متغیرهای

تصادفی هستند و یال‌های جهت‌دار آن بیانگر وابستگی بین رأس‌ها

می‌باشند. پارامترهای موجود در شبکه میزان این وابستگی را

مشخص می‌کنند. شبکه بیزی به یادگیری ارتباطات سببی کمک

می‌کند و به خاطر ساختار گرافیکی‌ای که دارد به صورت شهودی

قابل درک است. برای ساختن شبکه بیزی می‌توان از اطلاعات

پیشین و اطلاعات افراد خبره استفاده کرد و با کمک تکنیک‌های

آمار بیزی، داده‌ها و اطلاعات موجود در آن زمینه را با هم ترکیب

زنگامی که تعداد نمونه‌ها کم باشد، روش‌های محدودیت‌گرا با

خطاهای آماری زیادی مواجه هستند و هنگامی که تعداد متغیرها

زیاد باشد روش‌های امتیازگرا (به علت کند شدن سرعت انجام

روش‌های گشت) به مشکل برمنی خورند. در این مقاله بر مدل‌بندی

نظریه احتمال، علوم کامپیوتر و آمار می‌باشد و در بسیاری از زمینه‌ها

با شبکه بیزی تمرکز می‌کنیم و پس از یادگیری ساختار شبکه بیزی

از جمله تشخیص، پیش‌بینی، دسته‌بندی و تصمیم‌گیری کاربرد

مورد مطالعه به یادگیری پارامتری در آن شبکه می‌پردازیم.

<sup>۱</sup>دانشجوی کارشناسی ارشد آمار ریاضی، گروه آمار دانشگاه تهران

<sup>۲</sup>دانشجوی دکتری آمار، گروه آمار دانشگاه تهران

<sup>۳</sup>استادیار گروه آمار دانشگاه تهران.

شبکه می‌پردازیم. در نهایت در بخش هفتم، از مطالب بیان شده بیز<sup>۴</sup> (۱۷۲۰) ارائه شده و به همین دلیل آن را شبکه بیزی می‌نامند. نتیجه‌گیری می‌کنیم.

## ۲ تعاریف و مفاهیم مورد نیاز

در این بخش به بیان تعاریف و مفاهیم مورد نیاز بخش‌های آتی مقاله می‌پردازیم. یک گراف بدون دور جهت‌دار را دگ<sup>۱۳</sup> می‌نامیم و شبکه بیزی دگی است که رأس<sup>۱۴</sup>‌های آن متغیرهای تصادفی هستند. مجموعه‌ی همه‌ی رأس‌هایی که از رأس مورد نظر یک یال<sup>۱۵</sup> جهت‌دار به آن‌ها وجود داشته باشد را مجموعه اولاد آن رأس و مجموعه‌ی همه‌ی رأس‌هایی که از آن‌ها یک یال جهت‌دار به رأس مورد نظر وجود داشته باشد را مجموعه اجداد آن رأس می‌نامند. شبکه‌ای که هیچ یالی در آن وجود نداشته باشد را شبکه تهی و شبکه‌ای که در آن بین هر رأس و همه‌ی رأس‌های دیگر یال وجود داشته باشد را شبکه کامل می‌نامند. یک شبکه بیزی از دو بخش  $B_s$  (ساختار بیزی<sup>۱۶</sup>) و  $B_p$  (احتمال بیزی<sup>۱۷</sup>) تشکیل شده و معمولاً آن را با  $(B_s, B_p)$  نشان می‌دهند.  $B_s$  و  $B_p$  به ترتیب به ساختار و پارامترهای شبکه بیزی اشاره می‌کنند. لازم به ذکر است که معمولاً عبارت‌های ساختار شبکه بیزی  $(B_s)$  و دگ<sup>(G)</sup> را به جای یکدیگر به کار می‌برند. یکی از ویژگی‌های مهم شبکه بیزی این است که هر رأس به شرط اجدادش از مجموعه رأس‌های غیراولاد مستقل است. از طرف دیگر ساختار شبکه‌های بیزی دارای دور جهت‌دار نیست و می‌توان  $X_i$  را طوری مرتب کرد که اجداد  $X_i$  در مجموعه‌ی  $\{X_1, \dots, X_{i-1}\}$  و اولاد آن در مجموعه‌ی  $\{X_{i+1}, \dots, X_n\}$  قرار

ایده اصلی شبکه‌های بیزی بر پایه قانون بیز است که توسط توماس بیز<sup>۵</sup> طبق قانون بیز، محققین شبکه‌های از قبل ساخته شده برای متغیرها (که شبکه‌های پیشین<sup>۶</sup> نامیده می‌شوند) را در نظر می‌گیرند و آن‌ها را با مجموعه داده‌ها ترکیب می‌کنند تا شبکه پسین، که شبکه محتمل تری نسبت به شبکه‌های پیشین است، ساخته شود. اصطلاح شبکه بیزی اولین بار توسط پرل<sup>۷</sup> (۱۹۸۸) مطرح شد. امروزه با گذشت ۲۷ سال از پیدایش شبکه بیزی، این شبکه‌ها به عنوان یک ابزار قوی در علوم مختلف از جمله بیوانفورماتیک برای مدل‌بندی شبکه بیان ژن به کار می‌روند. در بین مطالعات انجام شده در این زمینه می‌توان به آنگ<sup>۸</sup> و همکاران (۲۰۰۲) و فریدمن<sup>۹</sup> (۲۰۰۴) اشاره کرد. آنگ و همکاران (۲۰۰۲) یک شبکه بیزی پویای زمان پیوسته برای مدل‌بندی داده‌های بیان ژن ارائه دادند. فریدمن (۲۰۰۴) داده‌های تجربی میزان بیان ژن را به صورت شبکه بیزی ایستاده مدل‌بندی کرده است که در آن شبکه امکان پیش‌بینی دقیق وجود ندارد. برای مطالعه بیشتر در شبکه‌های بیزی می‌توان به جنسن<sup>۱۰</sup> و نیلسن<sup>۱۱</sup> (۲۰۰۷) و کوسکی<sup>۱۲</sup> (۲۰۰۹) مراجعه کرد.

ساختار مطالب بیان شده در این مقاله به شرح زیر است: در بخش دوم تعاریف و مفاهیم مورد نیاز را بیان می‌کنیم. در بخش سوم روش‌های مختلف یادگیری ساختاری شبکه بیزی را معرفی می‌کنیم. در بخش چهارم به یادگیری پارامتری شبکه بیزی می‌پردازیم. سپس، در بخش پنجم یک مجموعه داده‌ی شبیه‌سازی شده را با شبکه بیزی مدل‌بندی می‌کنیم و در بخش ششم با استفاده از یک مجموعه داده واقعی، به یادگیری ساختار و پارامترهای

<sup>۵</sup>Prior Networks

<sup>۶</sup>Pearl

<sup>۷</sup>Ong

<sup>۸</sup>Friedman

<sup>۹</sup>Jensen

<sup>۱۰</sup>Nielsen

<sup>۱۱</sup>Koski

<sup>۱۲</sup>Noble

<sup>۱۳</sup>DAG (Directed Acyclic Graph)

<sup>۱۴</sup>Node

<sup>۱۵</sup>Edge

<sup>۱۶</sup>Bayesian Structure

<sup>۱۷</sup>Bayesian Probability

- رتبه‌بندی) انجام می‌شود که در ادامه به اختصار به معرفی هر یک می‌پردازیم.

### ۱.۳ روش‌های محدودیت‌گرا

در این روش‌ها ساختار شبکه بیزی با مشخص کردن رابطه‌ی استقلال شرطی بین گره‌ها به دست می‌آید و برای مشخص کردن واپسگی یا عدم واپسگی بین متغیرها از آزمون‌های استقلال استفاده می‌کنند. تعیین ساختار شبکه بیزی با روش‌های یادگیری مبتنی بر قید یک فرآیند دو مرحله‌ای است. در مرحله اول، ساختار شبکه تعیین می‌شود به طوری که یال‌های ساختار بدون جهت هستند و اصطلاحاً به این ساختار کالبد دگ گفته می‌شود. در این مرحله برای تعیین یال‌ها از آزمون‌های استقلال شرطی مانند آزمون استقلال  $\chi^2$ , آزمون Z فیشر و آزمون استقلال نسبت درستنمایی استفاده می‌شود. سپس، در مرحله دوم یال‌های ساختاری که در مرحله اول به دست آمده جهت‌دار می‌شوند. در الگوریتم‌های مختلف بر اساس قیود متفاوتی یال‌ها جهت‌دار می‌شوند (جنسن و نیلسن، ۲۰۰۷).

الگوریتم‌های یادگیری ساختاری زیادی مانند الگوریتم استقلال شرطی، الگوریتم اجداد و اولاد<sup>۲۱</sup> (PC) (جنسن و نیلسن، ۲۰۰۷) و الگوریتم دوگانه اجداد و اولاد (ابراهیمی، ۱۳۹۰) وجود دارد. در اینجا تنها به معرفی مهم‌ترین الگوریتم یادگیری ساختاری، یعنی الگوریتم PC می‌پردازیم. برای درک بهتر این الگوریتم گزاره زیر را بیان می‌کنیم.

**گزاره ۱.۳.** در یک شبکه بیزی متغیرهای تصادفی X و Y به شرط مجموعه Z از یکدیگر مستقلند ( $X \perp Y | Z$ ) اگر و تنها اگر X و Y توسط مجموعه Z از یکدیگر جدا شده باشند.  
 $(d - sep_G(X, Y | Z))$

بر طبق این گزاره، اگر X و Y به شرط مجموعه Z از یکدیگر مستقل باشند، آنگاه به ازای هر  $W \in Z$  جهت یال‌ها به صورت

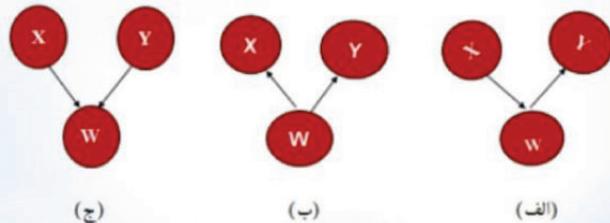
بگیرند، بنابراین طبق قانون احتمال کل، تابع احتمال توأم را می‌توان به صورت زیر نوشت:

$$\begin{aligned} P(X_1, \dots, X_n) &= \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1}) \\ &= \prod_{i=1}^n P(X_i | X_i), \end{aligned}$$

دو شبکه بیزی دارای توزیع احتمال توأم یکسان را معادل گویند. دو متغیر تصادفی X و Y را به شرط مجموعه‌ی Z مستقل می‌نامند و این مطلب را با نماد  $X \perp Y | Z$  نشان می‌دهند هرگاه برای هر Z داده شده، مستقل باشند و در این صورت داریم:

$$P(X, Y | Z) = P(X | Z) P(Y | Z).$$

مسیر بین دو رأس X و Y توسط مجموعه‌ی Z بلوک شده خوانده می‌شود. هرگاه به ازای هر  $W \in Z$  جهت یال‌ها به صورت ساختار (الف) یا (ب) شکل ۱ باشد و یا اینکه برای هر  $W \notin Z$  جهت یال‌ها به صورت ساختار (ج) شکل ۱ باشد.



شکل ۱. بلوک شدن مسیر

مجموعه‌ی Z رأس Y را از رأس X در گراف G جدا می‌کند اگر و تنها اگر هر مسیری که از Y به X وجود دارد توسط مجموعه‌ی Z بلوک شود. این مطلب با نماد  $1^{18} d - sep_G(X, Y | Z)$  نشان داده می‌شود.

## ۳ یادگیری ساختاری

هدف از یادگیری ساختاری یافتن بهترین ساختار برای شبکه بیزی است که با داده‌ها یا اطلاعات موجود مطابقت داشته و از لحظه پیچیدگی بهینه باشد (دارای کمترین پیچیدگی باشد). این یادگیری

<sup>۱۸</sup>“d” for “directed graph”

<sup>۱۹</sup>Constraint-Based

<sup>۲۰</sup>Scoring-Based

<sup>۲۱</sup>Parents and Children

یک دور جهت‌دار شود به حالت مخالف جهت‌دهی خواهد شد.

نهایتاً ممکن است تعدادی از یال‌ها بدون جهت باقی بمانند که در این صورت با یک گراف بدون دور جزئی جهت‌دار روبرو هستیم که متناظر با آن یک خانواده از گراف‌های بدون دور جهت‌دار وجود دارد و باید یکی از ساختارها را به تصادف انتخاب کیم. الگوریتم اجداد و اولاد برای ساختن یک دگ در مرحله زیر را طی می‌کند:

- ساختن کالبد دگ:** الگوریتم ابتدا یک گراف کامل روی تمام متغیرها در نظر می‌گیرد. سپس برای هر دو متغیر، مانند  $X$  و  $Y$ ، مشخص می‌کند که دارای وابستگی هستند یا خیر. این کار با استفاده از آزمون استقلال<sup>۲</sup>  $\chi^2$  به این صورت انجام می‌شود که ابتدا مجموعه  $Z$  را تهی در نظر گرفته و درست بودن رابطه  $Y \perp\!\!\!\perp X$  را بررسی می‌کند. اگر با  $p$ -مقدار در نظر گرفته شده برای آزمون<sup>۲</sup>  $\chi^2$ ، فرض صفر رد نشد یعنی دو رأس از هم مستقل بودند یال بین آن‌ها حذف می‌شود، در غیر این صورت یکی از همسایه‌های  $X$  را به مجموعه  $Z$  اضافه کرده و نتیجه را بررسی می‌کند. به همین ترتیب برای همهی زیرمجموعه‌های تک عضوی از مجموعه‌ی رأس‌هایی که با  $X$  همسایه هستند، این کار را انجام می‌دهد. در هر مرحله که استقلال دو رأس تأیید شد الگوریتم متوقف و یال بین  $X$  و  $Y$  حذف می‌شود، اما اگر یال باقی بماند الگوریتم از آزمون‌های استقلال مرتبه دو استفاده می‌کند، یعنی با انتخاب  $Z$  از زیرمجموعه‌های دو عضوی از همسایه‌های  $X$ ، شرط  $Y \perp\!\!\!\perp X | Z$  را بررسی می‌کند. الگوریتم این کار را ادامه می‌دهد تا جایی که تعداد عناصر مجموعه‌ی  $Z$  برابر تعداد همسایه‌های  $X$  شود (در نظر داشته باشید که  $Z$  زیرمجموعه‌ای از همسایه‌های  $X$  است). اگر در مرحله‌ای رابطه  $Y \perp\!\!\!\perp X | Z$  تأیید شود یال بین  $X$  و  $Y$  حذف خواهد شد در غیر این صورت یال مربوطه باقی می‌ماند. از آنجایی که در الگوریتم ترتیب در نظر گرفتن  $X$  و  $Y$  مهم است، فقط زمانی یال بین  $X$  و  $Y$  باقی می‌ماند که نه در بین همسایه‌های  $X$  و نه در بین همسایه‌های  $Y$ ، مجموعه  $Z$  که  $X$  و  $Y$  را از هم جدا کند

یکی از دو ساختار (الف) یا (ب) شکل ۱ است.

در قسمت (الف)تابع احتمال تؤمن متغیرهای  $W$ ,  $X$  و  $Y$  عبارت از:

$$\begin{aligned} P(W, X, Y) &= P(Y|W)P(W|X)P(X) \\ &= P(Y|W)P(W, X), \end{aligned}$$

است و در قسمت (ب) تابع احتمال تؤمن متغیرهای  $W$ ,  $X$  و  $Y$  عبارت از:

$$\begin{aligned} P(W, X, Y) &= P(Y|W)P(X|W)P(W) \\ &= P(Y|W)P(W, X), \end{aligned}$$

است. از آنجایی که تابع احتمال تؤمن متغیرها در هر دو ساختار یکسان است این دو ساختار با یکدیگر معادلند. این دو ساختارهای  $V$ -شکل نامیده می‌شوند. بر طبق این گزاره، اگر  $X$  و  $Y$  به شرط مجموعه  $Z$  از یکدیگر مستقل نباشند آنگاه به ازای هر  $W \in Z$  جهت یال‌ها به صورت ساختار (ج) شکل ۱ است و تابع احتمال تؤمن متغیرها عبارت از:

$$P(W, X, Y) = P(W|X, Y)P(Y|X)P(X).$$

است.

### ۱.۱.۳ الگوریتم یادگیری ساختاری اجداد و اولاد (PC)

این الگوریتم با آزمون استقلال<sup>۲</sup>  $\chi^2$ ، استقلال یا وابستگی متغیرها را بررسی می‌کند و برای تشخیص استقلال یا وابستگی دو رأس، به دنبال پیدا کردن مجموعه جدایت‌نده در بین تمام همسایه‌های دو رأس است و با استفاده از مفهوم مجموعه جدایت‌نده و دو اصل زیر به جهت یابی یال‌ها اقدام می‌کند (جنسن و نیلسن، ۲۰۰۷):

- الگوریتم از به وجود آمدن ساختار  $V$ -شکل در گراف جلوگیری می‌کند، یعنی ساختار  $Z \leftrightarrow Y \rightarrow X$  را به شکل  $Z \leftarrow Y \leftarrow X$  جهت‌دهی می‌کند.

- الگوریتم از به وجود آمدن دور در شبکه بیزی جلوگیری می‌کند. از آن جایی که بنابر تعریف، یک شبکه بیزی گرافی بدون دور جهت‌دار است در جهت‌دهی به یال‌ها این نکته همواره مدنظر است و هر زمان که جهت یالی باعث تشکیل

روش گشت مشخص می‌شود که تمام ساختارهای دگ ممکن را بررسی کند. در مرحله دوم، یک متر مناسب در نظر گرفته می‌شود که میزان تطابق هر ساختار را با داده‌ها یا مجموعه اطلاعات ارزیابی کند. این دو مرحله را تا جایی ادامه می‌دهند که هیچ ساختار ممکنی دارای تطابق بیشتری نباشد (ابراهیمی، ۱۳۹۰).

### ۱.۲.۳ انواع مترها

مترها به دو دسته زیر تقسیم‌بندی می‌شوند. مترهای وابسته به توزیع پیشین (متر بیزی): این مترها با در نظر گرفتن توزیع داده‌ها  $P(D|G)$  و توزیع پیشین شبکه  $P(G)$ ، توزیع احتمال پسین شبکه  $P(G|D)$  را محاسبه می‌کنند و شبکه‌ای که مقدار احتمال پسین آن ماقسیم باشد را به عنوان بهترین شبکه انتخاب می‌کنند. هدف پیدا کردن شبکه  $G$  است به طوری که  $\frac{P(G,D)}{P(D)}$  را ماقسیم کند. چون  $P(D)$  در همه‌ی شبکه‌ها یکسان است کافی است شبکه‌ای را پیدا کنیم که  $P(G,D)$  را ماقسیم کند. چون کار کردن با لگاریتم آسان‌تر است معمولاً در مترها به جای  $P(G,D)$ ،  $\log(P(G,D))$  را محاسبه می‌کنیم. از جمله مترهای بیزی می‌توان به متر بیزی دیریکله (BD) و حالت‌های خاص آن، مترهای  $K_2$ ،  $BD_{eu}$  و  $BD_e$  اشاره کرد (هکرمن<sup>۲۲</sup> و همکاران، ۱۹۹۵؛ هکرمن، ۱۹۹۶).

مترهای وابسته به مقاهیم نظریه اطلاعات (متر غیربیزی): این مترها میزان تطابق دگ با اطلاعات مجموعه داده‌ها را اندازه‌گیری می‌کنند و ساختاری که تطابق بیشتری با مجموعه داده‌ها داشته باشد به عنوان بهترین ساختار انتخاب می‌کنند. از جمله مترهای غیربیزی می‌توان به متر کوتاهترین توصیف<sup>۲۳</sup> (MDL) که شامل مترهای معیار اطلاع بیزی<sup>۲۴</sup> ( $BIC$ )، معیار اطلاع آکائیک<sup>۲۵</sup> (AIC) و لگاریتم درستنمایی<sup>۲۶</sup> (LL) است، اشاره کرد (ابراهیمی، ۱۳۹۰). در این روش‌ها با استفاده از یک متر، میزان تطابق شبکه‌ها را با اطلاعات موجود محاسبه می‌کنند و به دنبال شبکه‌ای هستند که توجه کنید که یک متر باید حداقل دارای این دو ویژگی باشد: تعادلی بین دقت ساختار و پیچیدگی آن برقرار کند و از نظر محاسباتی قابل حل باشد. یکی از مترهایی که دو ویژگی بیان

وجود نداشته باشد. به این ترتیب تعدادی از یال‌ها حذف خواهد شد و کالبد دگ مشخص می‌شود (جنسن و نیلسن، ۲۰۰۷).

۲. جهت‌دار کردن یال‌های ساختاری که در مرحله‌ی اول مشخص شده: در این مرحله پس از مشخص شدن ساختار گراف، الگوریتم سه تایی  $X$  و  $Y$  را پیدا می‌کند با این ویژگی که در ساختاری که در مرحله‌ی قبل به دست آمده، دو رأس  $X$  و  $W$  همسایه باشند و دو رأس  $Y$  و  $W$  نیز همسایه باشند اما دو رأس  $X$  و  $Y$  همسایه نباشند. در این صورت اگر رأس  $W$  متعلق به مجموعه‌ی جداگانه  $separator_{X,Y}$  باشد جهت یال  $XW$  از  $X$  به  $W$  و جهت یال  $YW$  از  $Y$  به  $W$  تعیین می‌شود (دلیل این امر با توجه به این که  $separator_{X,Y}$  همه‌ی مسیرهای بین  $X$  و  $Y$  را بلوك می‌کند واضح است). زمانی که تمام ساختارهایی که به این شکل هستند جهت‌یابی شدند، الگوریتم یال‌های باقی‌مانده در ساختار دگ را با استفاده از دو اصل مهم ذکر شده جهت‌دهی می‌کند (جنسن و نیلسن، ۲۰۰۷).

### ۲.۱.۳ محدودیت‌های روش‌های محدودیت‌گرا

اگر تعداد نمونه کم باشد یا داده گم شده داشته باشیم، روش‌های محدودیت‌گرا با خطاهای آماری مواجه‌اند. این روش‌ها، در جهت‌دهی به بعضی از یال‌ها ناتوان هستند و روش‌های ناپایداری هستند که یک اشتباه کوچک اولیه می‌تواند کار را به جایی سوق دهد که نتیجه با گراف اصلی بسیار متفاوت باشد.

### ۲.۳ روش‌های امتیاز‌گرا

در این روش‌ها با استفاده از یک متر، میزان تطابق شبکه‌ها را با اطلاعات موجود محاسبه می‌کنند و به دنبال شبکه‌ای هستند که بیشترین تطابق را با داده‌ها داشته باشد. ساختن دگ با استفاده از روش‌های امتیاز‌گرا شامل دو مرحله می‌شود. در مرحله اول، یک

<sup>۲۲</sup>Heckerman

<sup>۲۳</sup>Minimum Description length

<sup>۲۴</sup>Bayesian Information Criterion

<sup>۲۵</sup>Akaike Information Criterion

<sup>۲۶</sup>Log-Likelihood

مترا BIC مثالی از یک مترا تجزیه‌پذیر برای مجموعه داده‌های کامل است. با استفاده از متراهای تجزیه‌پذیر می‌توان افزایش یا کاهش رتبه ساختاری که با تغییر دادن یک یال به دست آمده را نسبت به ساختار اولیه اندازه‌گیری کرد. به عنوان مثال، اگر یک یال از رأس  $X_i$  به رأس  $X_j$  وارد کنیم آنگاه تنها رتبه‌ی رأس  $X_j$  تغییر خواهد کرد یعنی برای اندازه‌گیری تغییر رتبه ساختار جدید نسبت به ساختار اولیه کافی است اختلاف متراهای رأس  $X_j$  را از رابطه زیر محاسبه کنیم:

$$\begin{aligned}\Delta(X_i \rightarrow X_j) = & score(X_j, Pa(X_j) \cup \{X_i\}, D) \\ & - score(X_j, Pa(X_j), D)\end{aligned}$$

اگر  $0 > X_j \rightarrow X_i$  باشد آنگاه ساختار جدید رتبه بیشتری نسبت به ساختار اولیه دارد و در نتیجه بهتر به داده‌ها برآش داده می‌شود. به طور کلی، همه‌ی روش‌های گشت دو مرحله دارند: مرحله شروع: در این مرحله یک گراف اولیه (گراف تهی، شبکه پیشین و ...) در نظر گرفته می‌شود.

مرحله جستجو: بر اساس مترا انتخاب شده، میزان تغییر تطابق ساختارهایی که به عنوان مثال با اضافه شدن، حذف شدن یا بر عکس شدن یک یال به دست می‌آیند، محاسبه می‌شود. این فرایند تا جایی ادامه پیدا می‌کند که هیچ ساختار ممکنی دارای تطابق بیشتری نباشد (جنسن و نیلسن، ۲۰۰۷).

### ۳.۲.۳ الگوریتم گشت: $K_2$

یکی از روش‌های پایه‌ای و پرکاربرد گشت است. ورودی‌های این الگوریتم متغیرهای رتبه‌بندی شده و ماکسیمم تعداد اجدادی است که می‌خواهیم رأس‌ها داشته باشند. در این الگوریتم ابتدا یک ساختار تهی برای شبکه در نظر می‌گیریم سپس گره‌ها را طوری که گره ولد بیشترین امتیاز را داشته باشد و اجداد  $i$  در مجموعه‌ی  $\{X_1, \dots, X_{i-1}\}$  قرار بگیرند، وارد شبکه می‌کنیم. رتبه بندی می‌تواند بر اساس متراهای مختلف مانند  $K_2$ ,  $BDeu$  و  $BIC$  انجام شود. این الگوریتم به دنبال پیدا کردن مجموعه اجداد هر رأس به صورتی است که در ترتیب ایجاد شده برای هر رأس  $X_i$ ، رأسی که بیشتر از بقیه رتبه شبکه را افزایش می‌دهد از مجموعه  $\{X_1, \dots, X_{i-1}\}$  به مجموعه اجداد  $i$  اضافه می‌کنیم. این کار را

شده را دارد مترا BIC است که شامل عبارتی برای اندازه‌گیری میزان برآزندگی مدل به داده‌ها و عبارت دیگری برای اندازه‌گیری پیچیدگی مدل است.

### ۲.۲.۳ روش‌های گشت

به منظور انجام یادگیری ساختاری با روش‌های امتیازگرا، باید پس از در نظر گرفتن یکتابع رتبه‌بندی (مترا)، ساختار شبکه بیزی دارای بیشترین امتیاز (رتبه) را بین مجموعه‌ی همه‌ی ساختارهای شبکه‌ای ممکن جستجو کنیم. یعنی، مسئله‌ی یادگیری ساختاری به یک مسئله‌ی گشت محدود می‌شود. چالش این مسئله آن است که تعداد کل ساختارهای ممکن با تعداد رأس‌ها رابطه‌ی نمایی دارد. اگر  $n$  تعداد رأس‌های یک شبکه بیزی باشد آنگاه تعداد کل ساختارهای ممکن از رابطه‌ی زیر به دست می‌آید:

$$f(n) = \sum_{i=1}^n \frac{n!}{i!(n-i)!} 2^{i(n-i)} f(n-1),$$

که در آن  $f(0) = 1$  است. برای  $n$ ‌های بزرگ، در نظر گرفتن همه‌ی ساختارها غیرممکن است. بنابراین، محققین روش‌های گشت اکتشافی را در نظر می‌گیرند و مکررا با ایجاد تغییرات کوچک روی ساختار کنونی، مناسب‌ترین ساختار را در فضای ساختارهای ممکن جستجو می‌کنند. به دگ‌هایی که با یک تغییر در دگ کنونی به وجود می‌آیند همسایه‌های آن دگ می‌گویند (جنسن و نیلسن، ۲۰۰۷).

عملگرها یی که یک تغییر در دگ کنونی ایجاد می‌کنند به شرح زیر است:

۱. عملگر افزاینده‌ی یال: یک یال بین دو رأس ایجاد می‌کند.
۲. عملگر حذف‌کننده‌ی یال: یال بین دو رأس را حذف می‌کند.
۳. عملگر معکوس‌کننده‌ی یال: جهت یال بین دو رأس را بر عکس می‌کند.

یک خصوصیت مهم این عملگرها این است که آن‌ها تنها یک تغییر موضعی در ساختار موجود می‌دهند. به عنوان مثال، اگر یک یال از رأس  $X_i$  به رأس  $X_j$  وارد شود یا یال بین آن‌ها حذف شود آنگاه تنها خانواده (مجموعه اجداد) رأس  $X_j$  تغییر می‌کند و اگر یال بین آن‌ها معکوس شود آنگاه خانواده‌های هر دو رأس  $X_i$  و  $X_j$  تغییر می‌کند. از این خصوصیت در متراهای تجزیه‌پذیر استفاده می‌شود.

## ۱.۴ روش‌های یادگیری پارامتری

اگر مقادیر همهٔ متغیرها از پیش معلوم باشد (داده گمشده نداشته باشیم) انجام چنین استنتاجی ساده است ولی عموماً فقط بخشی از متغیرها مشاهده می‌شود. بنابراین، فرآیند یادگیری پارامتری در شبکه‌های بیزی با مشخص بودن یا نبودن ساختار شبکه و همچنین قابل مشاهده بودن یا نبودن متغیرها روندهای متفاوتی را طی می‌کند. این روندها در جدول ۱ نشان داده شده‌اند (مورفی، ۲۰۰۱).

جدول ۱: روش‌های یادگیری پارامتری

روش	مشاهده‌پذیری	ساختار
ماکسیمم درستنمایی و بیزی	کامل	مشخص
الگوریتم EM	ناقص	مشخص
گشت در فضای ساختارهای ممکن + ماکسیمم درستنمایی و بیزی	کامل	نامشخص
گشت در فضای ساختارهای ممکن + الگوریتم EM	ناقص	نامشخص

- اگر ساختار شبکه معلوم باشد و تمامی متغیرها قابل مشاهده باشند می‌توان برآورد ML، MAP و یا PM احتمال‌های شرطی هر رأس به شرط اجدادش را از روی داده‌های آموزشی و دانش پیشین در مورد پارامترها به دست آورد.

- در صورتی که ساختار شبکه از قبل معلوم باشد ولی فقط برخی از مقادیر متغیرها قابل مشاهده باشند، یادگیری مشکل‌تر خواهد بود. در این صورت از الگوریتم EM برای برآورد پارامترها استفاده می‌کنیم.

اگر ساختار شبکه معلوم نباشد یادگیری مشکل بوده و بسته به اینکه تمامی متغیرها و یا بخشی از آن‌ها قابل مشاهده باشند به صورت زیر عمل می‌کنیم.

- اگر ساختار شبکه معلوم نباشد اما تمامی متغیرها قابل مشاهده باشند، ابتدا با روش‌های گشت (مانند روش گشت

تا جایی ادامه می‌دهیم که تعداد اعضای مجموعه اجداد  $X_i$  از عدد از پیش تعیین شده بیشتر نشود (جنسن و نیلسن، ۲۰۰۷).

## ۴ یادگیری پارامتری

یادگیری پارامتری به معنی برآورد پارامتر است. هدف از یادگیری پارامتری شبکه بیزی برآورد احتمال شرطی هر رأس شبکه به شرط اجدادش می‌باشد. پس از ساخت یک شبکه بیزی، لازم است که یک سری از مقادیر احتمال از مدل طراحی شده استخراج شود که به این فرآیند، استنتاج می‌گویند. انواع استنتاج عبارتند از استنتاج دقیق و استنتاج تقریبی. در استنتاج دقیق مجموعه داده کامل است و برآورد میزان احتمال شرطی هر رأس شبکه از روش‌های متداول برآوردهای مانند روش ماکسیمم درستنمایی و در برخی از موارد روش بیزی دقیقاً محاسبه می‌شود. در روش ماکسیمم درستنمایی برآورد پارامترها با ماکسیمم کردن چگالی احتمال توأم متغیرها یعنی بر اساس مجموعه داده به دست می‌آید در حالی که در روش بیزی برآورد پارامترها با ماکسیمم کردن چگالی احتمال توأم و چگالی احتمال پیشین پارامترها یعنی با ترکیب کردن اطلاع نتیجه شده از داده‌ها با دانش پیشین در مورد پارامترها به دست می‌آید. معیارهای متداول برای برآورد نقطه‌ای پارامترها در روش ML<sup>۲۷</sup> برآورد ML و در روش بیزی برآورد ماکسیمم احتمال پسین<sup>۲۸</sup>(MAP) و میانگین احتمال پسین<sup>۲۹</sup>(PM) هستند.

در استنتاج تقریبی مجموعه داده‌ها کامل نیست و برای برآورد پارامترها از روش‌های تقریبی مانند الگوریتم امید ریاضی-ماکسیمم سازی<sup>۳۰</sup>(EM) استفاده می‌شود. الگوریتم EM یکی از تکنیک‌های پر کاربرد در به دست آوردن برآورد نقطه‌ای<sup>۳۱</sup> پارامترهای یک توزیع از یک مجموعه داده‌ی ناکامل(دارای مقادیر گمشده) است. در هر تکرار الگوریتم EM دو گام وجود دارد. ابتدا گام  $E$  که مرحله امیدگیری از تابع درستنمایی است. سپس گام  $M$  که مرحله ماکسیمم‌سازی امید تابع درستنمایی است (فریدمن، ۱۹۹۷).

<sup>۲۷</sup>Maximum Likelihood

<sup>۲۸</sup>Maximum A Posteriori

<sup>۲۹</sup>Posterior Mean

<sup>۳۰</sup>Expectation-Maximization

<sup>۳۱</sup>Point Estimation

$$\theta_{ijk} = P\left(X_i = x_i^{(k)} | \Pi_i = \pi_i^{(j)}\right),$$

که در آن  $1 = \sum_{k=1}^{r_i} \theta_{ijk}$  است.

تابع احتمال برای یک نمونه‌ی  $\underline{x}_{(l)}$  از دگ  $(V, E)$  به صورت زیر است:

$$p_{\underline{X}|\Theta}(\underline{x}_{(l)}|\theta, E) = \prod_{i=1}^n \prod_{j=1}^{q_i} \prod_{k=1}^{r_i} \theta_{ijk}^{n_{ijkl}}.$$

بنابراین تابع احتمال توأم برای نمونه‌های مستقل  $\underline{x}_{(1)}, \dots, \underline{x}_{(m)}$  به صورت زیر خواهد بود:

$$p_{\underline{X}|\Theta}(\underline{x}|\theta, G) = \prod_{i=1}^n \prod_{j=1}^{q_i} \prod_{k=1}^{r_i} \prod_{l=1}^m \theta_{ijk}^{n_{ijkl}} = \prod_{i=1}^n \prod_{j=1}^{q_i} \prod_{k=1}^{r_i} \theta_{ijk}^{n_{ijk}}$$

که در آن  $n_{ijkl} = \sum_{l=1}^m n_{ijkl}$  است. بنابراین برآورد ML پارامتر  $\theta_{ijk}$  از رابطه‌ی زیر به دست می‌آید:

$$\delta_{ijk}^{ML} = \frac{n_{ijk}}{n_{ij}}, \quad (1)$$

که در آن  $n_{ij} = \sum_{k=1}^{r_i} n_{ijk}$  است (برای مشاهده‌ی جزئیات بیشتر کوسکی و نوبل، ۲۰۰۹ را بینید).

به منظور به دست آوردن برآورد MAP و PM پارامتر  $\theta_{ijk}$  برای

$k = 1, \dots, r_i$  توزیع پیشین مزدوج

$$(\theta_{ij1}, \dots, \theta_{ijr_i}) \sim Dir(\alpha_{ij1}, \dots, \alpha_{ijr_i})$$

را با تابع چگالی احتمال زیر را در نظر بگیرید:

$$\pi(\theta_{ij1}, \dots, \theta_{ijr_i}) = \frac{\Gamma(\alpha_{ij})}{\prod_{k=1}^{r_i} \Gamma(\alpha_{ijk})} \prod_{k=1}^{r_i} \theta_{ijk}^{\alpha_{ijk}-1},$$

که در آن  $1 < \theta_{ijk} < 1$  و  $\alpha_{ij} = \sum_{k=1}^{r_i} \alpha_{ijk}$ .  $\sum_{k=1}^{r_i} \theta_{ijk} = 1$  و  $\alpha_{ijk} > 0$

تابع احتمال پسین برای مجموعه داده‌ی  $\underline{x}_{(1)}, \dots, \underline{x}_{(m)}$  از رابطه‌ی زیر به دست می‌آید:

$$\pi(\theta_{ij1}, \dots, \theta_{ijr_i} | \underline{x}) \propto \pi(\theta_{ij1}, \dots, \theta_{ijr_i}) p(\underline{x} | \theta_{ij1}, \dots, \theta_{ijr_i}) \propto \prod_{k=1}^{r_i} \theta_{ijk}^{n_{ijkl} + \alpha_{ijk}-1},$$

یعنی

$(\theta_{ij1}, \dots, \theta_{ijr_i}) | \underline{x} \sim Dir(n_{ij1} + \alpha_{ij1}, \dots, n_{ijr_i} + \alpha_{ijr_i})$  به وضوح توزیع پسین حاشیه‌ای به صورت  $\theta_{ijk} | \underline{x} \sim Beta(n_{ijk} + \alpha_{ijk}, n_{ij} + \alpha_{ij} - n_{ijk} - \alpha_{ijk})$  است.

بنابراین برآورد MAP برای  $\theta_{ijk}$  عبارت از:

$$\delta_{ijk}^{\pi, MAP}(x) = \frac{n_{ijk} + \alpha_{ijk} - 1}{n_{ij} + \alpha_{ij} - 2}. \quad (2)$$

$K_2$  بهترین ساختار را می‌یابیم و سپس برآورد ML، MAP و یا PM پارامترها را به دست می‌آوریم.

- در صورتی که ساختار شبکه معلوم نباشد و فقط برخی از مقادیر متغیرها قابل مشاهده باشند، ابتدا با روش‌های گشت (مانند روش گشت  $K_2$ ) بهترین ساختار را می‌یابیم و سپس از الگوریتم EM برای برآورد پارامترها استفاده می‌کیم.

## ۲.۴ یادگیری پارامتری شبکه بیزی با روش ماکسیمم درستنمایی و روش بیزی

دگ  $(V, E)$  را در نظر بگیرید که در آن  $V = \{X_1, \dots, X_n\}$  مجموعه‌ی از  $n$  متغیر تصادفی و  $E$  مجموعه‌ی از یال‌های جهت‌دار درون فضای  $V \times V$  است. فرض کنید متغیر  $X_i$  برای  $\mathcal{X}_i = \{x_i^{(1)}, \dots, x_i^{(r_i)}\}$   $i = 1, \dots, n$  مقادیرش را از مجموعه‌ی  $\{x_1^{(1)}, \dots, x_n^{(r_i)}\}$  بگیرد و مجموعه‌ی همه‌ی برآوردهای ممکن آزمایش به صورت زیر باشد:

$$\begin{aligned} \mathcal{X} &= \mathcal{X}_1 \times \dots \times \mathcal{X}_n \\ &= \left\{ (x_1^{(i_1)}, \dots, x_n^{(i_n)}) | i_j = 1, \dots, r_j, j = 1, \dots, n \right\}. \end{aligned}$$

مجموعه‌ی داده‌ها را به صورت یک نمونه‌ی  $m$  تایی از متغیرهای تصادفی  $X_1, \dots, X_n$  به صورت زیر در نظر بگیرید:

$$x = \begin{pmatrix} \underline{x}_{(1)} \\ \vdots \\ \underline{x}_{(m)} \end{pmatrix},$$

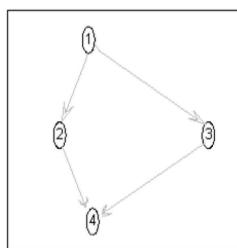
که در آن  $\underline{x}_{(l)} = (x_{l,1}, \dots, x_{l,n})$  بیانگر  $l$ -امین نمونه است. فرض کنید مجموعه‌ی  $\Pi_i$  دارای  $q_i$  ترکیب از اجداد متغیر  $X_i$  به صورت  $\Pi_i = \{\pi_i^{(q_1)}, \dots, \pi_i^{(q_i)}\}$  باشد که در آن به این نکته اشاره می‌کند که  $j$ -امین ترکیب از  $\Pi_i$  مشاهده شده است.

متغیر برنولی  $n_{ijkl}$  را به صورت زیر تعریف می‌کنیم:

$$n_{ijkl} = \begin{cases} 1, & \text{if } (x_i^{(k)}, \pi_i^{(j)}) \text{ is found in } \underline{x}_{(l)} \\ 0, & \text{otherwise} \end{cases}$$

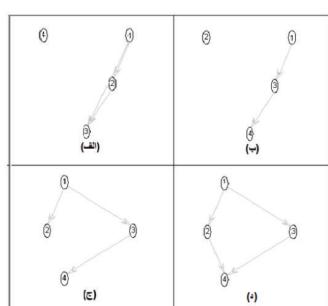
که در آن  $(x_i^{(k)}, \pi_i^{(j)})$  یک ترکیب از خانواده‌ی  $(X_i, \Pi_i)$  است. فرض کنید  $\theta$  مجموعه‌ی پارامترهای تعریف شده برای  $i = 1, \dots, n$  و  $j = 1, \dots, q_i$  به صورت زیر باشد:

از مقایسه‌ی این شبکه با شبکه واقعی (شکل ۲) ملاحظه می‌کنیم که این شبکه با شبکه واقعی اختلاف زیادی دارد. این نتیجه دور از انتظار نبود زیرا روش‌های محدودیت‌گرا در موقعی که تعداد نمونه کم باشد خطاهای آماری دارند. از آن جایی که روش‌های امتیاز‌گرا نسبت به روش‌های مبتنی بر قید جواب‌های دقیق‌تری دارند و موقعی که تعداد نمونه‌ها کم باشد قابلیت اجرا دارند، به تعیین ساختار شبکه بیزی این مجموعه داده با الگوریتم  $K_2$ ، با متر بیزی  $K_2$  و متر غیر بیزی BIC از روش‌های امتیاز‌گرا می‌پردازم. به منظور اجرای الگوریتم  $K_2$ ، متغیرها را (به دلخواه) به صورت  $C$ ،  $S = 3$ ،  $R = 4$  و  $W = 4$  شماره‌گذاری کرده و داده‌ها را به صورت  $10 \times 10$  مجموعه داده به حجم‌های  $5, 10, 15, \dots, 45$  و  $50$  در نظر می‌گیریم. سپس الگوریتم  $K_2$  را با در نظر گرفتن عدد  $2$  به عنوان ماکسیمم تعداد اجداد هر رأس اجرا می‌کنیم. به منظور مقایسه شبکه واقعی (شکل ۲) با شبکه‌های بازسازی شده توسط الگوریتم  $K_2$ ، متغیرها را به صورت  $C = 1, S = 2, R = 3$  و  $W = 4$  شماره‌گذاری کرده و ساختار شبکه شکل ۲ را با کدهای نرم افزار متلب، در شکل ۴ رسم کرده‌ایم.



شکل ۴: دگ متشكل از ۴ متغير دو وضعیتی  $S = 2, W = 1, C = 4$  و  $R = 3$

ابتدا الگوریتم  $K_2$  را بر اساس متر بیزی  $K_2$  اجرا می‌کنیم و مشخص می‌کنیم که شبکه حاصل از کدام حجم نمونه با شبکه اولیه یکسان است. خروجی الگوریتم به صورت شکل ۵ است.



شکل ۵: DAG حاصل از اجرای الگوریتم  $K_2$  با متر بیزی  $K_2$

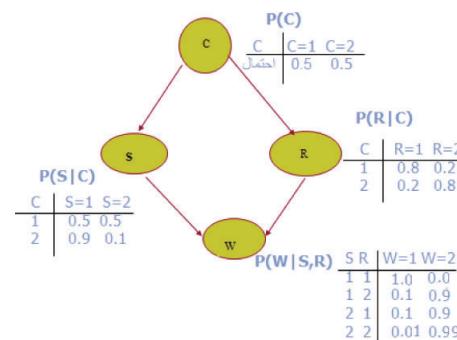
است و برآورد PM برای  $\theta_{ijk}$  عبارت از:

$$\delta_{ijk}^{\pi, PM}(x) = \frac{n_{ijk} + \alpha_{ijk}}{n_{ij} + \alpha_{ij}}. \quad (3)$$

است.

## ۵ مدل‌بندی یک مجموعه داده شبیه‌سازی شده با شبکه بیزی

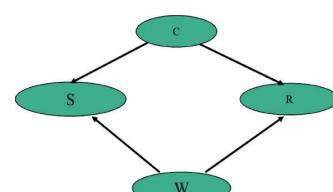
در این بخش یک مجموعه داده شبیه‌سازی شده را با شبکه بیزی مدل‌بندی می‌کنیم. برای این منظور ابتدا ساختار شبکه بیزی را با روش‌های محدودیت‌گرا و امتیاز‌گرا به ترتیب با الگوریتم‌های یادگیری ساختاری PC و  $K_2$  تعیین می‌کنیم و سپس به مسئله‌ی یادگیری پارامتری شبکه بیزی حاصل با روش ماکسیمم درستنمایی و روش بیزی می‌پردازیم. شبکه بیزی متشكل از ۴ رأس  $C = Wetgrass, R = Rain, S = Sprinkler, Cloudy$  و شکل ۲ را در نظر بگیرید.



شکل ۲: شبکه بیزی از ۴ متغير دو وضعیتی  $C = Cloudy$ ,  $.W = Wetgrass$  و  $R = Rain$ ,  $S = Sprinkler$

### ۱.۵ یادگیری ساختاری

ابتدا الگوریتم PC روش‌های محدودیت‌گرا را بر روی ۱۰۰ داده‌ی شبیه‌سازی شده از شبکه شکل ۲ اجرا می‌کنیم. خروجی الگوریتم به صورت شکل ۳ است.



شکل ۳: دگ حاصل از اجرای الگوریتم PC روی داده‌های شبیه‌سازی شده به حجم ۱۰۰ از شبکه شکل ۲

صورت متغیرهای  $X_1 = C = X_2$ ,  $X_3 = S = X_2$ ,  $X_4 = R = X_3$  و  $W = X_4$  در نظر بگیرید. همهٔ متغیرها دو وضعیتی هستند و مقادیر ۱ و ۲ را می‌گیرند. بنابراین طبق نمادهای ذکر شده در بخش قبل،  $n = 4$ ,  $q_4 = 4$ ,  $q_2 = 2$ ,  $r_1 = r_2 = r_3 = r_4 = 2$  است. با توجه به شکل ۲،  $q_1 = 0$  است یعنی مجموعهٔ جد برای متغیرها به صورت زیر است:

$$\Pi_1 = \{\}, \quad \Pi_2 = \{\pi_2^{(1)}, \pi_2^{(2)}\}$$

$$\Pi_3 = \{\pi_3^{(1)}, \pi_3^{(2)}\}, \quad \Pi_4 = \{\pi_4^{(1)}, \dots, \pi_4^{(4)}\},$$

که در آن:

$$\pi_2^{(1)} = 1, \pi_2^{(2)} = 2, \pi_3^{(1)} = 1, \pi_3^{(2)} = 2, \pi_4^{(1)} = (1, 1),$$

$$\pi_4^{(2)} = (1, 2), \pi_4^{(3)} = (2, 1), \pi_4^{(4)} = (2, 2).$$

در اینجا با در نظر گرفتن مشاهدات متغیر  $X_3$  از ۱۰۰ نمونه  $m = 100$  شبیه‌سازی شده، برآورد ML, MAP و PM پارامترهای  $\theta_{311}$  و  $\theta_{321}$  را از رابطه‌های (۱)، (۲) و (۳) به دست می‌آوریم. این برآوردها در جدول ۲ نشان داده شده‌اند.

جدول ۲: مقدار واقعی و برآورد ML, MAP و PM پارامترهای

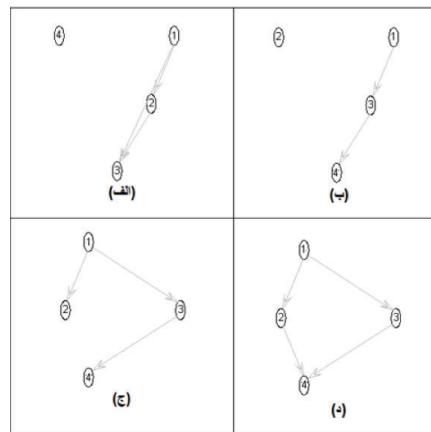
	$\theta_{321}$ و $\theta_{311}$	پارامتر	مقدار واقعی	مقدار برآورد
$\theta_{311}$	۰/۸۰۰	ML	۰/۸۱۱	۰/۸۱۸
$\theta_{312}$	۰/۲۰۰	MAP	۰/۲۸۶	۰/۲۵۹
		PM		۰/۲۶۵

لازم به ذکر است که به منظور به دست آوردن برآوردهای PM و MAP این پارامترها، توزیع پیشین مزدوج  $(\theta_{311}, \theta_{321}, \theta_{312}, \theta_{322}) \sim Beta(16, 4, 4, 16)$

را بر اساس دانش پیشین خود طوری در نظر گرفته‌ایم که میانگین توزیع حاشیه‌ای هر پارامتر با مقدار واقعی پارامتر یکسان باشد. به منظور کاهش محاسبات، برآورد ML, MAP و PM پارامترهای  $\theta_{312}$  و  $\theta_{322}$  را به جای محاسبهٔ مستقیم از رابطه‌های (۱)، (۲) و (۳)، به کمک رابطه‌های  $\theta_{312} = 1 - \theta_{311}$  و  $\theta_{322} = 1 - \theta_{321}$  محاسبه می‌کنیم. برآورد این پارامترها در جدول ۳ آمده است. بقیهٔ پارامترهای شبکه نیز به طور مشابه برآورد می‌شوند.

روی داده‌های شبیه‌سازی شده از شبکه شکل ۲، به (الف) حجم ۵، (ب) حجم ۱۰، (ج) حجم‌های ۱۵، ۲۰ و ۲۵، (د) حجم‌های ۳۰، ۳۵، ۴۰ و ۵۰.

از مقایسه این شبکه‌ها با شبکه واقعی ملاحظه می‌کنیم که شبکه‌های حاصل از نمونه‌های به حجم‌های ۳۰، ۳۵، ۴۰ و ۵۰ با شبکه واقعی یکسانند. پس برای ساختن یک شبکه بیزی با الگوریتم  $K_2$  با متر بیزی  $K_2$ ، از چهار متغیر دو وضعیتی، کافی است یک نمونه به حجم ۳۰ از متغیرها داشته باشیم. از آنجایی که متر بیزی  $K_2$  نا آگاهی‌بخش است یک بار دیگر الگوریتم  $K_2$  را بر اساس متر غیر بیزی BIC اجرا می‌کنیم و مشخص می‌کنیم که شبکه حاصل از کدام حجم نمونه با شبکه اولیه یکسان است. خروجی الگوریتم به صورت شکل ۶ است.



شکل ۶: دگ حاصل از اجرای الگوریتم  $K_2$  با متر غیر بیزی BIC روی داده‌های شبیه‌سازی شده از شبکه شکل ۲، به (الف) حجم ۵، (ب) حجم ۱۰، (ج) حجم‌های ۱۵، ۲۰ و ۲۵ (د) حجم‌های ۳۰، ۳۵، ۴۰ و ۵۰.

از مقایسه این شبکه‌ها با شبکه واقعی ملاحظه می‌کنیم که شبکه‌های حاصل از نمونه‌های به حجم‌های ۳۰، ۳۵، ۴۰ و ۵۰ با شبکه واقعی یکسانند. پس برای ساختن یک شبکه بیزی با الگوریتم  $K_2$  با استفاده از متر غیر بیزی BIC، از چهار متغیر دو وضعیتی، کافی است یک نمونه به حجم ۴۰ از متغیرها داشته باشیم.

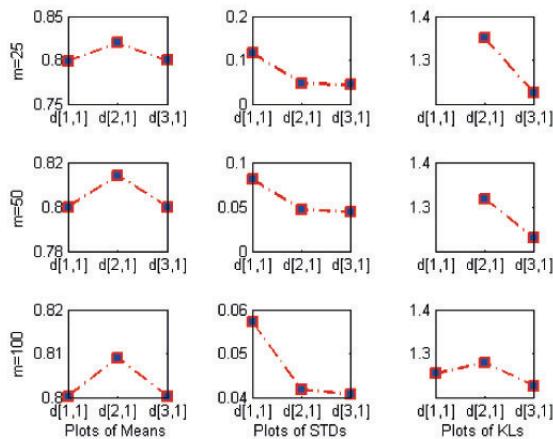
## ۲.۵ یادگیری پارامتری

اکنون که ساختار شبکه بیزی مشخص شد، به مسئلهٔ یادگیری پارامتری شبکه بیزی حاصل از پردازیم. رأس‌های شبکه را به

جدول ۴: آماره‌های کمی میانگین و انحراف استاندارد و معیار واگرایی KL برآوردهای پارامتر  $\theta_{311}$  برای مقادیر متفاوت  $m$ .

	$m$	$d[1, 1]$	$d[2, 1]$	$d[3, 1]$
Mean	25	0/7991	0/8195	0/7997
		0/1163	0/0464	0/0435
		*	1/3507	1/2261
STD	50	0/8000	0/8141	0/8000
		0/0812	0/0465	0/0444
		*	1/3171	1/2295
KL	100	0/8000	0/8090	0/8001
		0/0571	0/0418	0/0406
		1/2542	1/2793	1/2258

نمودارهای شکل ۷ نشان می‌دهند که عملکرد برآوردها با تغییر حجم نمونه به مقادیر چشمگیری تغییر نمی‌کند. به عبارت دیگر، دقت برآوردها تحت تأثیر حجم نمونه نیست و حجم نمونه ۲۵ نیز برای به دست آوردن برآورد دقیق پارامتر  $\theta_{311}$  کافی است. به منظور انجام یک استنتاج استقرایی، نمودارهای میانگین تکرار برآوردهای پارامتر  $\theta_{311}$  برای مقادیر متفاوت  $N = 10000$  در شکل ۷ به تصویر کشیده شده است.



شکل ۷: نمودارهای میانگین برآوردهای ML، MAP و PM و معیار واگرایی KL پارامتر  $\theta_{311}$  برای مقادیر متفاوت  $m$ .

از جدول ۴ و شکل ۷ مشاهده می‌کنیم که برآوردهای PM دقیق‌تر از MAP پارامتر  $\theta_{311}$  را برآورد می‌کنند. این مطلب در مقایسه مقادیر انحراف استاندارد نیز به وضوح قابل مشاهده است. به علاوه معیار واگرایی KL برآوردهای PM خیلی کمتر از MAP است.

جدول ۳: مقدار واقعی و برآورد ML، MAP و PM پارامترهای  $\theta_{322}$  و  $\theta_{312}$

پارامتر	مقادیر واقعی	ML	برآورد	MAP	PM
$\theta_{321}$	0/200	0/189	0/182	0/193	0/193
$\theta_{322}$	0/800	0/714	0/741	0/735	0/735

اکنون می‌خواهیم عملکرد برآوردهای ML، MAP و PM پارامتر  $\theta_{311}$  را با آماره‌های میانگین (Mean) و انحراف استاندارد (STD) برآوردها و معیار واگرایی کولبک-لایبلر (KL) با یکدیگر مقایسه کنیم. برای این منظور، مراحل زیر را در نظر می‌گیریم:

گام ۱. متغیرهای  $(x_1, \dots, x_4)$  را در  $m = 25, 50, 100$  نمونه شبیه‌سازی می‌کنیم.

گام ۲.  $d[3, k] = \delta_{31k}^{\pi, PM}$  و  $d[2, k] = \delta_{31k}^{\pi, MAP}$ ،  $d[1, k] = \delta_{31k}^{ML}$  را برای هر  $k = 1, 2$ ، با در نظر گرفتن توزیع پیشین مزدوج  $Beta(16, 4)$  برای پارامترهای  $(\theta_{311}, \theta_{312})$  محاسبه می‌کنیم.

گام ۳. گام‌های ۱ و ۲ را  $N = 10000$  بار تکرار می‌کنیم. سپس بر اساس داده‌های تولید شده، آماره‌های کمی میانگین و انحراف استاندارد و واگرایی کولبک-لایبلر هر برآورد را از رابطه‌های زیر محاسبه می‌کنیم:

$$\begin{aligned} \text{Mean } d[i, k] &= \frac{1}{N} \sum_{r=1}^N d[i, k, r], \\ \text{STD } d[i, k] &= \left( \frac{1}{N-1} \sum_{r=1}^N (d[i, k, r] - \text{Mean } d[i, 1])^2 \right)^{\frac{1}{2}}, \\ \text{KLD } d[i] &= \frac{1}{N} \sum_{r=1}^N (\theta_{311} \log_2(\theta_{311}/d[i, 1, r]) \\ &\quad + \theta_{312} \log_2(\theta_{312}/d[i, 2, r])) \end{aligned}$$

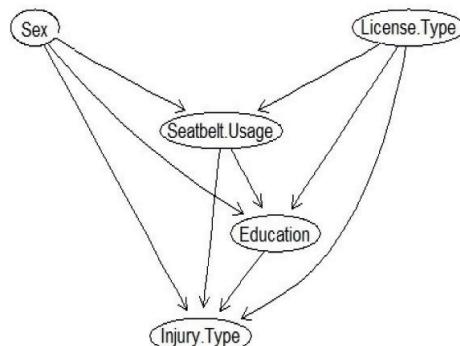
که در آن  $d[i, k, r]$  و  $k = 1, 2$ ،  $i = 1, \dots, 3$  برای  $d[i, k, r]$  در  $r = 1, \dots, N$  تکرار امین برآورد است.

آماره‌های کمی میانگین و انحراف استاندارد و معیار واگرایی KL برآوردهای پارامتر  $\theta_{311}$  برای مقادیر متفاوت  $m$ ، محاسبه شده و در جدول ۴ آورده شده است. علامت \* در این جدول بیانگر این است که معیار واگرایی KL قابل محاسبه نیست.

همان‌گونه که در شکل ۸ مشاهده می‌کنیم، گره‌های جنسیت، میزان تحصیلات، نوع گواهینامه و استفاده یا عدم استفاده از کمربند ایمنی، همگی جدهای گره نوع آسیب‌دیدگی هستند. حال فرض کنید یادگیری پارامتری در این شبکه مورد علاقه باشد و بخواهیم بدانیم که به طور مثال، احتمال عدم آسیب‌دیدگی راننده مردی با تحصیلات دیپلم و گواهینامه پایه دوم رانندگی در حالی که از کمربند ایمنی استفاده کرده است چقدر است. با توجه به ساختار شبکه تصادف به دست آمده در شکل ۸ و با استفاده از داده‌های واقعی در دسترس، برآورد ماکسیمم درستنایی این میزان احتمال را محاسبه نموده‌ایم که مقدار آن  $0.98299$  به دست آمده است. بدیهی است می‌توان برآورد ماکسیمم درستنایی سایر پارامترهای این شبکه را نیز به دست آورد. به همین ترتیب در صورت وجود اطلاعات قابل استناد و صحیح، می‌توان توزیع‌های پیشین مناسبی را (مشابه آنچه در زیربخش ۲.۵ بیان شد) تعیین نمود و برآوردهای MAP و PM را نیز گزارش نمود.

## ۶ پیاده‌سازی بر روی داده‌های واقعی

در این بخش به پیاده‌سازی روش‌های یادگیری ساختاری و پارامتری بر روی یک مجموعه داده واقعی تصادف می‌پردازیم. داده‌های تصادف شامل اطلاعات عمومی رانندگانی است که در شبکه جاده‌ای کشور دچار سانحه تصادف رانندگی شده‌اند. این داده‌ها مربوط به یک دوره زمانی ۵۶ ماهه از فروردین ۱۳۸۸ تا آذر ۱۳۹۲ است. در این دوره زمانی حدود ۶۵۵۰۰ تصادف جاده‌ای به ثبت رسیده است (کریم‌نژاد، ۱۳۹۳). متغیرهای مورد استفاده عبارت از جنسیت<sup>۳۳</sup> (مرد، زن)، میزان تحصیلات<sup>۳۴</sup> (بی‌سود، زیردیپلم، دیپلم، فوق دیپلم، کارشناسی، کارشناسی ارشد و بالاتر)، استفاده یا عدم استفاده از کمربند ایمنی<sup>۳۵</sup>، نوع آسیب‌دیدگی<sup>۳۶</sup> (آسیب ندیده، جرحی، فوتی) و نوع گواهینامه<sup>۳۷</sup> (بدون گواهینامه، گواهینامه پایه دوم، گواهینامه پایه اول، گواهینامه مشروط، گواهینامه ویژه) هستند. با توجه به اینکه حجم داده‌ها در این مطالعه زیاد است، برای یادگیری ساختار شبکه تصادف از الگوریتم PC بهره می‌گیریم (همان‌گونه که قبل از بیان شد، وقتی حجم داده‌ها کم باشد، روش‌های محدودیت‌گرا با خطاهای آماری مواجه هستند اما اگر حجم داده‌ها به اندازه کافی بزرگ باشد، می‌توان از روش‌های محدودیت‌گرا در یادگیری ساختار شبکه بهره گرفت). شکل ۸ ساختار نهایی حاصل از اجرای الگوریتم PC بر روی داده‌های تصادف نشان می‌دهد.



شکل ۸: ساختار به دست آمده با استفاده از الگوریتم PC

<sup>۳۳</sup>Sex

<sup>۳۴</sup>Education

<sup>۳۵</sup>Seatbelt Usage

<sup>۳۶</sup>Injury Type

<sup>۳۷</sup>License Type

کافی است. در آخر نشان دادیم که برآوردهای PM خیلی دقیق‌تر از بیزی با الگوریتم  $K_2$  از چهار متغیر دو وضعیتی، به ترتیب با متر بیزی  $K_2$  و متر غیر بیزی  $BIC$ ، کافی است نمونه‌های به حجم‌های ۴۰ و ۴۰ از متغیرها داشته باشیم. پس از مشخص شدن ساختار، به مسئله‌ی یادگیری پارامتری شبکه بیزی حاصل با روش ماکسیمم درستنمایی و روش بیزی پرداختیم. مشاهده کردیم که عملکرد برآوردها با تغییر حجم نمونه به مقدار چشمگیری تغییر نمی‌کند. به عبارت دیگر، دقت برآوردها تحت تأثیر حجم نمونه نیست و حجم نمونه‌ی  $m = 25$  نیز برای به دست آوردن برآورد دقیق پارامترها رانندگی شدنده، پرداختیم.

## مراجع

- [۱] ابراهیمی، علی. تنظیم شبکه بیان ژن بر مبنای شبکه بیزی؛ پایان‌نامه کارشناسی ارشد، دانشگاه شهید بهشتی، ۱۳۹۰.
- [۲] کریم‌نژاد، علی. تحلیل داده‌های تصادف جاده‌ای با استفاده از مدل‌های آماری؛ طرح پژوهشی دفتر تحقیقات کاربردی پلیس راهور، تهران، ۱۳۹۳.
- [۳] Friedman, N, (1997), *Learning Belief Networks in the Presence of Missing Values and Hidden Variables*; ICML.
- [۴] Friedman, N, (2004), *Inferring cellular networks using probabilistic graphical models*, Science, **303**:799-805.
- [۵] Heckerman, D, (1996), *A Tutorial on Learning with Bayesian Networks*, Technical Report, Microsoft Research.
- [۶] Heckerman, D, Geiger, D and Chickering, D. M, (1995), *Learning BNs: the combination of knowledge and statistical data*, Machine Learning, Vol. 20, pp. 197–243.
- [۷] Jensen, F. V and Nielsen, T. D, (2007), *Bayesian Networks and Decision Graphs*; Second Edition, Springer Science +Business Media, LLC.
- [۸] Koski, T and Noble, J.M, (2009), *Bayesian Networks - An Introduction*; Wiley.
- [۹] Murphy, K and Saira, M, (2001), *Modelling Gene Expression Data Using Dynamic Bayesian Networks*; Computer Science Division, University of California, Life Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720.
- [۱۰] Ong, I. M, Glasner, J. D and Page, D, (2002), *Modelling regulatory pathways in E. Coli from time series expression profiles*, Bioinformatics, Vol. 18, pp. S241–S248.
- [۱۱] Pearl, J, (1988), *Probabilistic Reasoning in Intelligent Systems*, Pacific Symposium, On Biocomputing.