



دانشگاه صنعتی امیر کبیر
(پلی تکنیک تهران)

گروه مستقل مهندسی رباتیک

گزارش تمرین اول درس مدل‌های احتمالاتی گرافی

استاد درس:

دکتر نیک آبادی

تدریس بار:

مهندس طاهرخانی

نام دانشجو:

نوید خزاعی

۹۲۱۳۵۰۰۸

اردیبهشت ۹۴

فهرست مطالب

۱	بخش نظری	۱
۱	۱.۱ سوال اول	۱
۲	۲.۱ سوال دوم	۲
۲	۱.۲.۱ توزیع توام	۲
۲	۲.۲.۱ درستی و نادرستی	۲
۳	۳.۲.۱ Markov Blanket for X_3	۳
۳	بخش پیاده‌سازی	۳
۳	۱.۲ مقدمه	۳
۴	۲.۲ پایگاه داده و پردازش ها	۴

۱ بخش نظری

۱.۱ سوال اول

قضیه زیر را اثبات کنید:

Theorem 1. Let \mathcal{G} be a Bayesian Network structure over a set of random variables \mathcal{X} and let P be a joint distribution over \mathcal{X} . If P factorizes according to \mathcal{G} , then \mathcal{G} is an I-map for P .

پاسخ:

با توجه به این که P روی گراف \mathcal{G} فاکتورایز می‌شود، پس می‌دانیم که:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid \text{Pa}(X_i)), \quad (۱)$$

که در آن $\text{Pa}(X_i)$ والدین متغیر تصادفی i م هستند. برای آن که نشان دهیم گراف \mathcal{G} یک I-map برای P است، باید نشان دهیم $I(\mathcal{G}) \subseteq I(P)$. با توجه به آن چه در درس آمده بود، اگر مجموعه روابط استقلال $I(\mathcal{G})$ را به گونه‌ای با استفاده از نتایج فرمول ۱ نشان دهیم، یعنی این روابط از روابط استقلال توزیع مورد نظر استخراج شده‌اند پس زیرمجموعه‌ی آن نیز هستند. برای این کار، با توجه به این که مجموعه روابط استقلال در گراف ایجاب می‌کند که:

$$\{X_i \perp \text{ND}(X_i) \mid \text{Pa}(X_i); i = 1, \dots, n\} \quad (۲)$$

که در آن $\text{ND}(X_i)$ مجموعه‌ی غیرنسل^۱ متغیر (X_i) است. اگر بتوانیم نشان دهیم که:

$$P(X_i \mid \text{ND}(X_i)) = P(X_i \mid \text{Pa}(X_i)) \quad (۳)$$

آن‌گاه رابطه‌ی استقلال ۲ برقرار است. برای این کار $P(X_i \mid \text{ND}(X_i))$ را محاسبه می‌کنیم (نسل X_i با $\text{D}(X_i)$ نشان داده شده است):

$$\begin{aligned} P(X_i \mid \text{ND}(X_i)) &= \frac{P(X_i, \text{ND}(X_i))}{P(\text{ND}(X_i))} = \frac{\sum_{\text{D}(X_i)} P(X_1, \dots, X_n)}{\sum_{X_i, \text{D}(X_i)} P(X_1, \dots, X_n)} \\ &= \frac{\sum_{\text{D}(X_i)} \prod_{j=1}^n P(X_j \mid \text{Pa}(X_j))}{\sum_{X_i, \text{D}(X_i)} \prod_{j=1}^n P(X_j \mid \text{Pa}(X_j))} \end{aligned} \quad (۴)$$

^۱ Non-descendant

در ۴ در محاسبه‌ی سیگما، مواردی هستند که جزو نسل X_i نیستند، پس می‌توان آن‌ها را از سیگمای روی نسل آن خارج نمود. همچنین آن‌چه باقی می‌ماند:

$$\sum_{D(X_i)} \prod_{X_j \in D(X_i)} P(X_j | \text{Pa}(X_j)) = 1$$

مشابه همین استدلال برای مخرج نیز پاسخ‌گو است، لذا خواهیم داشت:

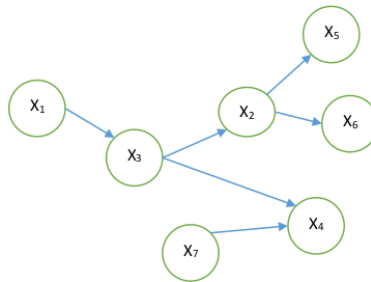
$$= \frac{\prod_{X_j \in (\text{ND}(X_i) \cup X_i)} P(X_j | \text{Pa}(X_j)) \times 1}{\prod_{X_j \in \text{ND}(X_i)} P(X_j | \text{Pa}(X_j)) \times 1} \quad (5)$$

$$= P(X_i | \text{Pa}(X_i)).$$

پس نشان دادیم ۳ و در نتیجه ۲ برقرار است، پس می‌توان نتیجه گرفت که $I(\mathcal{G}) \subseteq I(\mathcal{P})$ و حکم اثبات می‌شود.

۲.۱ سوال دوم

با توجه به شبکه‌ی شکل ۱ به سوالات پاسخ داده شده‌است.



شکل ۱: شکل سوال دوم

۱.۲.۱ توزیع توام

$$P(X_1, \dots, X_7) = P(X_1)P(X_3 | X_1)P(X_2 | X_3)P(X_5 | X_2)P(X_6 | X_2)P(X_4 | X_3, X_7)P(X_7)$$

۲.۲.۱ درستی و نادرستی

• $X_1 \perp X_5 | X_2$: درست است. اگر X_2 را بدانیم مسیر فعال بین X_1 و X_5 غیر فعال می‌شود و مستقل می‌شوند.

- $X_4 \perp X_7 \mid X_2$: نادرست است. در ساختار وی-شکل بین X_2 و X_7 با دانستن X_4 وابستگی ایجاد می‌شود، و چون X_2 به پدرش X_3 وابسته است، با X_7 نیز در این شرایط وابسته است.

- $X_1 \perp X_5 \mid X_2$: درست است، چرا که دانستن X_3 ، استقلال بین X_2 و X_4 را ایجاد می‌کند و بنابراین از فرزندان X_2 یعنی X_5 نیز مستقل می‌شویم.

۳.۲.۱ Markov Blanket for X_3

برای این کار ابتدا جهت‌ها را حذف می‌کنیم و سپس فساد^۱ ها را حذف می‌کنیم. یعنی نباید گره‌ای باشد که والد یکی از فرزندان همین گره باشد، در صورت وجود به آن گره وصل می‌کنیم تا فساد از بین برود. به تمامی اتصالات موجود نیز وصل می‌کنیم، بنا بر این $MB(X_3) = \{X_1, X_2, X_4, X_7\}$ خواهد بود.

۲ بخش پیاده‌سازی

۱.۲ مقدمه

در این بخش باید با استفاده از پایگاه داده‌ی معرفی شده، چند مدل بیزین ارایه کنیم تا قادر به تشخیص بیماری قلبی باشد. برای پیاده‌سازی این بخش، تصمیم به توسعه‌ی ابزار جدیدی برای کار با مدل‌های احتمالاتی گرافیکی گرفتیم، که کاملاً با ابزارهای تحت خط فرمان سیستم عامل لینوکس و پوسته‌ی Bash کار می‌کند. این گونه ابزارها، معمولاً برای استفاده در مواردی که فایل‌های حجیمی نیاز به پردازش دارند و می‌خواهیم درگیر مسایل مدیریت حافظه و Caching نشویم کاربرد دارند. همچنین، دید مبتنی بر فایل در بسیاری موارد به کمک توسعه‌دهنده می‌آید و توسعه‌ی ابزارها تحت قالب Shell Script، انعطاف بسیار زیادی را برای برنامه‌نویس فراهم می‌کند. این گونه ابزارها معمولاً در فازهای یادگیری به دلیل نیاز به نوشتن مفرط در فایل‌ها هستند، به گونه‌ای که ابزار تولید شده در این تمرین عملاً کارایی برای آموزش روی داده‌ها ندارد که در تمرینات آتی به رفع این ایرادها می‌پردازیم. در عوض، از آن‌جا که

^۱ Immorality

به کل مسالهی آموزش و تست به دید مسالهی پردازش متن نگاه کرده‌ایم، در فاز تست سرعت پاسخ‌گویی بسیار بالا خواهد بود چرا که ابزارهای پروژه‌ی GNU مانند grep که برای جست‌وجو استفاده می‌شوند و یا awk، sed و cut که به دفعات بسیاری از آن‌ها استفاده نموده‌ایم، کارایی به شدت بالایی در پردازش فایل‌های عظیم دارند. به این ترتیب در پردازش پایگاه‌های داده‌ی بیش از حد بزرگ و زمانی که CPD ها بزرگ می‌شوند و نیاز به جست‌وجوی موثر داریم، کارا خواهند بود. همچنین انعطاف در تعیین نوع داده‌ای ویژگی‌ها و متغیرها از جمله مزایای دیگر این ابزار است و ورودی و خروجی‌های تولید شده به آسانی با دیگر ابزارهای تحت خط فرمان لینوکس قابل استفاده هستند.

در ادامه پس از توضیح چگونگی استفاده از پایگاه داده، به بررسی ابزار توسعه داده‌شده و پاسخ سوالات تمرین می‌پردازیم. ☺

۲.۲ پایگاه داده و پردازش‌ها

با دریافت پایگاه داده‌ی شهر کیولوند^۱

توضیحات	مقدار گسسته	ویژگی/نام متغیر	شماره ستون پایگاه
نشانه پایین دهه	۲،۳،۴،۵،۶،۷	سن	۱
	۱ و ۰	جنسیت /	۲

^۱ Cleveland