

تمرین کامپیوتری دوم – استنتاج علی

بهراد منیری
۹۵۱۰۹۵۶۴
bemoniri@live.com

دانشکده‌ی مهندسی برق – دانشگاه صنعتی شریف

۱ بخش اول – پیاده‌سازی الگوریتم pc

در این بخش الگوریتم PC را با آزمون فرض استقلال شرطی مبتنی بر Partial Correlation پیاده‌سازی می‌کنیم. برای پیاده‌سازی از زبان پایتون استفاده کرده و گراف را به کمک کتابخانه‌ی networkx می‌سازیم. برای محاسبات نیز از کتابخانه‌های numpy و sklearn استفاده می‌کنیم.

۱.۱ پیاده‌سازی آزمون فرض

برای پیاده‌سازی آزمون فرضیه، تابع test را نوشتیم: $\text{test}(\text{data}, x, y, s)$ ، این تابع، ماتریس داده‌ها و اندیس x, y و مجموعه‌ای از اندیس‌ها، S ، را در ورودی گرفته و آزمون فرضیه‌ی

$$X \perp\!\!\!\perp Y | S$$

را انجام می‌دهد و مقدار p-value را در خروجی بر می‌گرداند. هر اندیس نماینده‌ی ستون مورد نظر در ماتریس داده‌ها، یعنی یک متغیر است.

۲.۱ پیاده‌سازی الگوریتم pc

تابع pc پیاده‌سازی الگوریتم pc است که شبکه‌کد آن ارائه شده.

این تابع، ماتریس داده‌ها که هر ستون آن یک متغیر است و همچنین مقدار significance level یعنی α را در ورودی دریافت کرده و اسکلت گراف مولد داده‌ها را تخمین می‌زند.

۳.۱ بررسی عملکرد الگوریتم با دو Toy Example

برای بررسی عملکرد الگوریتم، چند SCM گاوسی-خطی ساخته و الگوریتم را روی آنها اجرا می‌کنیم. از هر متغیر یک میلیون سمپل می‌سازیم.

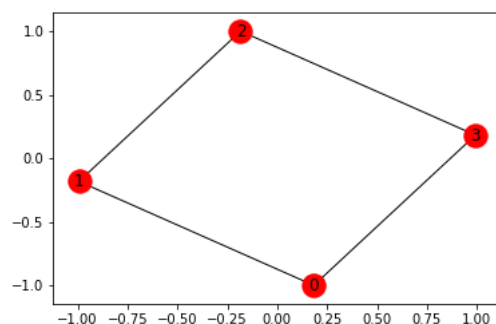
۱.۳.۱ مثال اول

$$\begin{cases} x_o = N_o \\ x_1 = x_o + N_y \\ x_2 = x_1 + N_z \\ x_3 = 2x_2 + 3x_o + N_w \end{cases} \quad N_1, N_2, N_3, N_4 \stackrel{iid}{\sim} \text{Normal}(0, 1) \quad (1)$$

برنامه‌ی ما، دو عبارت استقلال زیر را از داده استخراج می‌کند.

$$\begin{cases} x_o \perp\!\!\!\perp x_2 | x_1 \\ x_1 \perp\!\!\!\perp x_3 | x_o, x_2 \end{cases} \quad (2)$$

و دو یال متناظر را حذف می‌کند. شکل (۱) گراف تخمین زده شده است. تمام یال‌ها به درستی انتخاب شدند.



شکل ۱: گراف تخمین زده شده برای مثال اول

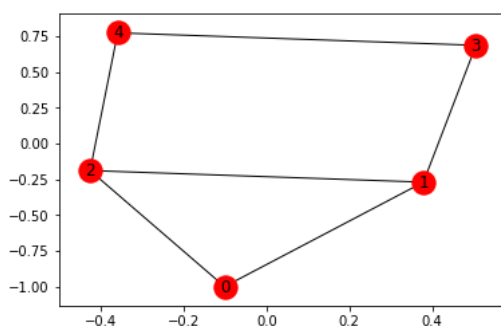
۲.۳.۱ مثال دوم

$$\begin{cases} x_0 = N_0 \\ x_1 = x_0 + N_1 \\ x_2 = 2x_0 + N_2 \\ x_3 = 3x_1 + N_3 \\ x_4 = 4x_2 + N_4 \end{cases} \quad N_0, N_1, N_2, N_3, N_4 \stackrel{ind.}{\sim} Normal(0, \sigma_i) \quad (3)$$

برنامه‌ی ما، دو عبارت استقلال زیر را از داده استخراج می‌کند.

$$\begin{cases} x_1 \perp\!\!\!\perp x_4 | x_2, x_3 \\ x_0 \perp\!\!\!\perp x_4 | x_1, x_2 \\ x_2 \perp\!\!\!\perp x_3 | x_1 \\ x_0 \perp\!\!\!\perp x_3 | x_1 \end{cases} \quad (4)$$

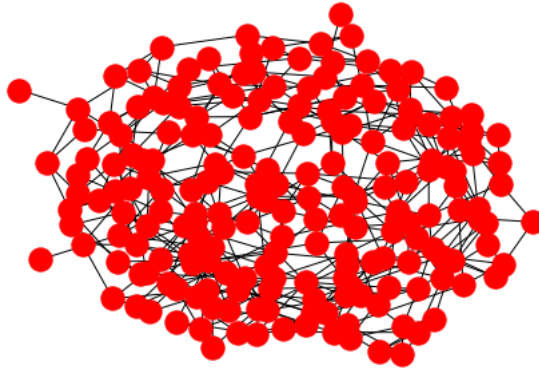
و یال‌های متناظر را حذف می‌کند. شکل (۲) گراف تخمین زده شده است. تمام یال‌ها به درستی انتخاب شدند.



شکل ۲: گراف تخمین زده شده برای مثال دوم

۲ بخش دوم – تخمین گراف مولد داده‌های ضمیمه

به کمک تابع pc اسکلت گراف را تخمین می‌زنیم. این الگوریتم بعد از سطح $l = 6$ متوقف شده و ۴۱۵ یال در این گراف باقی می‌ماند.



شکل ۳: خروجی الگوریتم

۳ بخش سوم – پیاده‌سازی الگوریتم pc پایدار

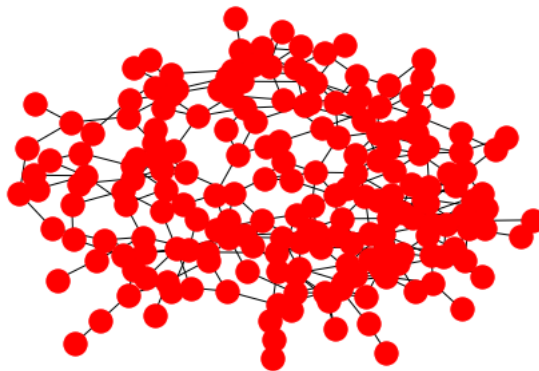
با تغییری کوچک در الگوریتم pc می‌توان به الگوریتم pc-stable رسید. در الگوریتم پایدار، بعد از اینکه تصمیم به قطع یک یال گرفتیم، یال مورد نظر را ذخیره می‌کنیم و در آخر هر سطح، تمام یال‌هایی که باید حذف شوند را یکجا حذف می‌کنیم. در الگوریتم اصلی این امکان وجود دارد که به دلیل حذف اشتباه یک یال، برخی از آزمون فرض‌ها را دیگر انجام ندهید و به ای دلیل خطا به شدت گسترش پیدا کند. در روش پایدار تا حدی این اثر کم‌رنگ می‌شود.

۱.۳ بررسی عملکرد الگوریتم پایدار با دو Toy Example

در این بخش الگوریتم pc-stable را بر روی داده‌های تولید شده توسط همان دو SCM بخش (۳.۱) اجرا می‌کنیم. مجدداً الگوریتم هر دو اسکلت را به درستی تخمین می‌زند.

۲.۳ اعمال الگوریتم پایدار بر داده‌های ضمیمه‌شده

با الگوریتم pc-stable اسکلت گراف مولد داده‌ها را تخمین می‌زنیم. این الگوریتم زمان بیشتری نسبت به الگوریتم pc می‌گیرد و در نهایت در آخر سطح ۶ متوقف شده و ۳۱۳ یال باقی می‌ماند.



شکل ۴: خروجی الگوریتم پایدار

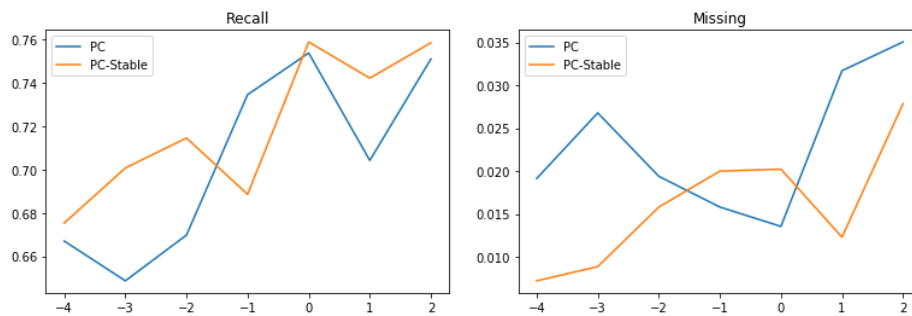
۴ بخش چهارم - بررسی اثر درصد اطمینان

در این بخش درصد اطمینان را در بازه‌ی ۲-۴٪ تا ۲۱٪ تغییر می‌دهیم و برای هر بازه‌ی اطمینان، ۲۰۰ بار الگوریتم مدنظر را اجرا می‌کنیم. در هر اجرا دو کمیت زیر را محاسبه می‌کنیم.

۱. Recall نسبت تعداد یال‌های گراف واقعی که درست تشخیص داده‌شده به تعداد کل یال‌های گراف واقعی

۲. Missing نسبت تعداد یال‌هایی که ما تشخیص دادیم درحالی که در گراف اصلی نبوده به تعداد کل یال‌های گراف واقعی

با افزایش ترشهود، مشاهده می‌شود که Recall و Missing افزایش می‌یابد. این موضوع بدیهی است زیرا با زیاد کردن ترشهود، با سخت‌گیری کمتری یال‌ها را حذف می‌کنیم پس درکل تعداد یال‌های حذف شده زیاد می‌شود و می‌توان به یال‌های باقی‌مانده اطمینان بیشتری داشت. دو نمودار زیر، نتایج این بخش هستند:



شکل ۵: مقایسه‌ی دو الگوریتم

الگوریتم پایدار عملکرد بهتری داشته زیرا در الگوریتم pc معمولی چندین مشکل وجود دارد. اول اینکه با تغییر دادن ترتیب آزمون‌ها، نتایج متفاوت خواهد بود. دوم اینکه با در صورت یک خطا در آزمون، تعدادی آزمون دیگر انجام نخواهند شد. این دو مشکل تا حد زیادی در pc-stable حل شده است. به همین دلیل تعداد یال‌های حذف‌شده در روش پایدار بیشتر از روش معمولی است.