

مانگنگر ز مدلها رگرسون

افشان فلاخ محسن محمدزاده

گروه مار دانشگاه تربیت مدرس

چکده در روشهای مدلسازی که سنته ه معارض انتخاب مدل و ز رتجل المگر مدلها می مقاومت انتخاب می شوند عدم حتمت در فریند مدلسازی نادله گرفته می شود مانگنگری زی مدلها شوهای قوانین در مدلسازی است که در نه هر مدل ه تناسب میزان حمایتی که از جاذب دادهها صورت می پذیرد وزنی اختصاص داده می شود و مانگن وزنی همه مدلها عنوان مدل نهایی کارگرفته می شود در ان مقاله نشان داده می شود که روش مانگنگری زی خای مدلسازی را کاهش و کارایی را افزایش می دهد

واژه ها کارد عدم حتمت توزیع پشن مانگنگری زی

مقدمه

معمولا در مدلسازی کلاسی از مدلها در ز رگرفته می شود سپس کی از نهای راساس که اچند معارض ارزای عنوان هترن مدل انتخاب می شود در حالی که ممکن است رقبای سار خوی رای مدل انتخابی در فای مدل وجود داشته باشد این شوه مدلسازی دارای نارسای های است که مهمترین نهای در ز رنگرفتن عدم حتمت مدلها انتخابی در فریند مدلسازی است مانگنگری زی مدلها¹ BMA شوهای است که در ن از تمام مدلها موجود در فای مدل رای دسته ای ه مدلی مناس استفاده می شود در ان روش راساس میزان حمایت دادهها از هر مدل وزنی ه ن اختصاص داده می شود سپس مدل حاصل از مانگن وزنی همه مدلها رای انجام استنباط و پشن نی ه کارگرفته می شود رشته تاریخی این روش ه مقاله لمر رمی گردد که ه دلیل محاسبات دشوار و پیچیده در ن زمان چندان مورد توجه قرار نگرفت ا پشرفت رانه ها و فنون محاسبات تقریبی در اولی دهه نود این روش دو امره بح و در کانون توجهات قرار گرفت مادگان و رفتاری و مادگان و ورک دو روش پاهای رای اجرای این شوه مدلسازی ارا کردند هوتنگ و همکاران ه نحوه انجام محاسبات و اجرای این روش پرداختند رورت و لاپکووج مانگنگری زی را در حالت چند متغیره مورد توجه قرار دادند در ان مقاله روش مانگنگری زی مدلها نحوه محاسبه مولفه های و روشهای اجرای این در بخش شرح داده

1) Bayesian Model Averaging

مجموعه مقالات

می شود و سپس در بخش BMA ارزایی روش پرداخته و نهاداً حد و نتیجه‌گری در بخش اراه خواهد شد

مانگنگر ز مدلها

فرض کند رای متغیر واسطه Y و مجموعه‌ای از متغیرهای پیش ن X_1, \dots, X_k هدف افتخاری مدل از ن همه مدلها خی

$$Y = \beta_0 + \sum_{j=1}^p \beta_j X_j + \varepsilon \quad (1)$$

ماشند که در ن X_{ip}, X_1, \dots, X_k ز مجموعه‌ای از متغیرهای X_1, \dots, X_k هستند وسیله $\hat{\Delta}$ و را عدم حور هر متغیر پیش ن در مدل k مدل رگرسونی وجود دارد که فرض می شود همگی در فای مدل \mathcal{M} قرار دارند اگر $\{M_1, \dots, M_T\} = \mathcal{M}$ نشان دهنده فای مدل Δ که تی مورد علاقه مانند مشاهدهای در نموده اشد در اینصورت استنباط زی راساس توسع پسند Δ صورت می پندرد که ما فرض مشاهده مجموعه داده D نا ر قاعده احتمال کل مختهای از احتمالهای پسند ن همه مدلها صورت

$$Pr(\Delta|D) = \sum_{k=1}^T Pr(\Delta|M_k, D).Pr(M_k|D) \quad (2)$$

است مانگن و وارانس پسند Δ ترتیب صورت

$$\begin{aligned} E[\Delta|D] &= \sum_{k=1}^T E[\Delta|M_k, D].Pr(M_k|D) \\ &= \sum_{k=1}^T \hat{\Delta}_k.Pr(M_k|D) \end{aligned} \quad (3)$$

و

$$\begin{aligned} Var[\Delta|D] &= E_M[Var(\Delta|D, M)] + Var_M[E(\Delta|D, M)] \\ &= \sum_{k=1}^T (Var[\Delta|D, M_k] + \hat{\Delta}_k^2) \times Pr(M_k|D) - E[\Delta|D]^2 \end{aligned} \quad (4)$$

هستند رفتري که در ن مؤلفه $Var_M[E(\Delta|D, M)]$ عدم حتمت ن مدلها را نشان میدهد عبارت $\hat{\Delta}_k$ مانگن وزنی احتمال پسند است که در ن هر مدل احتمال

۸ هفتمن کنفرانس هادا دان

پسند متنا ر خود عني $P(M_k|D)$ وزن دار شده است و توز پشن ن Δ ه شر مدل M_k صورت

$$Pr(\Delta|M_k, D) = \int Pr(\Delta|\theta_k, M_k, D).Pr(\theta_k) d\theta_k \quad ()$$

است که در ن θ_k ردار پارامترهای مدل M_k را نشان می دهد احتمال پسند مدل M_k استفاده از قاعده ز صورت

$$Pr(M_k|D) = \frac{Pr(D|M_k).Pr(M_k)}{\sum_{j=1}^T P(D|M_j).Pr(M_j)} \quad ()$$

دست می د که در ن $Pr(M_k)$ احتمال پشن درست ودن مدل M_k و

$$Pr(D|M_k) = \int Pr(D|\theta_k, M_k).Pr(\theta_k|M_k) d\theta_k \quad ()$$

درستنمای جم سته ۲ مدل M_k و $Pr(D|\theta_k, M_k)$ درستنمای مدل M_k است رای محاسبه $Pr(\Delta|D)$ اید مولفه های تشکیل دهنده ن را مشخص و جاگز ن نمود ناران لازم است توز پشن پارامترها احتمالهای پشن و پسند هر مدل و توز پشن که مت مورد علاقه را رای همه مدل های موجود در فای مدل مشخص کرد

الف تع ن احتمال پشن پارامترها مدل کی از مسائل دشوار در BMA تخصص پشن ه پارامترهای مدل است استفاده از پشن های نامناسب ۳ منجر ه توز های پسند نامناسب می شود که در انصورت نمی توان ان احتمالها را ه عنوان احتمال مدل و نسبت نهرا را ه عنوان ز فاکتور تعبیر نمود ه هم دلیل ساری از محققان گونه های مختلفی از پشن های گاهی بخش ۴ را پشن هاد کرده اند هوتونگ استفاده از پشن های مناسبی ۵ را پشن هاد کرده است که در قسمت های از فای پارامتر ا درستنمای از رگ هموار اشند مدل را می توان صورت

$$Y = X\beta + \varepsilon \quad ()$$

نوشت که در ن $X_{n \times (p+1)}$ ماترس مشاهدات Y ردار n عدی متغرهای وابسته و ε ردار خواست فرض می شود ε ها رای مشاهدات مختلف مستقل و دارای توز نرمال ام انگان صفر و وارانس σ^2 هستند ردار $(\beta_0, \dots, \beta_p) = \beta$ و پارامتر σ^2 نامعلوم هستند توز عهای پشن ان پارامترها ایستی گونه ای تع ن شوند که عدم حتمت نهرا ه خوی منعکس سازند هوتونگ رده پشن های مزدوج نرمال گاما را صورت

$$\beta \sim N_{p+1}(\mu, \sigma^2 V), \quad \frac{\nu \lambda}{\sigma^2} \sim \chi^2_\nu \quad ()$$

2) Integrated Likelihood 3) Improper Prior 4) Informative 5) Proper Prior

مجموعه مقالات

در نظر گرفت که در نهایت λ ماتریس $V_{(p+1) \times (p+1)}$ و ردار p عدی μ از پارامترهای مستند که امروز شوند

تعیین احتمال پشن هر مدل وقتی هچ‌الا پشنی در مورد مدلها وجود ندارد اما زان اطلاعات پشن در مورد مدلها اندک است معمولاً فرض می‌شود توزیع مدلها گنواخت است یعنی $Pr(M_j) = \frac{1}{T}$ صورت

$$Pr(M_k|D) = \frac{Pr(D|M_k)}{\sum_{j=1}^T Pr(D|M_j)} \quad (1)$$

خواهد ود اگر μ از مدلها در مقامه اسارت دارای احتمال شتری باشد اما اطلاعات خوی در مورد نهایا در دسترس نداشت لازم است این اطلاعات رای تعدل احتمالهای پشن مدلها بکار گرفته شوند تا از پشن های گاهی خشن تر استفاده شود در مسائل مربوطه انتخاب متغیرهای پشن اطلاعات پشن شکل شواهد قبلي رای در نظر گرفتن که متغیر مورد استفاده قرار می‌گیرد فرض کند تنها این تو اطلاعات پشن در اختار باشد و مدل M_k توزیع ردار $(\delta_{kp}, \dots, \delta_{k1})$ مشخص شود که در ن

$$\delta_{ki} = \begin{cases} X_i \in M_i \\ X_i \notin M_i \end{cases} = 1, \dots, p$$

توان نشانگر مستند حال اگر π_i احتمال موثر ودن متغیر X_i را نشان دهد و پذیرم که اطلاعات پشن در مورد متغیرهای متفاوت تقریباً مستقل هستند می‌توان

$$Pr(M_k) = \prod_{i=1}^p [\pi_i^{\delta_{ki}} \times (1 - \pi_i)^{1 - \delta_{ki}}] \quad (2)$$

را عنوان احتمال پشن صحیح ودن مدل M_k در نظر گرفت چون این توزیع متغیری که مهمتر است احتمال زرگری تخصیص می‌دهد توزیع مادگان و فترتی توزیع پشن متغیر نامده شده است

ج) تعیین احتمال پشن هر مدل رای محاسبه احتمال پسن هر مدل لازم است $Pr(D|M_k)$ از رای که انتگرال اعداد را با تعداد پارامترهای مدل M_k است محاسبه شود نهاران محاسبه دقیق این احتمال دلیل پژوهش ودن انتگرال مربوطه تنها در حالات سهار خاص و ساده امکان پذیر است و در سار موارد از روش‌های تقریبی و محاسباتی استفاده می‌شود در اینجا از معماری BIC رای ان من ور هرمه می‌رم BIC معمار رگرسیون خودی صورت

$$BIC_j = n \log(-R_j^2) + k_j \log n \quad (3)$$

6) Variable Prior 7) Bayesian Information Criteria

..... هفتمن کنفرانس هادا دان

تعزیز می شود که در ن n تعداد مشاهدات R_j^2 رفع ن تعديل شده مدل و k_j تعداد متغیرهای پوشش ن موجود در مدل زام و را نشان می دهد در انصورت درستنمای جمه سنته مدل زام را می توان صورت تقریبی

$$Pr(D|M_j) \propto e^{-\frac{1}{5}BIC_j} = e^{-\frac{1}{5}(n \log(1 - r_j^2) + k_j \log n)} \quad (1)$$

خواهد و د ران اساس احتمال پس ن مدل k استفاده از قاعده ز صورت

$$Pr(M_k|D) = \frac{\exp\{-\frac{1}{5}BIC_k\} \cdot Pr(M_k)}{\sum_{j=1}^T \exp\{-\frac{1}{5}BIC_j\} \cdot Pr(M_j)} \quad (2)$$

نوشت چون معیار BIC رای ساری از مدلها دارای شکل سنته و ساده ای است استفاده از آن تقریب موج سهولت محاسبه و افزایش سرعت می شود د تغییر توزیع پوشش ن انتگرال تشکیل دهنده توزیع پوشش ن از دو جزء تشکیل شده است معمولا رای جز دوم ن از تقریب

$$Pr(\Delta|M_k, D) \approx Pr(\Delta|M_k, \hat{\theta}_k, D) \quad (3)$$

استفاده می شود رفتري و همكاران که در ن $\hat{\theta}_k$ رورد حداکثر درستنمای ردار پارامترهای مدل M_k است اکنون با فرض مشخص و دن همه مولفه های لازم رای اجرای روش BMA محاسبه مجمو واسه ه تعداد زاد جملات ن عمل امکان پذربنست و لازم است زر مجموعه ای از محتمل تر ن مدلها انتخاب شود مادگان و رفتري روش پنجه اه اوکام⁸ را رای ان منه و پشنهداد کرده اند که از دو اصل کلی پروی می کنند ناراصل اول مدلهاي که در مقامه ما محتمل تر ن مدل خلی کم شناسی هستند کنار گذاشته می شوند ناراصل دوم که تغییر اوکام⁹ نام ده می شود مدلهاي که نسبت ه ز رمدهای ساده تر خود کمتر از جاز داده ها حماست می شوند کنار گذاشته می شوند با اجرای اصل اول مدلهاي که در مجموعه

$$\mathcal{A}' = \left\{ M_k : \frac{\max_{M_l \in \mathcal{M}} \{Pr(M_l|D)\}}{Pr(M_k|D)} > C_1 \right\}$$

قرار دارند و همچنان ناراصل تغییر اوکام مدلهاي که در مجموعه

$$\mathcal{B} = \left\{ M_k : \exists M_l \in \mathcal{A} \text{ s.t } M_l \subset M_k, \frac{Pr(M_l|D)}{Pr(M_k|D)} > C_2 \right\}$$

8) Occam's Window 9) (OW)Occam's Razor

مجموّعه مقالات

اشتند از مجموع خارج می‌شوند که در ن C_1 و C_2 توسه تحلیل‌گرته ن می‌شوند در انصورت مجموع را می‌توان صورت

$$Pr(\Delta|D) = \frac{\sum_{M_k \in \mathcal{A}} Pr(\Delta|M_k, D).Pr(D|M_k).Pr(M_k)}{\sum_{M_k \in A} Pr(D|M_k).Pr(M_k)} \quad ()$$

از نویسی کرد که در ن $\mathcal{A} = \mathcal{A}' - \mathcal{B} \in \mathcal{M}$ مجموعه پذرش است

مثال کار رد

در این بخش کارای روش BMA در شخوص متغیرهای پشن ن موثر در مدل و رورد پارامترهای مدل با روش مدلسازی گام به گام عنوان کی از روش‌های مدلسازی مرسوم مورد مقاسه قرار می‌گردید اینکه ورتعی کارای روش BMA افزایش تعداد متغیرهای پشن ن افزایش می‌آید رای نماش هنرکاری روش BMA از مثالی تعداد متغیرهای پشن ن کم استفاده شده تا حتی المکان از الای ودن کارای ان روش دلیل تعداد زاد متغیر پشن ن اجتناب شود

جدول داده‌های مر و کارای شغلی پرستاران

Y	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}	X_{11}	X_{12}	X_{13}	X_{14}
X_1														
X_2														
X_3														
X_4														
X_5														
X_6														
Y	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}	X_{11}	X_{12}	X_{13}	X_{14}
X_1														
X_2														
X_3														
X_4														
X_5														
X_6														

جدول مشاهدات که مالعه رگرسونی را نشان می‌دهد که هدف ن ررسی تأثیر متغیر مستقل شامل قاعع X_1 علاقه‌مندی X_2 اند همچنین X_3 روا اجتماعی X_4 هنرمندی در حل مشکلات X_5 و انکار X_6 رکارای شغلی پرستاران است مدل وجود متغیر پشن ن فای مدل شامل $= 6$ مدل است روش پنجره اوکام مدل را راساس مقاسه احتمالهای پسون مدلهاه عنوان رترن مدلهاه موجود مورد استفاده قرار می‌دهد مدلهاهی حاصل همراه احتمال پسون مقدار R^2_{adj} درصد و معمارا لا ز نها در جدول اراه شده‌اند تعداد زاد مدلهاهی که در پنجره اوکام قرار گرفته‌اند مش از $\frac{1}{3}$ مدلها نشان

..... هفتمن کنفرانس مادا دان

جدول مدل‌های انتخابی توسعه روش پنجره اوکام

BIC	R^2_{adj}	احتمال پسن	مدل	شماره مدل
—	—	—	X_3	X_6
			$X_2 X_3$	X_6
			$X_1 X_2 X_3$	X_6
			$X_1 X_3$	X_6
			X_3	
			$X_3 X_5 X_6$	
			$X_2 X_3 X_5 X_6$	
			$X_3 X_4 X_6$	
			$X_2 X_3 X_4$	
			$X_1 X_3 X_4$	
			$X_1 X_2 X_3 X_4$	
			$X_1 X_2 X_3 X_5 X_6$	
			$X_1 X_2 X_3 X_4 X_6$	
			$X_1 X_3 X_4 X_5 X_6$	
			$X_1 X_3 X_4 X_5$	
			$X_2 X_3 X_4 X_5$	
			$X_3 X_4 X_5 X_6$	
			$X_2 X_3 X_4 X_6$	
			$X_2 X_3 X_5$	
			$X_2 X_3 X_4$	

دهنده زاد و دن عدم حتمت است عینی ش از ک مدل صحیح وجود دارد ناران انتخابی ک مدل ه عنوان مدل نهایی من قی ن رنمی رسد ملاحه می شود که مقدار R^2_{adj} و معار ا لا ز BIC رای مدل‌های که در پنجره اوکام قوارگرفته‌اند نزدیک ه بکدیگر هستند شترن مقدار R^2_{adj} مرو ه مدل شماره است در حالی که راساس مقدار که ه معار ا لا ز مدل هترن است انتخابی ک مدل و مبنا قرار دادن ن ه معنی ن است که احتمال صحیح و دن ن ک می‌اشد در حالی که شترن احتمال پسن مرو ه مدل شماره و حدود است و احتمال پسن سار مدلها همگی از کتر است جدول نشان می‌دهد که متغیر X_3 در همه مدلها و متغیر X_6 در مدل متنخ روش پنجره اوکام ه و دارند که نشان دهنده تأثیرگذار و دن ان دو متغیر است رورد را متغرهای مستقل که در جدول اراه شده نز اهمت دو متغیر X_3 و X_6 را تا د می‌کند همان ورکه در جدول ملاحه می‌شود روردهای ز را ان دو متغیر که ا استفاده از روش BMA دست مدها ند ه خوی اهمت ان دو متغیر را منعکس می‌سازند جدول رورد پارامترها و خای معار نهارا که ه روش گام ه گام دست مدها ند را نشان می دهد همان ورکه ملاحه می شود خای معار متنا را روردهای ای دو پارامتر در روش BMA صورتی قابل ملاحه ه از مقدار مشاه در روش گام ه گام زرگتر هستند ان بکی از رترهای روش BMA است که مزان واقعی عدم حتمت را نشان می دهد ه ل دگری که می‌توان ه ن اشاره کرد ان است

مجموّعه مقالات

۳

جدول رورد را در روش <i>BMA</i>				
پارامتر رگرسیونی	$Pr[\beta_i \neq D]$	مانگن پسند	انحراف معنادل	پسند
β_1				
β_2				
β_3				
β_4				
β_5				
β_6				

جدول رورد را در مدل کامل و روش گام به گام					
خای معنادل	رورد پارامترها	روش گام به گام	خای معنادل	رورد پارامترها	مدل کامل
		*			β_0
		*			β_1
		*			β_2
					β_3
					β_4
					β_5
					β_6

که بخت ح ور متغیر X_2 در مدل ش از است که راساس قواعد سرانگشتی جفرز شواهد مثبتی رای ح ور متغیر X_2 در مدل اراه میکند در حالی که ا لاعات موجود در ان متغیر احذف ن در مدل حاصل از روش گام به گام نادمه گرفته میشوند رورد را متغرهای در روش *BMA* در مقامه سار روشها و خصوصا در مورد متغرهای X_1 و X_4 کاملا کوچک و نزدیک صفر هستند ا ان وجود کاملا نادمه گرفته نمیشوند و ا لاعات نها استفاده میشود

بحث و نتیجهگیر

در هنگام مدلسازی غالبا نمیتوان مدلی اافت که صورت کامل ه دادهها را زش داشته اشد در ان حالات انتخاب ک مدل ه معنی از ن رفتن ا لاعات سار مدلها است مانگنگری زی مدلها ازا لاعات همه مدلها ا دستهای از هترن مدلها استفاده میکند و رخلاف روشهای مرسوم عدم حتمت را ه خوی منعکس میکند از نتیجه ن رکارای پشن ن نز نتیجه حاصل از روش *BMA* از نتایج حاصل از هر ک از مدلها موجود در فای مدل ملو تراست

..... هفتمین کنفرانس هادا دان

عملکرد این روش افزایش عدم حتمت هبود می‌آید مثلاً هر چه تعداد متغیرهای پوش ن در مدل رگرسونی شتر اشد کارای این روش شتر می‌شود در مدترين حالت که عدم حتمت اندک اشد هتر ودن روش BMA نسبت به سار روشها چشمگر نمی‌اشد و کارگری روشها مرسوم دلیل سادگی مقرنون به صرفه‌تر است

مراج

- [1] Good, I. J., (1950). "Probability and The Weighing of Evidence", Griffin London.
- [2] Hoeting, J., (1994). "Accounting for Uncertainty in Linear Regression Models", Ph.D Dissertation, Department of Statistics, University of Washington.
- [3] Hoeting, J., Raftery, A. and Madigan, D., (1996). "A Method for Simultaneous Variable Selection and Outlier Identification in Linear Regression Models". *J. Comput. Statist.* **22**, 251-271.
- [4] Hoeting, J., Raftery, A. and Madigan, D., (1999). "Bayesian Simultaneous Variable Selection and Transformation Selection in Linear Regression Models". Technical Report 9905, Dept. Statistics, Colorado State Univ. Available at www. Colostate. edu.
- [5] Lipkovich, I. A., (2002). "Bayesian Model Averaging and Variable Selection in Multivariate Ecological Models", Ph.D Dissertation, Faculty of The Virginia Polytechnic Institute and State University, Blacksburg, Virginia.
- [6] Madigan, D., Gavrin, J. and Raftery, A. E., (1995). "Eliciting Prior Information To Enhance The Performance of Bayesian Graphical Models", *Comm. Statist. Theory Methods*, **24**, 2271-2292.
- [7] Madigan, D and Raftery, A., (1991). "Model Selection and Accounting for Model Uncertainty in Graphical Models Using Occam's Window", Technical Reports 213, Univ. Washington, Seattle.
- [8] Madigan, D. and Raftery, A., (1994). "Model Selection and Accounting for Model Uncertainty in Graphical Models Using Occam's Window", *Journal of American Statistical Association*, **89**, 1535-1546.
- [9] Miller A., (1990). "Subset Selection in Regression Variables", New York, Chpman-hal.
- [10] Raftery, A. E., (1995). "Bayesian Model Selection in Social Research(With Discussion)", in *Sociological Methodology 1995*(P. V. Marsden, ed.) 111-195. Blakwell. Cambridge, MA.

..... مجموعه مقالات

- [11] Raftery, A. E., (1996). "Approximate Bayes Factor and Accounting for Model Uncertainty in Generalized Linear Models", Technical Report, Biometrika, **83**, 351-266.
- [12] Robert, B. Nobel, J. (2000). " Multivariate applications of Bayesian Model averaging", Ph.D dissertation, Faculty of the virginia polytechnic institue and state university, Blacksburg, virginia.

Archive of SID