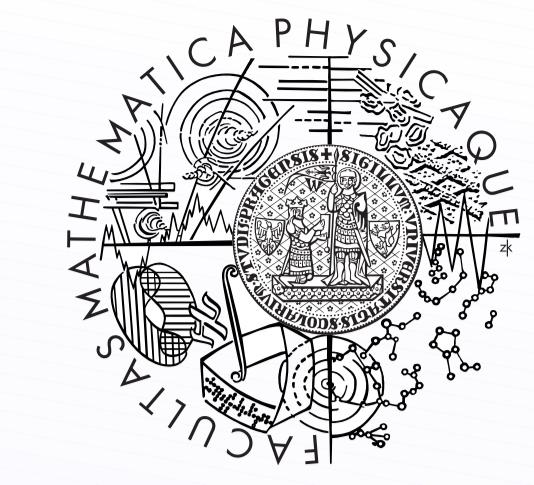




Robustifikace hřebenové regrese - první pokusy

Tomáš Jurczyk

Matematicko-fyzikální fakulta, UNIVERZITA KARLOVA V PRAZE, Katedra pravděpodobnosti a matematické statistiky,
Sokolovská 83, 18675 Praha 8, ČR

1 Abstrakt

Jeden z problémů, se kterými se můžeme při datové analýze setkat, je multikolinearita (též lineární závislost mezi regresory). Multikolinearita vede k nestabilitě řešení normální rovnice a může způsobit velký rozptyl odhadu regresních koeficientů. Jednou z metod, která se snaží tento problém řešit je hřebenová regrese. Tato metoda ale není imunní vůči přítomnosti kontaminace, která může multikolinearitu skrýt nebo ji naopak uměle vytvořit. Tento příspěvek se věnuje prvním jednoduchým pokusům o robustifikaci metody hřebenové regrese. V práci jsou představeny dva návrhy založené na robustní metodě nejmenších vážených čtverců (LWS).

3 Hřebenová regrese

Jak řešit problém s multikolinearitou?

Místo klasického odhadu se v případě podezření na multikolinearitu raději použije vychýlený odhad:

Definice hřebenové regrese

Pro $\delta \geq 0$ definujeme **hřebenový odhad** vektoru β jako

$$b_\delta = (\mathbf{X}'\mathbf{X} + \delta \mathbf{I})^{-1}\mathbf{X}'\mathbf{y}.$$

- Tento odhad má při malých hodnotách δ (přesněji $0 < \delta < 2\sigma^2 \|\beta\|^{-2}$) menší čtvercovou chybu než klasický odhad b .

Volba δ

Možnosti, jak volit δ :

- a) Odhadne se horní mez (viz předešlá vlastnost): $\delta = 2s^2 \|\beta\|^{-2}$, kde s^2 je odhad σ^2 . Někdy se také doporučuje $\delta = ks^2 \|\beta\|^{-2}$.

- b) Grafickou metodou zvanou **ridge trace**, kde se vynášeji do grafu hřebenové odhadu koeficientů proti proměnné δ . Za vhodnou volbu δ se pak považuje podle autorů metody taková hodnota, při níž jsou již znaménka a hodnoty složek b_δ stabilizovány.

Problémy

Hřebenová regrese není díky své konstrukci uzpůsobena pro práci s odlehlymi pozorováními. Navíc odlehlá pozorování mohou multikolinearitu skrýt nebo ji naopak uměle vytvořit – už i jedno přidané odlehle pozorování může výrazně změnit, jak výběrový korelační koeficient, tak například i výše zmíněné indexy podmíněnosti.

Detekce multikolinearity v kontaminovaných datech je velmi důležitá, protože odhaluje pravé závislosti v datech.

Reference

[Rousseeuw and Leroy (1987)] Rousseeuw, P. J., and Leroy, A. M., Robust Regression and Outlier Detection, John Wiley & Sons, 1987.

[Víšek (2001)] Víšek, J. Á., Regression with high breakdown point, Robust 2000 (eds. Jaromír Antoch & Gejza Dohnal, published by Union of Czech Mathematicians and Physicists), Prague: matfyzpress, 324-356, 2001.

[Víšek (2008)] Víšek, J. Á., Consistency of the Instrumental Weighted Variables, To appear in Annals of the Institute of Statistical Mathematics, 2008.

[Zvára (1989)] Zvára, K., Regresní analýza, Academia, Praha, 1989.

5 Robustifikace hřebenové regrese

V této části budou představeny dva první jednoduché návrhy na robustifikaci metody hřebenové regrese za použití nejmenších vážených čtverců. Také bude shrnuto, jak si tyto metody vedly v případech, kdy byla multikolinearita kontaminací zakryta.

Hřebenová regrese a WLS

Nejdříve je potřeba si uvědomit, že volba vah u WLS neřeší problém se závislostí mezi regresory. Jednotlivá pozorování jsou převážena, to ale nepotačí závislosti mezi jednotlivými složkami pozorování. Proto pokud je problém s multikolinearitou u regresní matice \mathbf{X} , bude tento problém i u $\mathbf{W}^{\frac{1}{2}}\mathbf{X}$. Hřebenovou regresi pro WLS má tedy smysl definovat jako $b_{\delta}^{WLS,w} = (\mathbf{X}'\mathbf{W}\mathbf{X} + \delta \mathbf{I})^{-1}\mathbf{X}'\mathbf{W}\mathbf{y}$.

Hřebenová regrese a LWS (robustní hřebenová regrese)

Popis metody

Dvoukrokový odhad:

1.krok: Nejdříve se spočítá LWS

2.krok: V prvním kroku získané $\mathbf{W}(\pi)$ využijeme pro výpočet klasické hřebenové regrese na data $\mathbf{W}(\pi)^{\frac{1}{2}}\mathbf{y}, \mathbf{W}(\pi)^{\frac{1}{2}}\mathbf{X}$. Tedy

$$b_{\delta}^{LWS,w} = (\mathbf{X}'\mathbf{W}(\pi)\mathbf{X} + \delta \mathbf{I})^{-1}\mathbf{X}'\mathbf{W}(\pi)\mathbf{y}.$$

Idea aneb proč by to mohlo fungovat

Máme data, kde odlehlá pozorování zakrývají problém s multikolinearitou.

Indexy podmíněnosti toto neodhalí, také ridge trace vypadá odlišně než u případu s téměř závislými sloupcemi.

Pokud použijeme LWS s vahami, kde $w_i = 0$ pro $i > h$ pro nějaké h a tato metoda odhalí všechna odlehlá pozorování, pak na datech $\mathbf{W}(\pi)^{\frac{1}{2}}\mathbf{y}, \mathbf{W}(\pi)^{\frac{1}{2}}\mathbf{X}$ bude již kolinearita patrná a na tato nová data se pak použije hřebenová regrese.

Poznámka: Pokud by toto fungovalo, pak by se v takových případech dala LWS výhodně použít jako robustní detektor multikolinearity.

Výsledky

Ačkoli se zdají být tyto úvahy správné, na simulovaných a také na datech z knížky [Zvára (1989)] (kapitola 9.) se ukázalo, že takto navržená metoda nefunguje, jak bychom čekali.

Za přítomnosti odlehly pozorování zakrývajících multikolinearitu LWS neodhalí správně tato odlehlá pozorování, ačkoli je na to „stavěná“. Navíc

dává přiřazením vah vždy alespoň jednomu z odlehlymi velmi vysokou důležitost.

Interpretace tohoto jevu

Jsem přesvědčen, že tento problém je způsoben tím, že v datech je vlastně jakási volnost způsobená téměř závislostí. V datech jsou v zásadě dvě skupiny regresorů, které se snaží vysvětlit to samé.

Při přidání odlehly pozorování k datům se potom může jedna z těchto skupin soustředit výhradně na vyrovnaní jen tohoto pozorování – přitom může být totiž vyrovnaný velmi přesně. Což vede k tomu, že LWS vybere raději tento model s přesně vyrovnaným odlehlym pozorováním místo modelu na skutečných datech, která jsou ovšem nestabilní.

Závěr

► Vlivem skryté multikolinearity nedokáže navržená metoda ve svém prvním kroku úspěšně odhalit všechna odlehlá pozorování, a proto není vhodným kandidátem pro robustní verzi hřebenové regrese.

► V knize [Rousseeuw and Leroy (1987)] doporučují autoři řešit možný problém skryté multikolinearity tím, že také nejdříve očistí data od odlehlych pozorování metodou nejmenšího mediánu čtverců a poté provedou hřebenovou regresi.

Vzhledem k předešlým úvahám a provedeným analýzám se obávám, že tato metoda také nemůže fungovat správně.

Návrh 1

Popis metody

Pro pevný δ se do vnitřního cyklu LWS implementuje iterativně

$$\mathbf{b}_{(k+1)} = (\mathbf{X}'\mathbf{W}(\pi_{(k)})\mathbf{X} + \delta \mathbf{I})^{-1}\mathbf{X}'\mathbf{W}(\pi_{(k)})\mathbf{y}$$

místo

$$\mathbf{b}_{(k+1)} = (\mathbf{X}'\mathbf{W}(\pi_{(k)})\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}(\pi_{(k)})\mathbf{y}.$$

Vnější cyklus pak vyhodnocuje nejlepší model podle klasického kritéria LWS.

Nevýhody

Pro zkonstruované ridge trace se provede takto upravená metoda pro nějakou míru hodnoty δ . Nic ale nezaručuje, že s rostoucím δ zůstane výsledná ideální permutace vah stále stejná, což potvrzily i provedené výpočty. V ridge trace se tedy mohou vyskytovat skoky.

Tím se komplikuje výběr vhodného δ . Navíc je také otázka, co vůbec znamenají návrhy a) na volbu δ v části o hřebenové regresi, protože odhad z celých dat za přítomnosti kontaminace může být touto kontaminací výrazně ovlivněn, nehledě na další problémy s tím spojené.

Výhody

S rostoucím δ se narodí od návrhu 1 této metody dařilo na zkoumaných datových souborech správně odhalit všechna odlehlá pozorování zakrývající multikolinearitu.

Závěr

- Možné problémy s volbou δ .
- S rostoucím δ se metoda daří správně identifikovat všechna odlehlá pozorování zakrývající multikolinearitu!

6 Plány do budoucna

V tomto posteru byly prezentovány zatím jen úplně základní návrhy, je možné hledat další metody, ať už vycházejí z hřebenové regresy či například z regrese s lineárními omezeními, která se také používá při výskytu kolinearity.

První výsledky prezentovaného návrhu 2 se zdají být velmi slibné, proto je potřeba se touto metodou zabývat podrobněji.

Také zůstává nezodpovězeno mnoho dalších otázek, například v práci nebyl diskutován problém dat, kde je multikolinearita uměle vytvořena odlehlymi pozorováními, apod.