

رگرسیون ژرفا در حالت چندگانه

حمید شریف

بانک مرکزی جمهوری اسلامی ایران

چکیده در این مقاله به آن اشاره ژرفا در رگرسیون می‌پردازیم. رگرسیون ژرفا^۱ که عدد صحیح n و m است به عنوان خاصی از رازش رگرسیون در n گرفته می‌شود که رتبه رازش رگرسیون تعبیر می‌شود و معاری رای مقایسه رازشهای مختلف رگرسیون می‌اشد. در این مقاله ما آن رگرسیون ژرفا در حالت چندگانه به بررسی خواص n پرداخته و روش ژرفتر n رگرسیون را که نسبت به تبدلات کنوا n متغیر پاسخ هم‌ورد است ارائه می‌دهیم. ما معرفی معیار مقدار فرورزش نشان می‌دهیم که آن روش روشی استوار می‌اشد. در پایان ورودکننده ژرفتر n رگرسیون را با دیگر ورودکننده‌های رگرسیون و محاسبه رگرسیون ژرفای آنها رای داده‌های نمونه‌ای موجود مقایسه می‌کنیم.

واژه‌ها کلید واژه‌ها: ژرفا نارازا^۲ ژرفتر n رگرسیون مقدار فرورزش رگرسیون ژرفای ماکس حال هم‌وردی

مقدمه

وقتی رگرسیون کمترین توانهای دوم را استفاده از m مشاهده رای یک مدل پارامتری $Y = X\beta + \varepsilon$ به کار می‌ریم فرم‌های را در مورد ردار ε در n می‌گیریم. یکی از این فرم‌ها نرمال بودن توزیع ε است. در عمل انحرافات ε از این فرم‌ها رخ می‌دهد. مثلاً ممکن است توزیع ε نرمالی ε متقارن اما غیر نرمال باشد. ε نیز از نرمال ε و دمهای کوتاهتری داشته باشد. کشیدگی کمتر از نرمال ε دمهای پهن‌تر داشته باشد و ε ممکن است توزیع ε صورت نرمال باشد ولی دارای دورافتاده‌های باشد.

در این موارد به جای روش کمترین توانهای دوم از روشهای رگرسیون استوار استفاده می‌شود که در مقایسه با روش کمترین توانهای دوم نسبت به این انحرافات حساسیت کمتری دارند. یک روش رگرسیون استوار که اخرا توسط روسف و هورت^۳ مطرح شده است روش ژرفتر n رگرسیون نام دارد که رپا به n رگرسیون ژرفا نام شده است. رگرسیون ژرفا کمترین رازش را اندازه‌گیری کرده و m زمان دوری n را از هر نارازا اندازه می‌گیرد و آن می‌کند که ارفصفحه رازشی در توصیف داده‌ها به چه اندازه خوب عمل می‌کند. نارازان رازشی n ژرفای بزرگ نسبت به

1) regression depth 2) nonfit 3) Rousseeuw and Hubert

داده‌ها متعادلتر است و لذا یک رازش خو ژرفای زگرتی نسبت به یک رازش د دارد در این مقاله مرور رگرسون ژرفا در حالت چندگانه به بررسی خواص ن پرداخته و روش ژرفترین رگرسون را بیان می‌کند.

در بخش دوم این مقاله ورمختصر رگرسون ژرفا در حالت ساده مرور می‌کند. در بخش سوم رگرسون ژرفا در حالت چندگانه معرفی می‌کند. رگرسون ژرفای ماکسیمال و ژرفترین رگرسون به ترتیب در بخشهای چهارم و پنجم بیان می‌شود و در بخش ششم خواص ژرفترین رگرسون در حالت چندگانه بررسی می‌گردد. این خواص شامل هم‌وردایی و استواری ژرفترین رگرسون است و در انتها با ارائه مثال ورودکننده ژرفترین رگرسون را با دیگر ورودکننده‌های رگرسونی مقایسه می‌کند.

مرور رگرسون ژرفا در حالت ساده

هدف در رگرسون ساده رازش یک خط راست $y = \theta_1 x + \theta_2$ است که مجموعه داده $Z_n = \{(x_i, y_i); i = 1, \dots, n\} \subseteq R^2$ است. یک رازش را به صورت $\theta = (\theta_1, \theta_2)$ نشان می‌دهیم که مولفه اول ن ورودش و مولفه دوم ن جمله عرض از مبدا است. مانده‌های مجموعه داده Z_n متناسب با رازش θ را به صورت $r_i(\theta) = y_i - \theta_1 x_i - \theta_2$ نشان می‌دهیم. رای معرفی ژرفای یک رازش ابتدا یک نارازا را تعریف می‌کنیم.

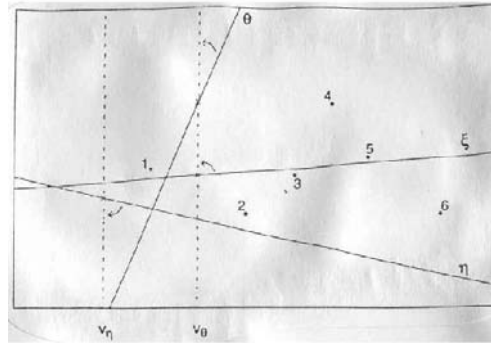
تعریف ۱ رازش $\theta = (\theta_1, \theta_2)$ رای مجموعه داده Z_n نارازا نامیده می‌شود اگر و فقط اگر یک عدد حقیقی $v = v_\theta$ مخالف همه x_i ها وجود داشته باشد و برای که

$$\forall x_i < v, r_i(\theta) < 0, \quad \forall x_i > v, r_i(\theta) > 0$$

$$\forall x_i < v, r_i(\theta) > 0, \quad \forall x_i > v, r_i(\theta) < 0$$

ذی که الا با ن همه مشاهدات واقع می‌شود هم‌شبه یک نارازا است. عدا خواهیم دید که نارازاها در واقع رازش‌های ژرفای صفر هستند.

شکل یک مجموعه داده را با مشاهده θ و نارازا نشان می‌دهد. مقدار v_θ و v_η نشان داده شده‌اند. با توجه به شکل وجود v متناظر با وجود یک برادر روی خط عمودی که از نقطه v رسم شده است. چون $v_\theta < v_\eta$ و همچنین $r_i(\theta) > 0$ و $r_i(\eta) < 0$ لذا با توجه به تعریف فوق θ نارازا است. همچنین چون $v_\theta < v_\eta$ در پان همه مشاهدات قرار دارد این خط نیز نارازا است. از ن جاکه تعریف فوق برای خط ξ صادق نیست این خط نارازا نیست. همان‌ور که ذکر شد مشخص شدن v نقطه متناظر با ن که محل تلاقی خط عمودی ترسیم شده از نقطه v با رگرسونی است و این رگرسون نشان داده شده است مشخص می‌گردد. به عبارت دیگر محل تلاقی خط رازش داده شده با خط عمودی که از نقطه v متناظر با ن خط ترسیم می‌گردد برادر مشخص شده است. از این



شکل مجموعه داده دو متغیره با نازاهای θ و η و ک رازش ξ ا رگرس ون ژرفای

رو می توان خ و نازازا را از روی شکل مشخص کرد به این ترة. که خ رازش داده شده را حول نقطه ای که ا بر روی ن خ مشخص شده است می گردانیم تا عمودی شود اگر در حین چرخاندن خ از هیچ مشاهده ای عبور نکند نازازا است بنابراین خ و θ و η نازازا هستند ولی خ ξ نازازا نیست چون وقتی ن را حول نقطه مذکور مشخص شده ا ردر می گردانیم تا عمودی شود از مشاهدات و عبور می کند به ورکلی ژرفای ک رازش θ رای ک مجموعه داده Z_n به حجم n صورت ز ر ان می شود

تعرف رگرس ون ژرفای ($rdepth$) ک رازش θ متناسب با مجموعه داده Z_n عبارت است از کمترین تعداد مشاهداتی که ا داشته شود تا θ نازازا شود. و معادل $rdepth(\theta, Z_n)$ عبارت است از کوچکترین تعداد مانده های که می است علامتشان تغییر کند تا θ نازازا شود

رای مثال خ ξ را در شکل در نر می گیریم این خ ا حذف کردن مشاهدات و نازازا می شود ز را ا قرار دادن v_ξ مساوی v_θ و حذف مشاهدات و می توان خ ξ را به صورت عمودی درورد دون ن که از مشاهده ای عبور کند چون خ ξ ا حذف حداقل مشاهده نازازا می شود لذا $rdepth(\xi, Z_n) =$

توجه ا تعارف و در مواقعی که در x_i ها تکرار وجود داشته اشد x_i ها ا هم مساوی اشدن ز رقرار است و x_i ها هیچ فرض توزعی ناز ندارند

مرتبه زمانی گو م $f(n)$ از مرتبه زمانی $g(n)$ است و ن را ا نماد $f(n) = O(g(n))$ نشان می ده م اگر و تنها اگر اعداد صحیح و مثبتی مانند n_0 و c وجود داشته اشد و ر که $f(n) \leq c g(n)$ رای تمام n های $n \geq n_0$ رقرار اشد

رای محاسبه $rdepth(\theta, Z_n)$ ابتدا مشاهدات را به صورت $x_1 \leq x_2 \leq \dots \leq x_n$ در

مرتبۀ زمانی $O(n \log n)$ مرتۀ می‌کنم تعرف می‌کنم

$$L^+(v) = \#\{j; x_j \leq v, r_j \geq \}$$

و

$$R^-(v) = \#\{j; x_j > v, r_j \leq \}$$

L^+ و R^+ و L^- و R^- و $L^+(x_i)$ و $R^-(x_i)$ و $L^-(x_i)$ و $R^+(x_i)$ سپس می‌شوند. برای هر $i = 1, \dots, n$ ، مشاهده تعارف فوق محاسبه می‌شوند $rdepth(\theta, Z_n)$ در $O(n)$ عمل و صورت زر محاسبه می‌شود

$$rdepth(\theta, Z_n) = \min_{1 \leq i \leq n} (\min\{L^+(x_i) + R^-(x_i), R^+(x_i) + L^-(x_i)\}) \quad ()$$

مثال زر نحوه محاسبه رگرسون ژرفای یک خ را با استفاده از ξ نشان می‌دهد
مثال ۱ شکل مشاهده ξ را نشان می‌دهد همان ور که داده می‌شود آن خ
 نارازا نیست با استفاده از ξ رگرسون ژرفای آن خ صورت زر محاسبه می‌شود

$$i = 1: \min\{L^+(x_1) + R^-(x_1), R^+(x_1) + L^-(x_1)\} = \min\{1 + 3, 3 + 1\} = 4$$

$$i = 2: \min\{L^+(x_2) + R^-(x_2), R^+(x_2) + L^-(x_2)\} = \min\{1 + 2, 3 + 2\} = 5$$

$$i = 3: \min\{L^+(x_3) + R^-(x_3), R^+(x_3) + L^-(x_3)\} = \min\{2 + 2, 2 + 2\} = 4$$

$$i = 4: \min\{L^+(x_4) + R^-(x_4), R^+(x_4) + L^-(x_4)\} = \min\{2 + 1, 2 + 3\} = 3$$

$$i = 5: \min\{L^+(x_5) + R^-(x_5), R^+(x_5) + L^-(x_5)\} = \min\{3 + 1, 1 + 3\} = 4$$

$$i = 6: \min\{L^+(x_6) + R^-(x_6), R^+(x_6) + L^-(x_6)\} = \min\{3 + 0, 1 + 4\} = 3$$

$$i = 7: \min\{L^+(x_7) + R^-(x_7), R^+(x_7) + L^-(x_7)\} = \min\{4 + 0, 0 + 4\} = 4$$

بنابراین

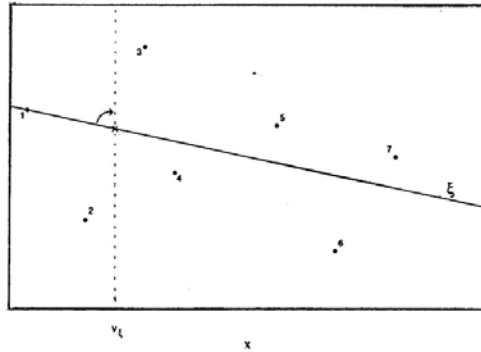
$$rdepth(\xi, Z_n) = \min(4, 5, 4, 3, 4, 3, 4) = 3$$

از روی شکل نز داده می‌شود که برداشتن مشاهدات و می‌توان خ ξ را بدون
 این که از مشاهدات دیگر عبور کند حول نقطه متنا بر v_ξ که بر روی ن خ مشخص
 شده است دوران داد تا به صورت عمودی در بد به عبارت دیگر برداشتن مشاهدات
 و خ ξ نارازا می‌شود رگرسون ژرفای خ ξ برابر با کران پایین رگرسون ژرفای
 ماکس مال یعنی $\lceil \frac{v}{p} \rceil = \lceil \frac{v}{p} \rceil$ است که در بخش‌های بعد آن می‌شود
 توجه $\lceil \lambda \rceil$ کوچکترین عدد صحیح بزرگتر مساوی λ است

رگرسون ژرفا در حالت چندگانه

در رگرسون چندگانه مجموعه داده Z_n و صورت

$$Z_n = \{(x_{i1}, \dots, x_{i,p-1}, y_i); i = 1, \dots, n\} \subset \mathbb{R}^p$$



شکل مجموعه داده دو متغیره و رازش ξ با رگرسیون ژرفای

است X را به عنوان قسمتی از مختصات هر نقطه a به صورت
 $X_i = (x_{i1}, \dots, x_{i,p-1}) \in \mathbb{R}^{p-1}$

در نرمی گرام اکنون می خواهیم y_i را به وسایله

$$\theta_1 x_{i1} + \dots + \theta_{p-1} x_{i,p-1} + \theta_p = (X_i, \theta)$$

معنی به وسایله یک ارفصفحه فن در \mathbb{R}^p که $\theta = (\theta_1, \dots, \theta_p) \in \mathbb{R}^p$ است رازش ده م
 در اینجا هیچ فرض توزعی وجود ندارد

قبل از آن تعارف و مایا ان بخش ابتدا ارفصفحه و ارفصفحه فن را تعرف می کنیم

تعرف ۳ فرض کنید $[f = \alpha]$ نشان دهنده مجموعه $\{x; f(x) = \alpha\}$ باشد یک
 ارفصفحه مجموعه ای به شکل $[f = \alpha]$ است که f تابعی غرضی غرضی روی فای برداری
 X و α عددی حقیقی است به عبارت دیگر مجموعه $H \subseteq \mathbb{R}^n$ یک ارفصفحه است اگر اعداد
 حقیقی a_1, \dots, a_n و C را $a_i \neq 0$ برای حداقل یک i وجود داشته باشند به طوری که H
 متشکل از همه نقاط $x = (x_1, \dots, x_n)$ ای باشد که در آن $\sum_i a_i x_i = C$ صدق می کنند

تعرف یک ارفصفحه فن در فای برداری X مجموعه ای مانند M است به طوری
 که برای هر x_0 در M یک ارفصفحه باشد به عبارت دیگر $M \subseteq X$ یک
 ارفصفحه فن است اگر یک تابعی غرضی $f: X \rightarrow \mathbb{R}$ وجود داشته
 باشد به طوری که $M = \{x \in X; f(x) = \alpha\}$ باشد \mathbb{R} میدان حقیقی \mathbb{C} است

مثال فای برداری $X = \mathbb{R}^2$ و تا $f(x, y) = y$ را در نرمی گرام در آن صورت
 مجموعه $[f = \alpha] = \{(x, y); f(x, y) = \alpha\} = \{(x, \alpha); x \in \mathbb{R}\}$

فای \mathbb{R}^2 را به دو نیم فای دسته $\{(x, y); f(x, y) \geq\}$ و $\{(x, y); f(x, y) \leq\}$ افزایش می‌کند مجموعه

$$C = \{(x, y); f(x, y) =\} = \{(x, y); x \in \mathbb{R}\}$$

ک ا ر صفحه فن است زیرا ازای هر $(x_0, y_0) \in C$

$$C - (x_0, y_0) = \{(x - x_0, y - y_0); x \in \mathbb{R}\} = \{(x - x_0, y - y_0); f(x - x_0, y - y_0) =\}$$

ک ا ر صفحه است

تعریف رازش $\theta = (\theta_1, \dots, \theta_p)$ یک نارازا نامیده می‌شود اگر وفقه اگر ک ا ر صفحه فن V در فای X وجود داشته باشد و وری که هیچ یک از x_i ها متعلق به V نباشند و رای همه x_i ها در یکی از نیم فاهای از ن داشته باشیم $r_i(\theta) >$ و همچنین رای همه x_i ها در نیم فای از دیگر $r_i(\theta) <$ باشد

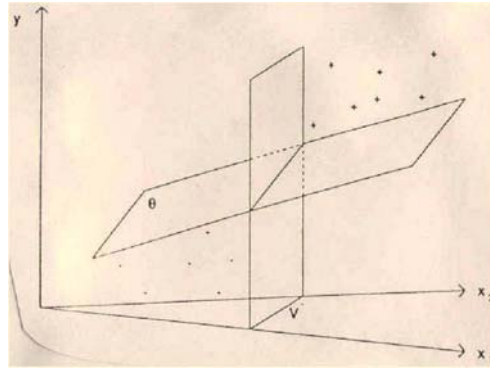
شکل مثالی از یک صفحه نارازاست در ان مثال $p =$ است ناران θ که صفحه متعلق است فای X و صورت صفحه افقی است $y \equiv$ که شامل V است داده می‌شود با ان تعریف هر η وری که همه مانده‌ها مثبت یا همه مانده‌ها منفی باشند یک نارازاست زیرا کافی است V و صورتی انتخاب شود که همه x_i ها در یک رف ن واقع شوند رگرس ون ژرفای هر $\theta \in \mathbb{R}^p$ متناسب است $Z_n \subset \mathbb{R}^p$ همانند تعریف و صورت ز ر تعریف می‌شود

تعریف رگرس ون ژرفای یک رازش $\theta \in \mathbb{R}^p$ متناسب است مجموعه داده $Z_n \subset \mathbb{R}^p$ عبارت است از کوچکترین تعداد مشاهداتی که وقتی از مجموعه داده‌ها خارج شود θ نارازا گردد و عبارت دیگر $rdepth(\theta, Z_n)$ عبارت است از کمترین تعداد مانده‌های که با دستگیری علامت داده شوند تا θ نارازا شود

$rdepth(\theta, Z_n)$ همانند را به ان می‌شود بدین صورت که هر ا ر صفحه فن V در فای X مشاهدات را به دو مجموعه تقسیم می‌کند که ن‌ها را $L(V)$ و $R(V)$ نشان می‌دهد مجموعه‌های $L^+(V)$ و $L^-(V)$ ترتیباً تعداد مشاهدات در $L(V)$ مانده‌های مثبت و تعداد مشاهدات در $L(V)$ مانده‌های منفی است $R(V)$ نیز و مشاهدات در مجموعه $R^+(V)$ و $R^-(V)$ تقسیم می‌گردد ناران

$$rdepth(\theta, Z_n) = \min_V \{\min\{L^+(V) + R^-(V), R^+(V) + L^-(V)\}\}$$

که V شامل همه ا ر صفحه‌های فن در فای X است رای محاسبه رگرس ون ژرفای یک رازش جستجو را به مجموعه‌ای متناهی از ا ر صفحه‌های V محدود می‌کنیم الگوریتم‌های مناسب رای محاسبه رگرس ون ژرفا در فصل بعد ارائه می‌شود لازم به ذکر است که حالت تعمیم یافته قبلی یعنی در حالت چندگانه p بعدی زمانی که $Z_n = \{(X_i, y_i); i = 1, \dots, n\}$ که $X_i = (x_{i1}, \dots, x_{i,p-1}) \in \mathbb{R}^{p-1}$ و



شکل مثالی از ناراها $\theta \in \mathbb{R}^3$ ارضفحه فن V در فای $X (\mathbb{R}^2)$ مشاهدات ا مانده‌های مثبت را از مشاهدات ا مانده‌های منفی جدا می‌کند

$\theta = (\theta_1, \dots, \theta_p) \in \mathbb{R}^p$ اشد نیز رقرار است

ق ا خاصیت دقیق برازش اگر تعداد مشاهداتی که روی θ قرار گیرند k اشد که $\leq k \leq n$ نگاه

$$k \leq rdepth(\theta, Z_n) \leq \left\lfloor \frac{n+k}{2} \right\rfloor$$

رای $k = n$ دار m $rdepth(\theta, Z_n) = n$

□ تو ح ۱ از ن جاکه روردکننده L^1 هم شه از حداقل p مشاهده می‌گذرد لذا رگرسیون ژرفای ن حداقل p است لومفلد واسته‌گر^۴ را بنده همچنین این خاصیت رای روردکننده کمتر ن توان‌های دوم پر استه (LTS) که ه سه له روسف م رح شده رقرار است روسف

□ تو ح روردکننده کمتر ن توان‌های دوم (LS) هرگز ناراها ن است در حقیقت اگر Z_n ه صورتی اشد که $X_n = \{x_1, x_2, \dots, x_n\}$ رتبه کامل داشته اشد و $n \geq p$ سپس $rdepth(\theta_{LS}, Z_n) \geq$

رگرسیون ژرفا ماکسه مال

کران‌های الا رای رگرسیون ژرفای ماکسه مال در ق ا ز راره می‌شود ز ر مجموعه‌ای از \mathbb{R}^p در موقعیت عمومی قرار دارد هرگاه شتر از p مشاهده در هر زرف ای فن $(p -)$ عدی قرار نگیرد

4) Bloomfield and Steiger

ق ۴ اگر (x_i, y_i) در موقعیت عمومی باشند نگاه

$$\max_{\theta} rdepth(\theta, Z_n) \leq \left\lfloor \frac{n+p}{2} \right\rfloor$$

□

حس ۱ رای هر مجموعه داده $Z_n \subset \mathbb{R}^p$ دارم

$$\max_{\theta} rdepth(\theta, Z_n) \geq \frac{n}{p+1}$$

ژرفترین رگرسیون

ژرفترین رگرسیون T_r^* بعنوان θ ای که $rdepth(\theta, Z_n)$ را بیشینه می‌کند تعریف شده است. ژرفترین رگرسیون T_r^* "متعادل‌ترین" رازش است T_r^* را می‌توان به وسیله محاسبه رگرسیون ژرفای همه رازش‌هایی که از میان p مشاهده می‌گذرد دست برد

تعریف ورودکننده ژرفترین رگرسیون $T_r^*(Z_n)$ عبارت است از رازش θ از بزرگترین رگرسیون ژرفا معنی

$$T_r^*(Z_n) = \arg \max_{\theta} rdepth(\theta, Z_n)$$

توجه ۳ در رگرسیون چندگانه ژرفترین رگرسیون برخلاف L^1 و LTS برای تبدلات کنوای y هم‌وردا است

مثال ۳ می‌خواهم را n میزان بخار مصرفی ماهانه توسط یک ماشین بخار Y را متوسه دمای ماهانه X_1 فارنهایت و تعداد روزهای کار در ماه X_2 را مورد بررسی قرار می‌دهم. n دن من و n مشاهده جموری شده است که در جدول نشان داده شده‌اند. n دراپر و اسمت^۵ مراجعه شود

استفاده از روش‌های رگرسیون مختلف را رگرسیون را ورود کرده و رگرسیون ژرفای هر یک از این رازش‌ها را دست بردم. نتایج در جدول مده است. مده می‌شود که ورودکننده ژرفترین رگرسیون رگرسیون بالاتری نسبت به سایر ورودکننده‌ها دارد و در نتیجه متعادل‌ترین رازش است

خواص پنج‌نگ ژرفترین رگرسیون

در این بخش n ای از خواص هم‌وردای ژرفترین رگرسیون را بیان می‌کنم. همچنین معرفتی کردن مقدار فرورزش n بعنوان معیار پایداری ورودکننده در رازش‌ها پرت نشان می‌دهم که روش ژرفترین رگرسیون n توجه به این معیار روشی استوار است

5) Draper and Smith

جدول داده‌های ماشین بخار

Y	X _۲	X _۱	Y	X _۲	X _۱

جدول مقایسه رگرسیون ژرفای ورودکننده‌های رگرسیون و ا ورودکننده ژرفترین رگرسیون

روش	LS	LAD	رورد هور	روردکننده ژرفترین رگرسیون
$\hat{\beta}_0$				
$\hat{\beta}_1$				
$\hat{\beta}_2$				
رگرسیون ژرفا				

۱ هم‌وردای ژرفترین رگرسیون

بررسی خواص فوق در قه زیر مده است
 قه ۴ ۳ مجموعه داده $Z_n = \{(X_i, Y_i); i = 1, \dots, n\}$ که در ن

$$X_i = (x_{i,1}, \dots, x_{i,p-1}) \in \mathbb{R}^{p-1}$$

و ورودکننده ژرفترین رگرسیون $T_r^*(Z_n) = \arg \max_{\theta} rdepth(\theta, Z_n)$ را در ن ر می‌گرم

a ورودکننده ژرفترین رگرسیون و رگرسیون هم‌وردا است معنی رای هر ردار ستونی V

$$T_r^*(Z'_n) = T_r^*(Z_n) + V$$

که $Z'_n = \{(X_i, y_i + X_i V); i = 1, \dots, n\}$ و $T_r^*(Z'_n) = \arg \max_{\theta} rdepth(\theta, Z'_n)$

b ورودکننده ژرفترین رگرسیون مقاس هم‌وردا است معنی رای هر c ثابت

$$T_r^*(Z'_n) = c T_r^*(Z_n)$$

که $Z'_n = \{(X_i, cy_i); i = 1, \dots, n\}$ و $T_r^*(Z'_n) = \arg \max_{\theta} rdepth(\theta, Z'_n)$

ورودکننده ژرفترین رگرسیون ف ن هم‌وردا است معنی رای هر ماتریس مر ناوژه A

$$T_r^*(Z'_n) = A^{-1} T_r^*(Z_n)$$

که $Z'_n = \{(X_i A, y_i); i = 1, \dots, n\}$ و $T_r^*(Z'_n) = \arg \max_{\theta} rdepth(\theta, Z'_n)$

استوار ژرفترین رگرسیون

ک معار معروف رای پاداری ورودکننده در رار نقا پرت مقدار فرورزش است نا ه
 تعرف داناوو و هو ر^۶ مقدار فرورزش هر ورودکننده T_n ه صورت ز راست

$$\varepsilon_n^*(T_n, Z_n) = \min \left\{ \frac{k}{n}; \sup_{Z'_n} \|T_n(Z'_n) - T_n(Z_n)\| = \infty \right\}$$

که Z'_n شامل همه مجموعه داده‌های دست مده با جاگزی نبی هر k مشاهده دلخواه Z_n ما
 مقادیر اختاری است بنابراین مقدار فرورزش کوچکترین کسری از مشاهدات تبدیل شده
 مغشوش شده ای است که می‌تواند ورودکننده را ه دلخواه منحرف کند لازم ه ذکر است که

6) Donoho and Huber

این اغتشاش فقط θ نقطه پرت در y_i محدود نمی‌شود بلکه Z'_n می‌تواند شامل نقاط پرت در x_i ها نیز باشد

حدس رای هر مجموعه داده $Z_n \subset \mathbb{R}^p$ x_i ها در موقعیت عمومی باشند نگاه

$$\varepsilon_n^*(T_r^*, Z_n) \geq \frac{1}{n} \left(\frac{n}{p+1} - p + 1 \right) \approx \frac{1}{p+1}$$

را θ آن می‌کند که مقدار فرورزش ژرفترین رگرسیون همواره مثبت است این مقدار زمانی برابر با کران پاین $\frac{1}{p+1}$ می‌شود که داده‌های اصلی خود ویژه باشند مثلاً وقتی n ها روی منحنی گشتاور قرار گیرند

میان ژرفترین رگرسیون در برابر نقاط نافذ همانند نقاط دورافتاده عمودی استوار است علاوه بر ژرفترین رگرسیون متفاوت از رگرسیون L^1 تعریف شده θ صورت

$$L^1(Z_n) = \arg \min_{\theta} \sum_{i=1}^n |r_i(\theta)|$$

است که در نتیجه سه پذیری n در برابر نقاط نافذ می‌باشد

مثال‌ها

در این بخش رویدکننده ژرفترین رگرسیون را برای داده‌های موجود دست ورده و نشان می‌دهیم که در مقایسه با رویدکننده‌های دیگر دارای شتاب رگرسیون ژرفاست این مقایسه با محاسبه رگرسیون ژرفای رویدکننده‌ها انجام می‌گردد. مدن مزور از الگوریتم‌هایی که برنامه مرور θ نه با استفاده از نرم‌افزار *Plus* و زبان برنامه‌نویسی فرترن نوشته شده است و از سایت <http://win-www.uia.ac.be/u/statis/> قابل دسترسی است استفاده می‌گردد

۱ مثال

در یک تحقیق در سال میزان لودگی رودخانه در ایالت نیویورک بررسی شده است هدف بررسی سهم هر نو از زمین‌های مورد استفاده در ایالت رودخانه بر روی لودگی رودخانه‌هاست این لودگی براساس ماندگن غلات نیتروژن مملی‌گرم اندازه‌گیری می‌شود نو زمین‌های مورد استفاده در ایالت رودخانه‌ها و میزان لودگی n ها در زیر بیان شده است داده‌ها در جدول نشان داده شده‌اند به چاترچی و همکاران^۷ مراجعه شود

Y ماندگن غلات نیتروژن مملی‌گرم را کمتر براساس نمونه‌هایی که در فواصل m م‌هاار تاستان و ماههای پایز گرفته شده است

7) Chatterjee and et al.

جدول داده‌های رودخانه‌های زوورک

Y	X_4	X_3	X_2	X_1	Y	X_4	X_3	X_2	X_1

جدول ورود پارامترهای مدل و رگرسیون ژرفای ورودکننده‌های مختلف رگرسیون

روش	LS	LAD	رورد هور	روردکننده ژرفترین رگرسیون
β_0				
β_1				
β_2				
β_3				
β_4				
رگرسیون ژرفا				
$\sum e_i^2$				
$\sum e_i $				

X_1 کشاورزی درصد مساحت زمین مورد استفاده رای کشاورزی

X_2 جنگلی درصد زمین جنگلی

X_3 مسکونی درصد مساحت زمین مورد استفاده رای سکونت

X_4 تجاری صنعتی درصد مساحت زمین مورد استفاده تجاری یا صنعتی

ا بررسی داده‌های فوق داده می‌شود که مشاهدات و دور افتاده و نافذ و در عن حال

موثر هستند بنابراین آن‌ها را می‌رود که روش‌های LS و LAD و ورود هور تحت تأثیر این نقا قرار گیرند و و قابل ملاحظه رگرسیون ژرفای کوچکتری نسبت به ورودکننده ژرفترین رگرسیون داشته باشند با نگاه کردن به جدول مشخص می‌شود که ورودکننده رگرسیون LAD و ورود هور نزدیک به هم هستند و دارای رگرسیون ژرفای هستند چون ورودکننده ژرفترین رگرسیون نسبت به نقا دور افتاده و نافذ استوار است لذا تحت تأثیر این نقا قرار نمی‌گیرد و نه ر می‌رسد که رازشی مناسب رای داده‌هاست با توجه به اینکه دارای شترین رگرسیون ژرفا معنی است

جدول جدول ورود پارامترهای مدل و رگرسیون ژرفای ورودکننده های مختلف رگرسیون

روش	LS	LAD	رورد هور	روردکننده ژرفترین رگرسیون
$\hat{\beta}_0$				
$\hat{\beta}_1$				
$\hat{\beta}_2$				
$\hat{\beta}_3$				
$\hat{\beta}_4$				
رگرسیون ژرفا				
$\sum e_i^2$				
$\sum e_i $				

مثال شبیه ساز

رای بررسی مناسب بودن ورودکننده ژرفترین رگرسیون و نقه در عدم افاق مدل $y = -x_1 + x_2 - x_3 + x_4 + e$ تولد می کنیم که x_1, x_2, x_3, x_4 و e از توزیع نرمال استاندارد گرفته شده اند پس از محاسبه y با استفاده از مقادیر تولد شده x_1, x_2, x_3, x_4 و e ورودکننده های رگرسیون دست از روش های مختلف و ورودکننده ژرفترین رگرسیون را برای این مقادیر دست می وریم نتایج دست مده در جدول زیر مده اند

ما مقایسه را ورود شده دست مده از روش های مختلف مده می شود که را دست مده از روش ژرفترین رگرسیون نسبت به سایر روش ها را مدل مفروض نزدیک تر است و در این حال دارای شترین رگرسیون ژرفا است عبارت دیگر ورودکننده ژرفترین رگرسیون

$$y = -x_1 + x_2 - x_3 + x_4$$

دارای رگرسیون ژرفای تقریبی است

نتیجه گیری

وقتی که استواری و تشخیص نقه دورافتاده اهمیت دارند روش LMS روشی مناسب است اما زمانی که نقه دورافتاده زیادی وجود ندارد می توان به وسیله روش های LMS یا LTS بررسی کرد و کتوایی اهمیت دارد روش ژرفترین رگرسیون T_n^* انتخابی مناسب است

مراج

- [1] Bloomfield, P., and Steiger, W. (1983). Least Absolute Deviations: Theory, Applications, and Algorithms, Boston: Birkhauser.
- [2] Chatterjee, S., and Hadi, A.S., and Price, B. (2000), Regression Analysis by Example, New York: John Wiley.
- [3] Donoho, D.L., and Huber, P.J. (1983), "The notion of breakdown point", in A Festschrift for Erich Lehmann, eds. Bickel, P., Doksum, K., and Hodges, J.L., Belmont: Wadsworth.
- [4] Draper, N., Smith, H. (1998) Applied Regression Analysis, third edition, John Wiley & Sons.
- [5] Rousseeuw, P.J., (1984) "Least Median of Squares Regression", Journal of the American Statistical Association, **79**, 871-880.
- [6] Rousseeuw, P.J., and Hubert, M. (1999), "Regression Depth", Journal of the American Statistical Association, **94**, 388-402.

Archive of SID