

# رگرسیون چندکی خطی

(تحقیق درس مدل های خطی)

1393/10/29

## ارائه دهندگان

راحله موسوی

مریم حسن پور

استاد مربوطه

دکتر ریحانه ریخته گران

## در چه مواقعی از رگرسیون چندکی استفاده می کنیم؟

- ❖ برای بررسی ارتباط دو متغیر، معمولاً رگرسیون حداقل مربعات به کار گرفته می شود. با رگرسیون حداقل مربعات ارتباط بین متغیر کمکی و میانگین پاسخ شرطی را می توان برآورد کرد، اما این روش مثل بسیاری از روش های آماری دارای کمبودهایی است.
- ❖ رگرسیون حداقل مربعات در صورتی که توزیع داده ها، برای تخمین ارتباط بین متغیر کمکی و متغیر پاسخ اعتبار کافی نباشد، مشکل ساز می شود.
- ❖ از طرفی رگرسیون حداقل مربعات تنها ارتباط بین متغیرهای کمکی و میانگین پاسخ را تخمین می دهد، در حالی که در بسیاری موارد هدف پیدا کردن ارتباط بین متغیرهای کمکی با سایر بخشهای توزیع به ویژه چندک های انتهایی توزیع می باشد.

**یک روش ارزشمند در مواجهه با چنین مشکلاتی استفاده از رگرسیون چندکی است. (Quantile Regression)**

رگرسیون چندکی یک روش آماری با قابلیت محاسبه و رسم منحنی های رگرسیونی متفاوت و منطبق با نقاط صدکی مختلف میباشد، که ضمن بیان تصویری کامل تر و جامع تر از داده ها، امکان سنجش ارتباط متغیرهای مستقل با چندک های مورد نظر متغیر وابسته را **بدون نیاز به نرمال بودن داده ها** و حتی در **حضور نقاط دور افتاده** فراهم میکند یعنی این رگرسیون نسبت به داده های دورافتاده نیرومند می باشد. از سوی دیگر برخلاف رگرسیون حداقل مربعات که روی میانگین شرطی یعنی پارامتر مکان متمرکز است، رگرسیون چندکی استراتژی منظمی را برای تعیین چگونگی تاثیر متغیرهای مستقل روی **مکان** و **مقیاس** و **شکل توزیع** پیشنهاد میکند.

به طور خلاصه رگرسیون چندکی مدلی است که به بیان چگونگی تاثیر متغیرهای مستقل بر چندک های دلخواه متغیر پاسخ می پردازد.



## مزایای استفاده از رگرسیون چندکی



### معایب

در ابتدای معرفی رگرسیون چندکی به دلیل **نبود** کامپیوتر و **نرم افزار های آماری** محاسبه رگرسیون چندکی بسیار **سخت و خسته کننده** بود که البته با آمدن نرم افزار های آماری این مشکل کاملاً برطرف شد.

❖ پیدا کردن ارتباط بین متغیرهای کمکی با سایر بخش های توزیع به ویژه چندکهای انتهایی توزیع میباشد.

❖ نیازی به شرط نرمال بودن ندارد

❖ نسبت به داده های دور افتاده نیرومند است

❖ چگونگی تاثیر متغیرهای مستقل روی مکان و مقیاس و شکل توزیع را نشان می دهد.

## آشنایی با چندک ها

• اگر  $Y$  یک متغیر تصادفی با تابع توزیع  $F(y) = p(Y \leq y)$  باشد  $\tau$  امین چندک  $Y$  برابر است با

$$Q(\tau) = F^{-1}(\tau) = \inf\{y: F(y) \geq \tau\}$$

اگر  $\{y_i, i=1,2,\dots,n\}$  نمونه ای تصادفی از متغیر تصادفی  $Y$  با تابع توزیع  $F$  باشد چندک  $\tau$  ام نمونه به صورت زیر محاسبه می شود:

$$\begin{aligned} & \min\left\{ \sum_{i \in \{i: y_i \geq u\}} \tau |y_i - u| + \sum_{i \in \{i: y_i < u\}} (1 - \tau) |y_i - u| \right\} \\ & = \min_u \left\{ \sum_{i \in \{1, \dots, n\}} \rho_\tau(y_i - u) \right\} \end{aligned}$$

$$\rho_\tau(u) = u(\tau - I_{(u < 0)}), 0 < \tau < 1$$

# مثال 1

فرض کنید  $Y$  یک متغیر تصادفی گسسته باشد که مقادیر  $1, 2, \dots, 9$  را با احتمال مساوی اختیار می کند، قصد داریم میانه  $Y$  را به دست آوریم. با استفاده از فرمول چندک نمونه ای ارائه شده داریم :

برای  $\tau = 0.5$

$$L(u) = \min_u \left\{ \sum_{i \in \{i: y_i \geq u\}} 0.5|y_i - u| + \sum_{i \in \{i: y_i < u\}} (1 - 0.5)|y_i - u| \right\}$$

**U=3** 

$$L(3) \propto \sum_{i=1}^2 -(i-3) + \sum_{i=3}^9 (i-3) = [(2+1) + (0+1+2+\dots+6)] = 24.$$

**U=5** 

$$L(5) \propto \sum_{i=1}^4 i + \sum_{i=0}^4 i = 20,$$



$u$	1	2	3	4	5	6	7	8	9
Expected loss	36	29	24	21	20	21	24	29	36

## رگرسیون چندکی (چندک های شرطی)

برخلاف رگرسیون حداقل مربعات که برای برآورد میانگین شرطی از مینیمم کردن مانده ها استفاده میکند، روش رگرسیون چندکی برای برآورد چندکهای شرطی از مینیمم کردن قدر مطلق موزون مانده های نامتقارن استفاده می کند.

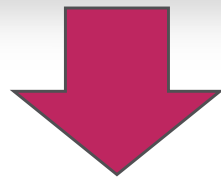
فرض کنید  $\tau$  امین چندک شرطی به صورت  $Q_{Y|X}(\tau) = X\beta_\tau$  باشد، در این صورت برای برآورد  $\beta_\tau$  از رابطه زیر استفاده می شود :

$$\min_{\beta \in \mathbb{R}^p} \left[ \sum_{i \in \{i: y_i \geq x_i' \beta\}} \tau |y_i - x_i' \beta(\tau)| + \sum_{i \in \{i: y_i < x_i' \beta\}} (1 - \tau) |y_i - x_i' \beta(\tau)| \right]$$
$$= \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \rho_\tau(y_i - x_i' \beta(\tau))$$

روش عددی

$$\rho_\tau(u) = u(\tau - I_{(u < 0)}), 0 < \tau < 1$$

$$\hat{\beta}_\tau = \arg \min_{\beta \in \mathbb{R}^k} \sum_{i=1}^n \rho_\tau(y_i - x\beta)$$



$$Q_{Y|X}(\tau) = X\hat{\beta}_\tau$$

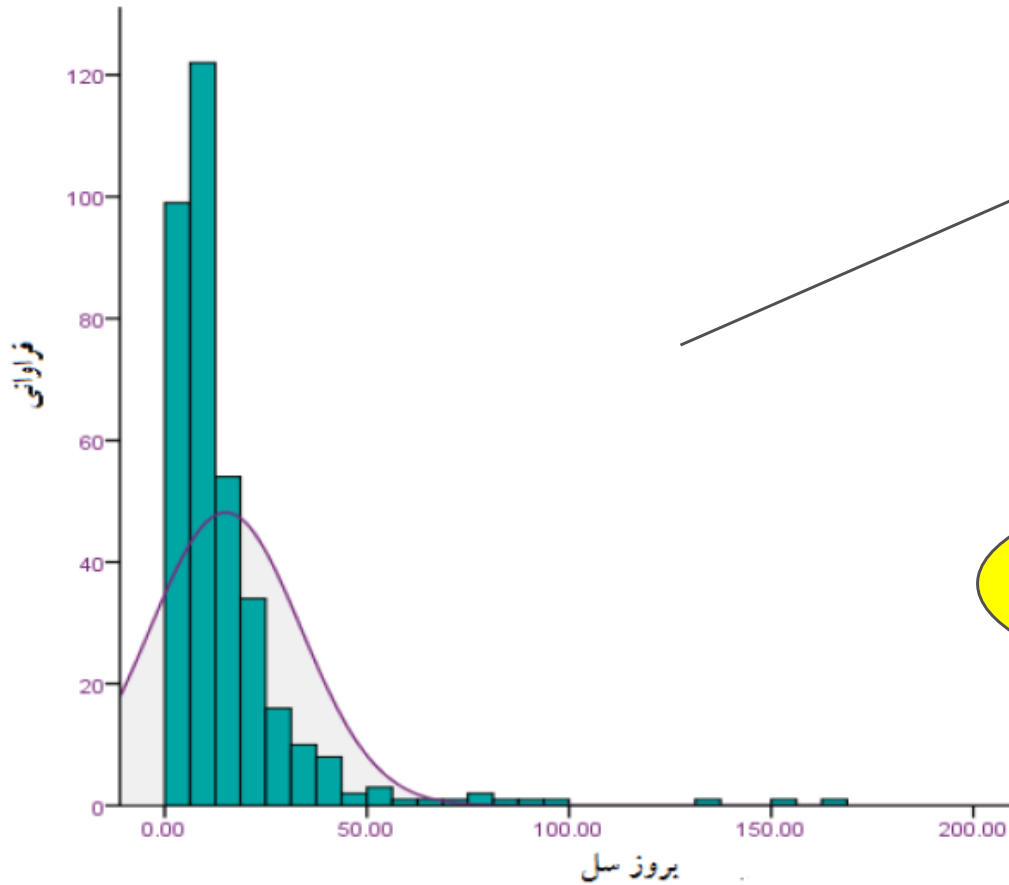


## مثال 2

- در مطالعه ای اطلاعات مربوط به 11320 فرد مبتلا به بیماری سل در 359 شهرستان ایران تجزیه و تحلیل شد. هدف از این مطالعه بررسی اثر نرخ مهاجرت در میزان بروز بیماری سل بوده است. هیستوگرام مربوط به مشخصه **میزان بروز سل** به صورت زیر رسم شده است.

چوله به راست میزان بروز سل  
دارای نقاط پرت حائز اهمیت

پس با توجه به این شرایط ترجیح می  
دهیم از رگرسیون چندکی استفاده کنیم



جدول ۱. نتایج برازش مدل رگرسیون چندکی برای بررسی ارتباط میزان بروز سل و نرخ مهاجرت

صدک	نرخ مهاجرت	
	$\hat{\beta}(\tau)$	خطای معیار
	$p$	
5	-۰/۰۰۵	۰/۰۷۲
10	۰/۰۵	۰/۰۷
20	۰/۰۵	۰/۰۷۱
30	۰/۰۸۹	۰/۰۹
40	۰/۱۷	۰/۰۹۸
50	۰/۲۶	۰/۱۲
60	۰/۲۹	۰/۱۲
70	۰/۴۲	۰/۱۵
80	۰/۵	۰/۱۵
90	۰/۴	۰/۵۷
95	۱/۹۱	۰/۸۳
رگرسیون حداقل مربعات	۰/۵۷	۰/۱۴
	<۰/۰۰۱	

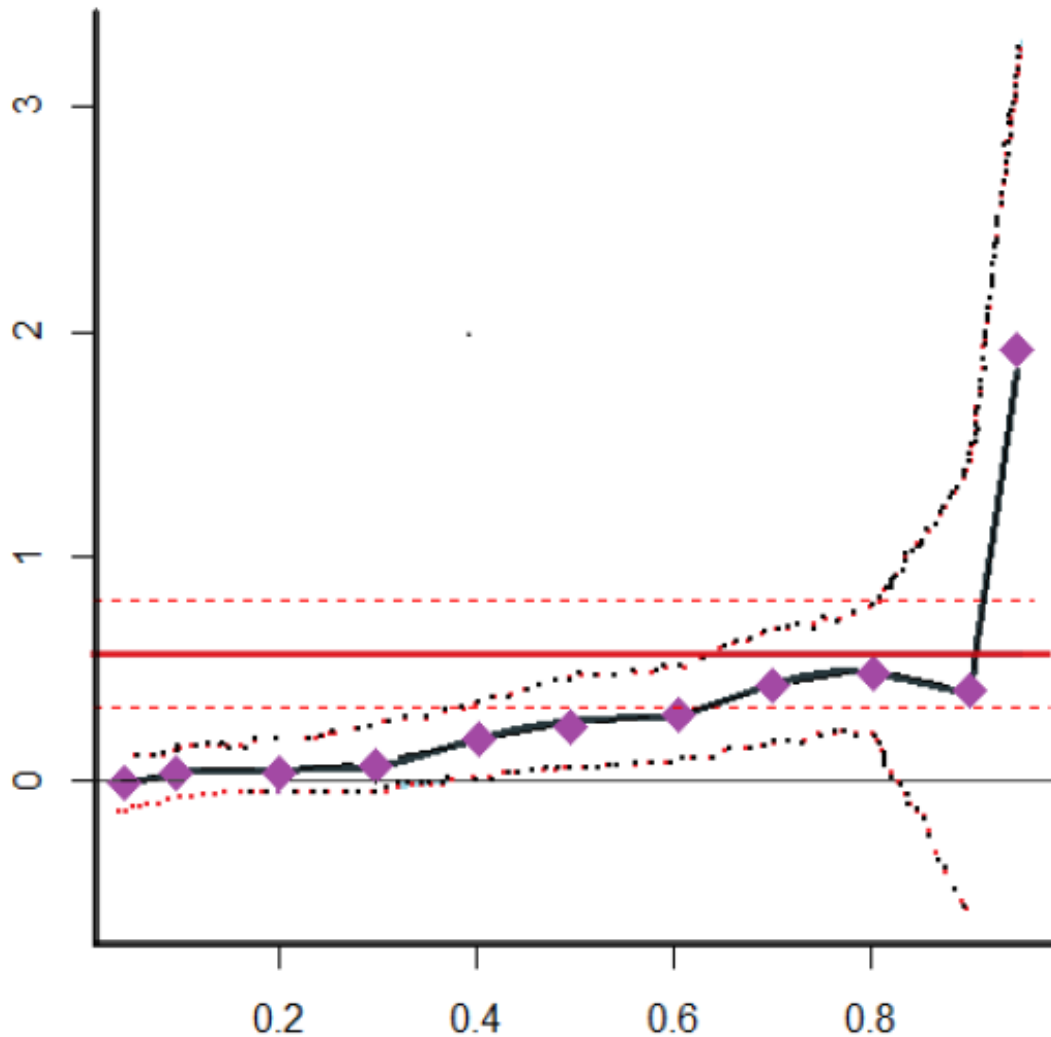
مقدار  $p < 0/05$  از نظر آماری معنی دار است.

مدل رگرسیون چندکی برای صدک های 5, 10, 20, 30, 40, 50, 60, 70, 80, 90, 95 به طور جداگانه برازش داده شد و ضرایب رگرسیون مربوط به هر صدک با استفاده از رابطه های تعریف شده برآورد شده اند. نتایج برازش مدل رگرسیونی چندکی در جدول روبرو نشان داده شده است.

$$Q_{Y|X}(\tau) = X\beta_{\tau}$$

## نمودار ضرایب مدل رگرسیون چندکی

این نمودار برآوردهای نقطه ای ضرایب رگرسیون چندکی را برای هر چندک شرطی نشان میدهد.



صدک های سطوح بالاتر میزان بروز سل اثر پذیری بیشتری نسبت به نرخ مهاجرت دارد. در واقع به ازای یک واحد افزایش نرخ مهاجرت صدک های سطوح بالاتر بروز سل بیشتر از صدک های پایینی افزایش میابند.

## مثال 3

در مطالعه ای هدف بررسی نقش سن در سلامت روانی افراد در یک جامعه می باشد. این پژوهش رگرسیون چندکی رابطه سن با سلامت روان مردان و زنان را به طور جداگانه به گونه ای نشان می دهد که با رگرسیون میانگین نمایش پذیر نیست.

● در این مطالعه شاخص سلامت روان با پاسخ گویی افراد مورد نظر به سوالات مربوطه با نمراتی بین 28 تا 118 مشخص شده است به طوری که **مقادیر پایین تر** نشان دهنده سلامت روان **بهتر** می باشند.

● قصد اصلی این مطالعه توصیف توزیع کل این شاخص با استفاده از نه چندک شرطی مختلف است. به عبارت دیگر اثر متغیر توضیحی سن را بر شکل توزیع شاخص سلامت روان بررسی می کنیم .

جدول زیر برآورد پارامترهای مدل های برازش داده شده در هر دو حالت رگرسیونهای چندک از روش بیان شده و رگرسیونهای میانگین از روش حداقل مربعات را فراهم می سازد. ستاره ها در این جدول نشان دهنده برآورد پارامترهایی هستند که در سطح 0/05 معنا دارند.

چندک	مردان	زنان
	سن	سن
0/1	-0/1	0/05
0/2	-0/047	0/07
0/3	-0/05	0/15*
0/4	-0/045	0/11*
0/5	-0/038	0/15*
0/6	-0/022	0/14*
0/7	0/073	0/21
0/8	0/04	0/21
0/9	0/031	0/16*
OLS	0/036	0/111*

$$Q_{\theta i} = \alpha_{\theta} + \beta_{\theta} x_i \quad i = 1, \dots, n$$

$$\theta \sum_{i/y_i \leq \alpha_{\theta} + \beta_{\theta} x_i} (y_i - \alpha_{\theta} - \beta_{\theta} x_i) + (1 - \theta) \sum_{i/y_i > \alpha_{\theta} + \beta_{\theta} x_i} (y_i - \alpha_{\theta} - \beta_{\theta} x_i)$$

یافته های رگرسیون چندک رابطه سن را با وضع سلامت روانی برای زنان در تمامی موارد اظهار می دارد اما برای مردان چنین رابطه ای معنادار نیست

مثلا مدل رگرسیون چندک برای  $\theta = 0.6$  برای زنان به صورت زیر برآورد شده:

$$Q_{i.6} = 31 + 0.14 \text{AGE}$$

این الگو نشان می دهد که صدک 60 ام شاخص سلامت روان با افزایش یک سال سن 14 درصد افزایش می یابد. پس چندک 0/6 شاخص سلامت به طور مثال برای زنان 40 ساله 36/6 برآورد شده است.

$$Q_{.6} = 31 + 0.14 \times 40 = 36/6$$

پس می توان نتیجه گرفت که 30 درصد زنان 40 ساله دارای شاخص سلامت بین 36/6 تا 53/4 هستند.

این امر آشکار می سازد که 60 درصد از زنان 40 ساله دارای شاخص سلامت روانی 36/6 و 40 درصد بالای آن می باشند.

با ادامه همین محاسبه ها، این نکته به دست آمده که چندک 0/9 شاخص سلامت روان برای زنان 40 ساله، 53/4 می باشد.



## شکل ۱: خطوط رگرسیون چندک تنظیمی برای زنان

- چنانکه ملاحظه می شود، خطوط پایین تر نسبت به خطوط بالاتر به یکدیگر نزدیک ترند این بدان معنا است که شکل توزیع متمایل به راست است .
- میزان اثر پذیری شاخص سلامت از سن زنان در صدک های بالاتر این شاخص بیشتر است.
- تحلیل ما نشان می دهد که نقش سن در سلامت روانی به طور آشکاری برای مردان و زنان متفاوت است به طوری که برای مردان این رابطه تقریباً در همه چندک ها معنادار نبود ولی برای زنان این رابطه مستقیم و برای چندک های پایینی نسبت به چندک های بالایی این میزان تاثیر پذیری کمتر بوده است.



شکل ۱: خطوط رگرسیون چندک تنظیمی برای زنان

○ توجه داشته باشید رگرسیون رایج نمی تواند میزان اثر پذیری متغیر وابسته از متغیر توضیحی را در سطوح مختلف توزیع نشان دهد و فقط با پارامتر میانگین (مکان) مرتبط است در حالی که رگرسیون چندکی ما را قادر به آشنایی با شکل توزیع می نماید.



## نرم افزار

برنامه زیر برای به دست آوردن رگرسیون چندکی در قالب یک مثال برای مدل های مورد نظر پیشنهاد شده است.

```
proc quantreg ci=sparsity/iid algorithm=interior(tolerance=1.e-4)
data=new;
class visit ed;
model weight = black married boy visit ed smoke
cigsper mom_age mom_age*mom_age
m_wtgain m_wtgain*m_wtgain /
quantile= 0.05 to 0.95 by 0.05
plot=quantplot;
run;
```

## مراجع

- Koenker, R.: Quantile regression for longitudinal data. *J. Multivar. Anal.* 91, 74–89 (2004)
- Koenker, R.: *Quantile Regression*. Cambridge University Press, New York (2005)
- Koenker, R., Bassett, G.: Regression quantiles. *Econometrica* 46, 33- 85(1978)
- ارتباط بروز بیماری سل و مهاجرت با استفاده از مدل رگرسیون چندکی در ایران سال 1389 (فاطمه سروی، یدا...محرابی ، علیرضا ابدی، مهشید ناصحی)
- کاربرد رگرسیون چندک در تحلیل سلامت روانی (محمدتقی انصاری، محمد بامنی مقدم، علیرضاخوشگویانفرد، عزت اله سام آرام)
- <http://www.wikipedia.org/>



خیلی ممنون از توجه شما