

ضریب همبستگی رشته‌ای

مهران نقیزاده قمی^۱ شکوفا کبیری^۲

چکیده:

در این مقاله، به بررسی رابطه‌ی بین دو متغیر در حالتی که یک متغیر در مقیاس فاصله‌ای یا نسبتی و متغیر دیگر با مقیاس اسمی اندازه‌گیری شده باشد می‌پردازیم. در چنین حالاتی از ضریب همبستگی رشته‌ای استفاده می‌شود. محاسبه‌ی چند ضریب رشته‌ای شامل ضریب دورشته‌ای و ضریب دورشته‌ای- نقطه‌ای با چند مثال و به کمک نرم‌افزار R انجام می‌شود.

واژه‌های کلیدی: ضریب دورشته‌ای، ضریب دورشته‌ای- نقطه‌ای.

۱ مقدمه

وقتی از یک جامعه‌ی آماری دو متغیر برای نمونه‌ای استخراج کنیم که یکی در مقیاس نسبتی یا فاصله‌ای و دیگری یک متغیر دوحالاتی باشد، برای محاسبه‌ی همبستگی متغیرها از ضریب دورشته‌ای می‌توان استفاده کرد. یکی از پذیره‌های این ضریب همبستگی این است که متغیر دو حالتی باید ذاتاً کمی باشد که عدم وجود وسیله‌ی درست اندازه‌گیری یا سایر محدودیت‌های مطالعه باعث می‌شود نتوان آن را دقیق اندازه‌گیری نمود. بنابراین با فرض این که توزیع واقعی آن نرمال باشد می‌توان آزمون معنی‌داری ضریب دورشته‌ای را انجام داد. پذیره‌ی دوم این است که رابطه‌ی بین متغیر پیوسته و متغیر دوحالاتی باید خطی باشد. ضریب دورشته‌ای به صورت

$$r_{bis} = \frac{d}{s} \frac{pq}{y} \quad (1)$$

محاسبه‌ی شود که در آن $d = \bar{x}_p - \bar{x}_q$ اختلاف بین میانگین‌های متغیر پیوسته در دو طبقه از متغیر دوحالاتی، s انحراف معیار نمونه برای متغیر پیوسته، p و q به ترتیب نسبت حالت‌ها در طبقه‌ی اول و دوم متغیر دو

یکی از مهمترین شاخص‌های رابطه‌ی بین دو متغیر کمی X و Y، ضریب همبستگی پیرسون است. برای توصیف همبستگی دو متغیر کمی بدون هیچ فرضی می‌توان از ضریب همبستگی خطی پیرسون استفاده کرد. اما برای انجام آزمون معنی‌داری ضریب همبستگی پیرسون و قضاؤت در خصوص همبستگی دو متغیر در جامعه، شرط نرمال بودن توزیع توان دو متغیر کمی در جامعه لازم است. این پذیره به کمک یکی از آزمون‌های شاپیرو-ولیک (قابل انجام در نرم‌افزار R به کمک بسته‌ی mvShapiro.Test)، دورنیک-هانسن، هتر-زیکلر و معیار چولگی و کشیدگی ماردیا (قابل انجام در نرم‌افزار STATA 12) قابل بررسی است. در مواردی که یکی از متغیرها دو یا چندحالاتی و متغیر دیگر کمی باشد و توزیع آن در جامعه نرمال فرض شود، از ضرایب همبستگی مستخرج از فرمول گشتاوری پیرسون به نام ضرایب رشته‌ای استفاده می‌شود. قصد داریم این ضرایب را با چند مثال مورد بررسی قرار دهیم.

^۱ استادیار گروه آمار دانشگاه مازندران m.naghizadeh@umz.ac.ir

^۲ کارشناس آمار، دانشگاه مازندران

```

d<-mean(x)-mean(y)
s<-sd(c(x,y))
n=length(x)+length(y)
p<-length(x)/n
q<-length(y)/n
y<-dnorm(qnorm(min(p,q)))
r<-d/s*p*q/y
z=y/sqrt(p*q)*r
cat("r.bis=",r,"\n")
t<-sqrt(z^2*(n-2)/(1-z^2))
cat("t.statistic=",t," nsig=",
2*min(pt(t,n-2),1-pt(t,n-2)), "\n")
y<-c(18.3,17.6,12.3,9.2,16.1,12.1,11.2,10.6)
x<-c(9.1,10.9,10.8,8.1,11.2,9.2,12.6,10.9,9.2,
9.4,9.2,11.4,8.9,14.6,11.3,8.1)
cor.bis(x,y)

```

خروجی R به صورت زیر است:

```

r.bis=-0.6826444
t.statistic= 2.904941
sig= 0.008211725

```

چون گروهی که ۶ ساعت در هفته تمرین می‌کنند ($\bar{x}_p = 10/31 = 10/31$) دارای میانگین زمان شنای کمتری در مقایسه با گروهی که ۳ ساعت در هفته تمرین می‌کنند ($\bar{x}_q = 13/43 = 13/43$) هستند، همبستگی بین زمان تمرین و زمان شنا $6/68 = 0/008 < 0/05$ (sig = ۰/۰۰۸) به دست آمد. چون زمان کمتر در نتیجه‌ی شنای سریع‌تر حاصل می‌شود، همبستگی مثبتی بین زمان تمرین و سرعت شنا وجود دارد. با توجه به مقدار معنی‌داری ($0/05 < \alpha = 0/008$)، فرضیه‌ی صفر در سطح معنی‌داری $r_{bis} = -0/68$ نشان می‌دهد که همبستگی مثبت معنی‌داری بین زمان تمرین و سرعت شنا وجود دارد.

تذکر ۲۰۲. توزیع دقیق نمونه‌گیری r_{bis} موجود نیست ولی تحت شرایطی می‌توان از تبدیل لگاریتمی فیشر استفاده کرد. اگر نسبت‌های p و q در جامعه تقریباً برابر با $1/2$ (تعداد کل آزمودنی‌ها) بزرگ باشد و r'_{bis} به r_{bis} نزدیک نباشد می‌توان نشان داد که

$$r'_{bis} = \frac{1}{\sqrt{n}} \ln \frac{1+r_{bis}}{1-r_{bis}}$$

حالی و y ارتفاع در نقطه‌ای است که منحنی نرم‌الاستاندارد به دو قسمت p و q تقسیم می‌شود که با دستور $dnorm(qnorm(min(p,q)))$ در نرم‌افزار R قابل محاسبه است. برای آزمون $H_0: \rho_{bis} = 0$ در مقابل آماره‌ی آزمون به صورت $H_1: \rho_{bis} \neq 0$.

$$t = \sqrt{\frac{\left[\frac{y}{\sqrt{pq}} r_{bis}\right]^2 (n-2)}{1 - \left[\frac{y}{\sqrt{pq}} r_{bis}\right]^2}} \quad (2)$$

دارای توزیع مجانی t -استودنت با $n-2$ درجه‌ی آزادی است که در آن n حجم نمونه است.

مثال ۱۰۲. می‌خواهیم همبستگی بین زمان تمرین در هفته و زمانی که برای شنا کردن ۲۵ متر لازم است را تعیین کنیم. تعداد ۸ نفر به مدت ۳ ساعت در هفته و تعداد ۱۶ نفر، ۶ ساعت در هفته تمرین می‌کنند. زمان تمرین یک متغیر دو حالتی غیرواقعی و زمان شنا ۲۵ متر یک متغیر پیوسته است. داده‌ها در جدول ۱ نشان داده شده است.

جدول ۱. زمان شنای ۲۵ متر و زمان تمرین در هفته

۳ ساعت تمرین هفتگی	۶ ساعت تمرین هفتگی
۳/۱۸	۱/۹
۶/۱۷	۹/۱۰
۳/۱۲	۸/۱۰
۲/۹	۱/۸
۱/۱۶	۲/۱۱
۱/۱۲	۲/۹
۲/۱۱	۹/۱۲
۶/۱۰	۹/۱۰
	۲/۹
	۴/۹
	۲/۹
	۴/۱۱
	۹/۸
	۶/۱۴
	۳/۱۱
	۱/۸

برای محاسبه‌ی ضریب همبستگی رشته‌ای، از دستورهای R بهره می‌گیریم.

```
cor.bis<-function(x,y){
```

برنامه‌ی R برای محاسبه‌ی ضریب همبستگی دورشته- نقطه‌ای به صورت زیر است:

```
cor.pbis.test<-function(x,y){
d<-mean(x)-mean(y)
s<-sd(c(x,y))
n=length(x)+length(y)
p<-length(x)/n
q<-length(y)/n
r<-(d/s)*sqrt(p*q)
cat("r.pbis=",r,"\\ n")
t<-sqrt(r^2*((n-2)/(1-r^2)))
cat("t.statistic=",t,"\\ nsig="
,2*min(pt(t,n-2),1-pt(t,n-2)),"\\ n")
x<-c(35,36,38,39,50,51,55)
y<-c(30,33,34,36,37,40,42)
cor.pbis.test(x,y)}
```

خروجی R به صورت زیر است:

```
r.pbis= 0.5044248
t.statistic= 2.023704
sig= 0.06585428
```

با توجه به مقدار معنی‌داری ($sig = 0.06 > 0.05$)، فرضیه‌ی صفر در سطح معنی‌داری $\alpha = 0.05$ رد نمی‌شود و نتیجه‌ی آن که رابطه‌ی معنی‌داری بین وسعت دبیرستان و توانایی در مهارت بسکتبال وجود ندارد.

قدکو ۲۰۳. می‌توان نشان داد که هر چه تفاوت p و q بیشتر باشد، تفاوت r_{pbis} بیشتر است. در عمل از r_{pbis} کمتر استفاده می‌شود و فقط موقعی کاربرد دارد که از نرمال بودن توزیع واقعی متغیر دوچالی اطمینان داشته باشیم. در صورت عدم اطمینان، استفاده از r_{pbis} مناسب‌تر است.

دارای توزیع نرمال با خطای استاندارد نمونه‌گیری $\frac{1}{\sqrt{n}}$ است. همچنین کارایی r_{pbis} برای برآورد ρ وقتی $0 = \rho$ باشد خوب و برای مقادیر ρ نزدیک به $1 \pm$ نامناسب می‌شود.

۳ ضریب دورشته‌ای- نقطه‌ای

ضریب دورشته‌ای- نقطه‌ای مانند ضریب دورشته‌ای، شاخص رابطه‌ی بین یک متغیر پیوسته در مقیاس حداقل فاصله‌ای و یک متغیر دوچالی واقعی مانند جنسیت را فراهم می‌کند. ضریب دورشته‌ای- نقطه‌ای به صورت

$$r_{pbis} = \frac{d}{s} \sqrt{pq} \quad (3)$$

محاسبه‌ی می‌شود که در آن $d = \bar{x}_p - \bar{x}_q$ اختلاف بین میانگین‌های متغیر پیوسته در دو طبقه از متغیر دوچالی، s انحراف معيار نمونه برای متغیر پیوسته و p و q به ترتیب نسبت حالت‌ها در طبقه‌اول و دوم متغیر دوچالی هستند. برای آزمون $H_0 : \rho_{bis} = 0$ در مقابل $H_1 : \rho_{bis} \neq 0$ آماره‌ی آزمون به صورت

$$t = \sqrt{\frac{r_{pbis}^2(n-2)}{1-r_{pbis}^2}} \quad (4)$$

دارای توزیع معجانی t - است و درجه‌ی آزادی است که در آن $n-2$ درجه‌ی آزادی است که در آن n حجم نمونه است.

مثال ۱۰۳. می‌خواهیم همبستگی بین نمره‌های آزمون مهارت بسکتبال و وسعت دبیرستانی که بازیکنان قبلًا در آن مسابقه داده‌اند را تعیین کنیم. جدول ۲ داده‌ها را نشان می‌دهد.

جدول ۲. نمره‌های آزمون مهارت و وسعت دبیرستان

وسعت کم	وسعت زیاد
۳۰	۳۵
۳۳	۳۶
۳۴	۳۸
۳۶	۳۹
۳۷	۵۰
۴۰	۵۱
۴۲	۵۵

مراجع

- [۱] میناسیان، و. (۱۳۸۸)، آمار در تربیت بدنی و علوم ورزشی، انتشارات علم و حرکت.
- [۲] Walker, H. M. (1953), *Statistical Inference*, New York:Holt, Rinehart and Winston.