

# روش درست‌نمایی نیمرخ در تحلیل شواهدی داده‌های ژنتیک

مهدی عمادی<sup>۱\*</sup>، انیس سنچولی<sup>۲</sup> مرتضی محمدی<sup>۳</sup>

<sup>۱</sup> دانشگاه فردوسی مشهد

<sup>۲</sup> دانشگاه فردوسی مشهد

<sup>۳</sup> دانشگاه زابل

۱۴۰۳/۲/۳۱

---

چکیده: در علم آمار سه مکتب به نام‌های مکتب فراوانگرایان، مکتب بی‌زین‌ها و مکتب شواهدی وجود دارد. در مکتب فراوانگرایان تنها به مشاهدات و بسامد رخدادها استناد می‌شود و بر حسب آن می‌توان مسائل را حل کرد، در حالی که در مکتب بی‌زین‌ها علاوه بر مشاهدات، اطلاعات و باورهای اولیه محقق نیز مهم است و در حل مسئله و نتیجه‌گیری مورد نظر قرار می‌گیرند و در مکتب شواهدی برای اندازه‌گیری میزان پشتیبانی داده‌ها از فرضیه مفروض، دو معیار، یکی برای اندازه‌گیری کفایت شواهد (شامل داده‌ها و اطلاعات مسئله) و دیگری اندازه‌گیری ارتباط شواهد (شامل فرضیه‌های مسئله) وجود دارد، که این دو معیار از هم جدا هستند. بنابراین، تنها معیاری که هر دو را اندازه‌گیری می‌کند، شاخص نسبت درست‌نمایی است.

---

\*نویسنده مسئول: جناب آقای دکتر مهدی عمادی emadi@um.ac.ir

واژه‌های کلیدی: رهیافت شواهدی، نسبت درست‌نمایی، درست‌نمایی نیمرخ، رگرسیون لجستیک.

## ۱ مقدمه

از اواخر سده بیستم با چاپ و انتشار مقالات و کتاب ریچارد رویال<sup>۱</sup> (۱) رهیافت فکری جدیدی تحت عنوان رهیافت شواهدی که مبتنی بر مفهوم درست‌نمایی است شکل گرفت و رقیبی برای دو رهیافت دیگر که در مقالات افران<sup>۲</sup> (۲) و گوستافسون<sup>۳</sup> (۳) وجود دارد شد. تابع درست‌نمایی باید برحسب پارامتر مورد علاقه باشد، ساده‌ترین راه حذف پارامترهای مزاحم از تابع درست‌نمایی است که روش‌های گوناگون در رویال (۱) وجود دارد، در این مقاله روش درست‌نمایی نیمرخ برای تجزیه و تحلیل داده‌های ژنتیک در رابطه با انتقال بیماری سیستمیک فیبروزیست بررسی شده و داده‌ها از مقاله استرانگ (۴) و از سایت<sup>۴</sup> انجمن ژنتیک بین المللی مدیفایر<sup>۵</sup> گرفته شده.

### ۱.۱ تابع درست‌نمایی نیمرخ

در این روش ماکسیمیم تابع درست‌نمایی را زمانی که مقدار پارامتر مورد علاقه ثابت در نظر گرفته شود، بدست می‌آید. در این روش می‌توان پارامتر  $\gamma$  را با روش جایگزینی با یک برآورد مقدار واقعی آن،  $\gamma_0$  به ازای هر مقدار  $\theta$ ، از مقدار  $\hat{\gamma}(\theta)$ ، که  $L(\theta, \gamma)$  را بیشینه می‌سازد، استفاده می‌شود. تابعی که از این طریق به صورت

---

Royall Richard<sup>۱</sup>

Efran<sup>۲</sup>

Gustafson<sup>۳</sup>

<http://locuszoom.sph.umich.edu><sup>۴</sup>

<sup>۵</sup>شامل گروه‌هایی از کشورهای فرانسه، کانادا و دانشگاه‌های جانز هاپکینز و سترن است.

زیر بدست می‌آید تابع درست‌نمایی نیم‌رخ نامیده می‌شود.

$$L_P(\theta) = L(\theta, \hat{\gamma}(\theta)) = \max_{\gamma} (L(\theta, \gamma))$$

مثال ۱. فرض کنید متغیر تصادفی  $X$  دارای توزیع نرمال با میانگین مجهول  $\mu$  و واریانس مجهول  $\delta^2$  باشد. برای محاسبه ی استنباط شواهدی درباره پارامتر مورد علاقه‌ی  $\mu$  بر اساس نمونه تصادفی به حجم  $n$  از توزیع فوق داریم:

$$L(\mu, \delta^2) = \left(\frac{1}{\sqrt{2\pi}\delta}\right)^n \exp\left(-\frac{1}{2\delta^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

ماکسیمیم تابع درست‌نمایی فوق را زمانی که  $\mu$  ثابت فرض شود برحسب  $\delta^2$  بدست می‌آوریم:

$$L_P(\mu) = \max_{\delta^2} \left[ \left(\frac{1}{\sqrt{2\pi}\delta}\right)^n \exp\left(-\frac{1}{2\delta^2} \sum_{i=1}^n (x_i - \mu)^2\right) \right]$$

ماکسیمیم فوق به ازای  $\delta^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$  رخ می‌دهد، پس داریم:

$$L_P(\mu) \propto \left( \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right)^{\frac{-n}{2}}$$

حال می‌توان استنباط شواهدی را بر حسب تابع درست‌نمایی نیم‌رخ فوق بدست آورد.

تعریف ۱. در بیماری سیستمیک فیبروزیست ( $CF$ ) یک آتوزمال را در نظر بگیرید که  $CFTR$  یک تنظیم کننده هدایت غشایی است. فرض کنید یک شرکت دارویی علاقمند به برآورد نسبت  $\theta$  برای افراد مبتلا به  $CF$  در یک کشور خاص است که حاملان ژنوتیپ هستند. شرکت فکر می‌کند این نسبت حدود ۵۰ درصد است. اما اگر بیشتر از این مقدار باشد، شرکت دارویی به دنبال توسعه درمانی هدفمند برای بیماری مورد نظر است. برای برآورد نسبت  $\theta$ ، شرکت دارویی می‌تواند از یک نمونه تصادفی از افراد جامعه استفاده کند و آنها را برای وجود جهش عامل ژنوتیپ آزمایش کند، سپس می‌تواند از روش‌های آماری مانند: نسبت درست‌نمایی یا عامل بیر برای برآورد نسبت  $\theta$  استفاده کند. بمنظور تعیین این که آیا نسبت افراد مبتلا به  $CF$  با نسبت حاملان ژنوتیپ

به‌طور قابل توجهی بیشتر از  $50^\circ$  درصد است یا خیر، شرکت دارویی می‌تواند از یک آزمون فرضیه استفاده کند بطوریکه:

فرضیه صفر ( $H_0: \theta \leq 50^\circ$ ) و فرضیه جایگزین ( $H_1: \theta > 50^\circ$ ) است. با تحلیل داده‌های جمع‌آوری شده و استفاده از آزمون‌های آماری می‌توان فرضیه صفر را پذیرفت یا رد کرد و بررسی کرد که آیا نسبت  $\theta$  به طور قابل توجهی بیشتر از  $50^\circ$  درصد است یا خیر. در صورتی که نسبت  $\theta$  بیشتر از  $50^\circ$  درصد باشد، شرکت دارویی می‌تواند به دنبال توسعه درمانی هدفمند برای جهش مورد نظر باشد.

## ۲.۱ تحلیل داده‌های ژنتیک در حضور پارامتر مزاحم توسط تابع درست‌نمایی نیمرخ

به عنوان مثال در نظر گرفتن مدل چند پارامتره  $L_n(\theta, \lambda)$  که پارامتر مورد علاقه  $\theta$  است و  $\lambda$  جز پارامترهای مزاحم محسوب می‌شود. درست‌نمایی شرطی و حاشیه‌ای همیشه در دسترس نیستند بنابراین همیشه می‌توان درست‌نمایی نیمرخ را محاسبه کرد. درست‌نمایی نیمرخ راه حل کلی برای نشان دادن قوت شواهد آماری برای پارامتر مورد علاقه در حضور پارامترهای مزاحم ارائه می‌دهد و در برنامه‌های کاربردی ژنتیک قابل استفاده است.

۱ در نمونه‌های بزرگ درست‌نمایی نیمرخ دو خصوصیت مهم درست‌نمایی را دارد بطوریکه در آن احتمال شواهد گمراه‌کننده توسط تابع کوهانی توصیف می‌شود.

۲ در نرم افزار  $R$  بسته  $PLikelihood$  برای محاسبه نسبت درست‌نمایی نیمرخ در بسیاری از مدل‌های آماری رایج: (مدل‌های خطی، خطی عمومی، ترکیبی و ترتیبی نسبتی) همراه با رسم نمودار تابع درست‌نمایی پارامتر مورد علاقه در دسترس است.

۳ و همچنین از بسته  $EVIAN$  نیز برای تجزیه و تحلیل رگرسیون خطی و لجستیکی مناسب برای مطالعات انجمن‌های ژنتیکی ارائه شده است.

به مثال شرکت دارویی بازمی‌گردیم، شرکت دارویی از نسبت درستنمایی نیمرخ استفاده کرده است تا قدرت شواهد آماری را برای ارتباط بین ایلئوس مکنونیوم و لوکوس  $SLC26A9$  در بیماری  $CF$  ارزیابی می‌کند. جمعیت مورد پژوهش شامل ۹۰۱ خواهر و برادر مبتلا به  $CF$  است که توسط انجمن بین المللی مدیفایر جمع آوری شده است. در اینجا از مدل رگرسیون لجستیک به عنوان مدل آماری استفاده شده است.

$$\log \frac{p_i}{1-p_i} = \beta_0 + \beta_1 G_{i1} + \gamma_2 Z_{i2} + \gamma_3 Z_{i3}$$

۱. که در آن  $p_i = E(y_i) = P(y_i = 1)$  و  $G_{i1} = 0, 1, 2$  است.

۲. برای مقدار ایلئوس مکنونیوم یک نتیجه دو دویی در نظر گرفته شده است، به این معنی که اگر فرد  $i$  ام با ایلئوس مکنونیوم دنیا آمده باشد  $y_i = 1$  در غیر این صورت  $y_i = 0$  است.

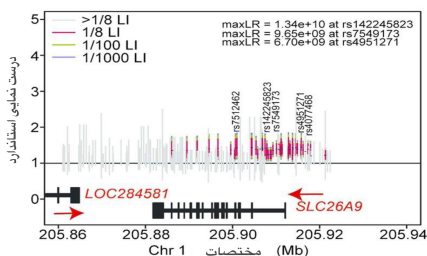
۳. متغیرهای  $\beta_0, \gamma_2, \gamma_3$  جزء پaramترهای مزاحم هستند و  $\beta_1$  پارامتر مورد علاقه است.

۴. متغیرهای  $Z_{i2}, G_{i1}, Z_{i3}$  به ترتیب به معنای تعداد آل‌های جایگزین در یک نکلوتید خاص در  $i$ -امین فرد بیمار شرکت کننده در اجرای تست ژنتیکی بر روی  $DNA$  آن‌ها است.

فرع ۱. طبق تعاریف گفته شده نمودار و جدول تحلیلی زیر را داریم:

۱. نمودار مقابل بازه شواهدی ارتباط بین ناحیه کروموزمی با ایلئوس مکنونیوم را نشان می‌دهد، که هر بازه خلاصه تابع درست‌نمایی برای هر  $SNP$  درج شده است.

۲. بازه‌های درست‌نمایی که با رنگ سبز و آبی مشخص شده‌اند ماکسیم نسبت درستنمایی نیمرخ را نشان می‌دهند، که نماینده‌های محدوده با داده‌های جمعیتی  $rs$  هستند که بازه درستنمایی باریکی را نشان می‌دهند که دور از خط  $OR = 1$  قرار دارند و دلیل قوی درباره ارتباط این منطقه ارائه می‌کند.



شکل ۱: تجزیه و تحلیل EP از لوکوس کروموزوم  $SLC26A9$  با ایلئوس مکنیوم در CF.

جدول ۱: خلاصه شبیه سازی داده های آماری برای تجزیه و تحلیل EP ساده در مقایسه با تجزیه و تحلیل

EP ساده برای نمونه‌های نامرتب و نسبتاً مرتبط CF

Robust factor ( $\frac{\sigma}{\beta}$ )	$\frac{1}{\sqrt{0.05}} LI$	$\frac{1}{\sqrt{0.01}} LI$	$\frac{1}{\lambda} LI$	MLE	maxLR	MAF	
<i>a. Unrelated n = ۵, ۸۶۹</i>							
NA	۱,۱۳۵۳, ۱,۶۵۳۰	۱,۱۷۶۳, ۱,۵۹۵۳	۱,۲۳۴۵, ۱,۵۱۶۳	۱,۳۶۹۹	۳۱۷, ۴۶۲, ۲۴۷	۰,۴۰۸۹	rs۷۵۱۲۴۶۲
NA	۱,۱۶۴۵, ۱,۶۸۲۷	۱,۲۰۶۶, ۱,۶۲۳۹	۱,۲۶۶۲, ۱,۵۴۷۴	۱,۳۹۸۰	۱۳, ۴۲۳, ۰۰۰, ۰۰۰	۰,۴۶۵۰	rs۱۴۲۲۴۵۸۲۳
NA	۱,۱۶۴۶, ۱,۶۷۴۴	۱,۲۰۳۷, ۱,۶۲۰۱	۱,۲۶۳۲, ۱,۵۴۳۸	۱,۳۹۸۲	۹, ۶۵۲, ۱۳۷, ۲۱۸	۰,۳۹۵۱	rs۷۵۴۹۱۷۳
NA	۱,۱۶۱۵, ۱,۶۸۱۷	۱,۲۰۴۵, ۱,۶۳۲۲	۱,۲۶۳۰, ۱,۵۵۱۴	۱,۳۹۸۰	۶, ۶۹۵, ۶۵۸, ۴۷۸	۰,۴۱۹۴	rs۴۹۵۱۲۷۱
NA	۱,۱۵۸۶, ۱,۶۸۱۷	۱,۲۰۰۵, ۱,۶۲۸۰	۱,۲۵۹۸, ۱,۵۴۷۴	۱,۳۹۴۵	۴, ۶۷۹, ۱۱۷, ۸۶۵	۰,۴۱۲۹	rs۴۰۷۷۴۶۸
<i>b. Related n = ۶, ۷۷۰</i>							
۰,۹۷۵	۱,۱۰۴۰, ۱,۵۸۱۲	۱,۱۴۱۰, ۱,۵۳۱۸	۱,۱۹۷۴, ۱,۴۵۹۷	۱,۳۲۲۱	۱۹, ۶۵۶, ۱۹۶	۰,۴۰۹۱	rs۷۵۱۲۴۶۲
۰,۹۶۶	۱,۱۳۲۴, ۱,۶۱۵۷	۱,۱۷۰۴, ۱,۵۶۳۲	۱,۲۲۸۲, ۱,۴۸۹۶	۱,۳۵۲۶	۵۰۲, ۳۱۸, ۲۸۲	۰,۴۶۵۲	rs۱۴۲۲۴۵۸۲۳
۰,۹۷۰	۱,۱۳۲۴, ۱,۶۰۷۸	۱,۱۷۰۶, ۰,۵۵۵۶	۱,۲۲۵۳, ۱,۴۸۶۱	۱,۳۴۹۴	۳۷۶, ۸۹۵, ۰۰۱	۰,۳۹۳۷	rs۷۵۴۹۱۷۳
۱,۰۰۴	۱,۱۲۶۶, ۱,۶۰۳۴	۱,۱۶۴۵, ۱,۵۵۱۴	۱,۲۲۲۰, ۱,۴۷۸۳	۱,۳۴۲۴	۲۷۸, ۹۵۸, ۸۰۸	۰,۴۱۹۳	rs۴۹۵۱۲۷۱
۰,۹۹۵	۱,۱۲۶۶, ۱,۶۰۷۵	۱,۱۶۴۵, ۱,۵۵۵۳	۱,۲۲۲۰, ۱,۴۸۲۱	۱,۳۴۵۸	۲۳۷, ۶۵۳, ۲۷۴	۰,۴۱۲۹	rs۴۰۷۷۴۶۸

$SNP$ ها با بیشترین  $maxLRs$  (ماکسیم نسبت درست‌نمایی نیم‌رخ) نمایش داده شده‌اند، و تحلیل‌های آماری برای واریانت‌های  $rs4077468$  و  $rs7512462$  که شواهد ارتباط قبلی با  $CF$  را نشان داده‌اند، نیز آمده است و فاکتور تنظیم قوی برای تجزیه و تحلیل نمونه مرتبط اعمال شده است.

$MAF$ : بیانگر فراوانی آلل‌های کمتر که همان مقادیر  $Z_{i2}$ ،  $G_{i1}$ ،  $Z_{i3}$  در مدل رگرسونی تعریف شده هستند.

$MLE$ : برآوردگر ماکسیم درست‌نمایی است.

## بحث و نتیجه‌گیری

رهیافت شواهدی یک روش آماری است که برای استخراج اطلاعات از داده‌های احتمالی استفاده می‌شود. این رویافت به دلیل خواص و ویژگی‌های خود قابلیت استفاده در مسائل پزشکی را دارد. به طور مثال: در آزمایشات بالینی برای تایید داروهای ژنتیک، رویافت شواهدی را می‌توان به عنوان یک روش برای ارزیابی یکسان بودن داروهای ژنتیک با داروهای اصلی استفاده کرد.

## مراجع

- [1] ROYALL, R. Statistical evidence a likelihood paradigm. CRC Press, New York, 2017.
- [2] EFRON, B. Why isn't everyone a bayesian. The American Statistician 40, 1 (1986), 1-5.
- [3] GUSTAFSON, P. Parameter restrictions for the sake of identification is there utility

in asserting that perhaps a restriction holds. *Statistical Science* 38, 3 (2023), 477–489.

- [4] Strug, L.J., 2018. The evidential statistical paradigm in genetics. *Genetic Epidemiology*, 42(7), pp.590607.



## **Profile likelihood evidence analysis of genetic data**

<sup>1</sup>mahdi Emadi, Ferdowsi University, Mashhad , Iran.

<sup>2</sup>Anis Sanchouli, Ferdowsi University, Mashhad , Iran.

<sup>3</sup>Morteza Mohammadi, Zabol University, Zabol , Iran.

---

**Abstract:** In the science of statistics, there are three schools of thought, the Bayesian school and the evidence school. In the frequentist school, only the observations and frequency of events are considered and problems can be solved according to that, while in the Bayesian school, in addition to the observations, information and creative beliefs of the researchers are important and they are involved in solving the problem and drawing conclusions.

**Keywords:** evidence paradigm, likelihood ratio, profile likelihood, logistic regression.

---