

## تقریب‌های قوی برای فرآیندهای چندکی بازنمونه‌گیری شده و کاربرد آن در روش‌شناسی *ROC*

سارا خزاعی<sup>۱</sup> - وحید فکور<sup>۲</sup>

<sup>۱</sup> دانشگاه آزاد اسلامی، واحد مشهد، دانشکده علوم پایه

<sup>۲</sup> دانشگاه فردوسی مشهد، دانشکده علوم ریاضی

**چکیده:** خم *ROC*<sup>۱</sup> از رسم مقادیر نرخ مثبت واقعی در مقابل نرخ مثبت کاذب حاصل می‌شود که به طور گسترده‌ای در علوم پزشکی برای توانایی در اندازه درستی آزمون-های تشخیصی استفاده می‌شود، از نظر ریاضی خم *ROC* ترکیب تابع بقای یک جامعه با تابع چندکی از جامعه دیگر است. در این مقاله ابتدا تقریب قوی برای فرآیندهای چندکی توزیع‌های بازنمونه‌گیری خودگردان‌ساز بیزی<sup>۲</sup> (*BB*) را مورد مطالعه قرار می‌دهیم و در ادامه با استفاده از این نتایج، تقریب‌های قوی برای *BB* فرآیند *ROC* در عبارت‌های با دو فرآیند کی فر مستقل را بررسی می‌کنیم (گو و گوسال، ۲۰۰۸).

**واژه‌های کلیدی:** تقریب قوی، خم *ROC*، خودگردان‌ساز بیزی، فرآیند چندکی، فرآیند کی فر

### ۱ مقدمه

خم *ROC*، نخستین بار در جنگ جهانی دوم برای تحلیل سیگنال‌های رادار استفاده شد و سپس در سال ۱۹۵۰ در رشتہ روانشناسی مورد استفاده قرار گرفت (گرین، ۱۹۶۶) و از سال ۱۹۷۰ به عنوان یک ابزار ضروری برای اندازه درستی آزمون‌های تشخیصی برای جداسازی افراد بیمار از سالم مطرح شد (متز، ۱۹۷۸). در مطالعات پزشکی که افراد به دو گروه سالم و بیمار طبقه‌بندی می‌شوند، ابتدا نقطه‌ی برش<sup>۳</sup> (*c*) را انتخاب کرده، افراد با نتایج آزمون بیشتر از *c* به عنوان بیمار در نظر گرفته می‌شوند، در غیر این صورت سالم هستند.

فرض می‌کنیم *X* دارای *T* تابع توزیع تجمعی *F* مربوط به جامعه‌ی افراد سالم و *Y* دارای

<sup>۱</sup> Receiver operating characteristic

<sup>۲</sup> Bayesian Bootstrap

<sup>۳</sup> Cut off point

تابع توزیع تجمعی  $G$  مربوط به جامعه افراد بیمار است. حساسیت<sup>۴</sup> یک آزمون را به صورت  $SE(c) = 1 - G(c)$  تعریف می‌کنیم که احتمال درست طبقه‌بندی شدن یک فرد بیمار است و ویژگی<sup>۵</sup> یک آزمون به صورت  $SP(c) = F(c)$  است که احتمال درست طبقه‌بندی شدن یک فرد سالم است (سیه و تارنیال، ۱۹۹۶). خم  $ROC$  نموداری از  $SE(c)$  یا نرخ مثبت واقعی روی محور عمودی در مقابل  $-1$  یا نرخ مثبت کاذب روی محور افقی است.

خم  $ROC$  را به صورت  $\{(P(X > c), P(Y > c)) : X \sim F, Y \sim G, c \in \mathbb{R}\}$  تعریف می‌کنیم، با فرض این که  $t = P(X > c)$  می‌توان تابع خم  $ROC$  را به صورت  $\bar{F}(x) = R(t) = \bar{G}(\bar{F}^{-1}(t))$  در نظر گرفت که  $1 < t < 0$  است،  $\bar{F}(x) = 1 - F(x)$  و  $(y) = 1 - G(y) = 1 - G(1 - F(x))$  تابع بقای متغیرهای مستقل  $X$  و  $Y$  هستند و مشتق آن به صورت  $R'(t) = g(\bar{F}^{-1}(t))/f(\bar{F}^{-1}(t))$  است.

از ویژگی‌های خم  $ROC$  می‌توان به موارد زیر اشاره کرد:

- ۱) یک خم  $ROC$  می‌تواند تمام اطلاعات مشخصه یک آزمون تشخیصی را در یک نمودار واحد نشان دهد که برای هدف‌های طبقه‌بندی شده مورد استفاده قرار می‌گیرد.
- ۲) آزمون‌های تشخیصی متفاوت توسط خم‌های  $ROC$  متناظرشان می‌توانند مقایسه شوند، بدین معنی که اگر یک خم  $ROC$  بالاتر از خم دیگر قرار گیرد، خم بالاتر را به عنوان خم  $ROC$  بهتر در نظر می‌گیریم (ویاند و همکاران، ۱۹۸۹). مقایسه‌ی خم‌های  $ROC$  در شکل (۱) نشان داده شده است. سطح زیر خم<sup>۶</sup> ( $AUC$ ) را می‌توان به صورت  $P(Y > X)$  تعریف کرد که بیانگر دقیقت تشخیص آزمون می‌باشد (بامبر، ۱۹۷۵).

از جمله روش‌های محاسبه‌ی  $AUC$  می‌توان به استفاده از قاعده ذوزنقه‌ای، برآزandن داده‌ها به مدل تبدیل دو نمونه‌ای با استفاده از برآورد ماکزیمم درستنمایی و استفاده از آماره‌ی من-وینتی نیز اشاره کرد.

برآورده‌گر ناپارامتری  $ROC$  توسط جایگزینی در همتای تجربی آن به دست می‌آید و تغییرپذیری آن با استفاده از خودگردان‌ساز<sup>۷</sup> برآورده شود (پپ، ۲۰۰۳)، که در ادامه به آن می‌پردازیم.

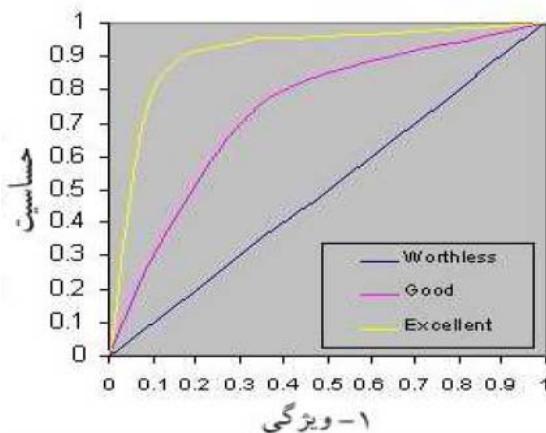
در این مقاله در بخش دوم مروری بر فرآیند چندکمی خواهیم داشت، در بخش سوم خودگردان‌ساز بیزی را معرفی می‌کنیم و در بخش پایانی تقریب‌های قوی گوسی برای

<sup>۴</sup> Sensitivity

<sup>۵</sup> Specificity

<sup>۶</sup> Area under curve

<sup>۷</sup> Bootstrap



شکل ۱: مقایسه خم‌های ROC

فرآیندهای چندکی  $BB$  و خم  $ROC$  را مورد مطالعه قرار می‌دهیم.

## ۲ فرآیند چندکی

ابتدا فرآیند تجربی را تعریف کرده و سپس فرآیند چندکی را معرفی می‌کنیم.

تعریف ۱۳ فرض کنید  $\{X_i, i \geq 1\}$  دنباله‌ای از متغیرهای تصادفی مستقل و هم توزیع، با توزیع مشترک  $F$  باشد، تابع توزیع تجربی این متغیرها عبارت است از:

$$\mathbb{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$$

که در آن  $I(\cdot)$  تابع نشانگر است.

اگر  $X_{1:n} \leq X_{2:n} \leq \dots \leq X_{n:n}$  آماره‌های ترتیبی نمونه تصادفی  $X_1, \dots, X_n$  باشد، آن‌گاه تابع توزیع تجربی را به صورت زیر نشان می‌دهیم

$$\mathbb{F}_n(x) = \begin{cases} 0 & x < X_{1:n} \\ \frac{j}{n} & X_{j:n} \leq x < X_{j+1:n}, j = 1, \dots, n-1 \\ 1 & X_{n:n} \leq x \end{cases}$$

فرآیند تجربی را به شکل زیر تعریف می کنیم:

$$\mathbb{J}_n(x) = \sqrt{n}(\mathbb{F}_n(x) - F(x))$$

تعریف ۱۴ فرض کنید  $X_1, \dots, X_n$  نمونه‌ی تصادفی از توزیع  $F$  باشند، آن‌گاه تابع چندک به صورت زیر تعریف می‌شود:

$$F^{-1}(y) = \inf\{x ; F(x) \geq y\}; \quad 0 < y < 1$$

و برآوردگر نمونه‌ای آن برابر است با:

$$\mathbb{F}_n^{-1}(y) = \inf\{x ; \mathbb{F}_n(x) \geq y\}; \quad 0 < y < 1$$

و یا

$$\mathbb{F}_n^{-1}(y) = X_{K:n}, \quad \frac{K-1}{n} < y \leq \frac{K}{n}, \quad K = 1, 2, \dots, n$$

فرآیند چندکی را به صورت زیر تعریف می کنیم:

$$\mathbb{Q}_n(y) = \sqrt{n}(\mathbb{F}_n^{-1}(y) - F^{-1}(y))$$

### ۳ خودگردان‌ساز بیزی (BB)

روش خودگردان‌ساز توسط افران در سال ۱۹۷۹ ارائه شد. خودگردان‌ساز یک روش برای برآورد واریانس و نیز یافتن توزیع آماره‌ها می‌باشد. همچنین می‌توانیم از خودگردان‌ساز برای ساختن بازه‌های اطمینان استفاده کنیم (واسمن، ۲۰۰۶). پس از روش خودگردان‌ساز، رابین در سال ۱۹۸۱ روش خودگردان‌ساز بیزی را معرفی کرد. فرض کنید  $X_1, \dots, X_n$  نمونه‌ی تصادفی از توزیع  $F$  باشند، در روش خودگردان‌ساز نمونه‌ی تصادفی  $X_1^*, \dots, X_n^*$  را به روش نمونه‌گیری با جایگذاری از نمونه‌ی  $X_1, \dots, X_n$  انتخاب می‌کنیم که هر مشاهده‌ی  $X_i^*$  با احتمال  $\frac{1}{n}$  انتخاب می‌شود. اما در روش BB به جای انتخاب با احتمال  $\frac{1}{n}$ ، یک توزیع احتمال پسین برای  $X_i$  ها به صورت زیر در نظر می‌گیریم:

ابتدا  $n-1$  متغیر تصادفی یکنواخت از بازه‌ی  $[0, 1]$  را در نظر می‌گیریم.  $V_{1:n-1}, \dots, V_{n-1:n-1}$  را به عنوان آماره‌های ترتیبی  $V_1, \dots, V_{n-1}$  فرض می‌کنیم که  $V_{n:n-1} = 1$  هستند. آن‌گاه  $\Delta_{j:n} = V_{j:n-1} - V_{j-1:n-1}$  را نمودار آماره‌های ترتیبی یکنواخت معرفی

می‌کنیم.  
 $\Delta$  برای تعیین احتمالات نمونه‌ی  $BB$  استفاده می‌شود (رابین، ۱۹۸۱).

با توجه به نمونه‌ی  $X_1, \dots, X_n$ ،تابع توزیع تجمعی  $BB$  به صورت  $\mathbb{F}_n^\#(x) = \sum_{1 \leq j \leq n} \Delta_{j:n} I(X_{j:n} \leq x)$  است.

به سادگی تابع چندکی  $BB$  به صورت زیر به دست می‌آید:

$$\mathbb{F}_n^{\#-\frac{1}{2}}(y) = \begin{cases} X_{j:n} & V_{j-1:n-1} < y \leq V_{j:n-1}, j = 1, \dots, n \\ X_{\circ:n} & y = \circ \end{cases}$$

فرآیند تجربی  $BB$  و فرآیند چندکی  $BB$  را به صورت زیر تعریف می‌کنیم:  
 $\text{فرآیند تجربی } BB : \mathbb{J}_n^\#(x) = \sqrt{n}(\mathbb{F}_n^\#(x) - \mathbb{F}_n(x))$

فرآیند چندکی  $BB$

$$\mathbb{Q}_n^\#(y) = \sqrt{n}(\mathbb{F}_n^{\#-\frac{1}{2}}(y) - \mathbb{F}_n^{-\frac{1}{2}}(y))$$

## ۴ تقریب‌های قوی گوسی برای فرآیندهای چندکی $BB$ و $ROC$

فرض می‌کنیم  $X_1, \dots, X_n$  نمونه‌ی تصادفی از توزیع  $F$  باشند و دامنه‌ی  $X$  را به صورت  $[a, b]$  تعریف می‌کنیم که

$$b = \inf\{x : F(x) = 1\}, \quad a = \sup\{x : F(x) = 0\}$$

است و  $l_n = n^{-\frac{1}{4}}(\log \log n)^{\frac{1}{4}}(\log n)^{\frac{1}{4}}$  را در نظر می‌گیریم.  
 شرایط  $A$  (چورگو و روز، ۱۹۷۸) و  $B$  را فرض می‌کنیم که برای یک کلاس بزرگ از تابع‌های توزیع برقرار هستند.

شرط  $A$ : تابع توزیع پیوسته  $F$  را در نظر می‌گیریم که مشتق مرتبه دوم آن روی  $(a, b)$  وجود دارد و  $F' = f \neq 0$  روی  $(a, b)$  است و برای بعضی  $\gamma > 0$  داریم:

$$\sup_{a < x < b} F(x)(1 - F(x))|f'(x)/f''(x)| \leq \gamma.$$

شرط  $F$  و  $G$  تحت شرط  $A$  برقرارند و شرایط زیر را نیز در نظر می‌گیریم:

$$\sup_{a < x < b} F(x)(1 - F(x)) \left| \frac{g'(x)}{f^*(x)} \right| < \infty$$

$$\sup_{a < x < b} F(x)(1 - F(x)) \left| \frac{g(x)}{f(x)} \right| < \infty$$

تقریب قوی دارای بیشینه‌ای طولانی می‌باشد، دستاوردهای بیشماری در این زمینه موجود است. قضایای زیر برخی از نتایج در این زمینه را نشان می‌دهند. کومولوش و همکاران (۱۹۷۵) نشان دادند که می‌توان  $\mathbb{J}$  را توسط دنباله‌ای از پلهای براونی<sup>۸</sup> تقریب زد. چورگو و روز (۱۹۷۸) ثابت کردند که  $\mathbb{Q}_n$  را می‌توان توسط یک فرآیند کی فر<sup>۹</sup> تقریب زد. لو (۱۹۸۷) قضیه‌ی تقریب قوی برای تابع توزیع تجمعی فرآیندهای خودگردان‌ساز و  $BB$  را مورد مطالعه قرار داد. در ادامه ابتدا فرآیند کی فر را تعریف کرده و در ادامه به بررسی قضایای تقریب قوی از (گو و گوسال، ۲۰۰۸) می‌پردازیم.

**تعریف ۱۵** فرآیند  $\{K(x, y) : 0 \leq x \leq 1, 0 \leq y < \infty\}$  یک فرآیند کی فر نامیده می‌شود هرگاه بتوان آن را به صورت  $K(x, y) = W(x, y) - xW(1, y)$  نوشت که  $W(x, y)$  یک فرآیند وینر دوپارامتری است.

لم ۱ فرض کنید  $X_1, \dots, X_n$  نمونه‌ی تصادفی از توزیع  $F$ ، تحت شرط  $A$  باشند،  $V_1, \dots, V_{n-1}$  نمونه‌ی تصادفی از توزیع  $(1, 0)$  و مستقل از  $X_i$  هاستند و  $\tilde{U}_{n-1}(x)$  را به عنوان تابع تجربی  $V_1, \dots, V_{n-1}$  در نظر می‌گیریم. آن‌گاه:

$$\sup_{0 < y < 1} \sqrt{n} \left| \mathbb{F}_n^{\#-1}(y) - \mathbb{F}_n^{-1}(\tilde{U}_{n-1}(y)) \right| = O(n^{-\frac{1}{4}} \log n), \text{ a.s.} \quad (1)$$

علاوه بر آن، یک فرآیند کی فر  $K$  وجود دارد به طوری که:

$$\sup_{\delta_n^\# \leq y \leq 1 - \delta_n^\#} \left| \sqrt{n} f(F^{-1}(y)) (\mathbb{F}_n^{-1}(\tilde{U}_{n-1}(y)) - \mathbb{F}_n^{-1}(y)) - n^{-\frac{1}{4}} K(y, n) \right| = O(t_n), \text{ a.s.} \quad (2)$$

□

<sup>۸</sup> Brownian Bridge

<sup>۹</sup> Kiefer Process

قضیه ۱۲ فرض کنید  $X_1, \dots, X_n$  نمونه‌ی تصادفی از توزیع  $F$  تحت شرط  $A$  باشند. آن‌گاه  $(y) \mathbb{Q}_n^\#$  را می‌توان توسط یک فرآیند کی فر  $K$  تقریب زد،

$$\sup_{\delta_n^\# \leq y \leq 1 - \delta_n^\#} |f(F^{-1}(y)) \mathbb{Q}_n^\#(y) - n^{-\frac{1}{\tau}} K(y, n)| = O(l_n), \text{ a.s.} \quad (3)$$

$$\delta_n = 25n^{-1} \log \log n \text{ و } \delta_n^\# = \delta_n + n^{-\frac{1}{\tau}} (\log \log n)^{\frac{1}{\tau}}$$

□

قضیه ۱۳ فرض کنید  $X_1, \dots, X_m$  نمونه‌ی تصادفی از توزیع  $F$  و  $Y_1, \dots, Y_n$  نمونه‌ی تصادفی از توزیع  $G$ ، تحت شرایط  $A$  و  $B$  باشند، آن‌گاه

$$\mathbb{R}(t) = R(t) + R'(t) \frac{K_1(t, m)}{m} + \frac{K_2(R(t), n)}{n} + O(\alpha_m^{-1} \tau_m), \quad t \in (\alpha_m, 1 - \alpha_m) \quad (4)$$

$$\mathbb{R}^\#(t) = \mathbb{R}(t) + R'(t) \frac{K_1(t, m)}{m} + \frac{K_2(R(t), n)}{n} + O(\alpha_m^{-1} \tau_m), \quad t \in (\alpha_m^\#, 1 - \alpha_m^\#) \quad (5)$$

که  $K_1$  و  $K_2$  فرآیندهای کی فر مستقل هستند و

$$\mathbb{R}(t) = \mathbb{R}_{m,n}(t) = \overline{\mathbb{G}}_n(\overline{\mathbb{F}}_m^{-1}(t)),$$

$$\mathbb{R}^\#(t) = \mathbb{R}_{m,n}^\#(t) = \overline{\mathbb{G}}_n^\#(\overline{\mathbb{F}}_m^{\# -1}(t)),$$

$$\alpha_m^{-1} = O\left(m^{\frac{1}{4-\epsilon}}\right) \text{ و } \alpha_m^\# = \alpha_m + (m-1)^{-\frac{1}{\tau}} (\log \log(m-1))^{\frac{1}{\tau}} \tau_m = \frac{l_m}{\sqrt{m}} \text{ را در نظر می‌گیریم.}$$

□

تذکر ۱. نتیجه‌ی مشابه (۴) در (سیه و تارنیال، ۱۹۹۶) اثبات شده است. تفاوت بین این دو یک موازن‌های متفاوت بین بازه‌ی اعتبار و نرخ خطای تقریب است.

تذکر ۲. تقریب قوی برای فرآیند چندکی خودگردان‌ساز و تقریب قوی برای برآوردهای خودگردان‌ساز  $ROC$  را می‌توان مشابهًا با همین نرخ خطای به دست آورد اما با دامنه تغییرات متمایز  $t$ .

## بحث و نتیجه‌گیری

در این مقاله، ابتدا به بررسی خم  $ROC$  پرداخته که عملکرد یک آزمون تشخیصی را توصیف می‌کند که نموداری از حساسیت در مقابل ۱ – ویژگی است. سپس فرآیندهای چندکی  $BB$  را توسط یک فرآیند کی فر با نرخ همگرایی  $O(l_n)$  تقریب زده و در ادامه با استفاده از این نتایج، تقریب‌های قوی برای  $BB$  فرآیند  $ROC$ ، مرتبط با دو فرآیند کی فر مستقل را مور بررسی قرار دادیم.

## مراجع

- Bamber, D. (1975), The area above the ordinal dominance graph and the area below the receiver operating characteristic graph, *J. Math. Psychol.* 12, pp. 387-415.
- Csőrgő ,M. and Révész, P. (1978) , Strong approximations of the quantile process, *Ann. Statist.* 6, pp. 882-894.
- Green, D.M. and Swets. J.A.(1966), Signal Detection Theory and Psychophysics, JohnWiley & Sons, NewYork.
- Gu, J. and Ghosal ,S.(2008), Strong approximations for resample quantile processes and application to ROC methodology. *Journal of Nonparametric Statistics*, 20:3, 229-240.
- Hsieh, F.S. and Turnbull, B.W.(1996) , Nonparametric and semiparametric estimation of the receiver operating characteristic curve, *Ann. Statist.* 24 ,pp. 25-40.
- Komlós, J. and Major,.P. and Tusnády ,G .(1975) , An approximation of partial sums of independent RV's and the sample DF. I., *Z. Wahrsch. Verw. Gebiete*, 32, pp. 111-131.
- Lo, A.Y ,(1987) A large sample study of the Bayesian bootstrap, *Ann. Statist.* 15, pp. 360-375.
- Metz, C.E.(1978) , Basic principles of ROC analysis, *Seminars Nuc. Med.* 8, pp. 283-298.
- Pepe, M.S.(2003), The Statistical Evaluation of Medical Tests for Classification and Prediction, Oxford University Press, Oxford.

Rubin, DB.(1981), The Bayesian bootstrap. The Annals of Statistics, 9:130-134.

Wasserman, L.(2006), All of Nonparametric Statistics, springer.

Wieand ,S.and Gail, M.H. and James, B.R. and James, K.L.(1989), A family of nonparametric statistics for comparing diagnostic markers with paired or unpaired data, Biometrika 76 ,pp. 585-592.