

بررسی آمارگان رن‌ها و بکارگیری آن در تعدادی تست آماری

محمد دخیل‌علیان

استادیار دانشکده برق و کامپیوتر

دانشگاه صنعتی اصفهان

mdalian@cc.iut.ac.ir

چکیده

آزمونهای آماری ابزار مناسبی برای بررسی خواص دنباله‌های شبه تصادفی خصوصاً در سیستم‌های رمزنگاری می‌باشد. در این راستا آزمونهای متنوعی ارائه شده‌اند که از جمله آنها می‌توان به آزمون رن‌ها^۱ اشاره نمود. در این مقاله ضمن اشاره به آزمون متداول برای بررسی رن‌های موجود در یک دنباله، با بررسی تئوری آمارگان رن‌ها در یک دنباله باینری تصادفی ایده‌آل، روشی برای آزمون رن‌ها و علاوه بر آن آزمونهایی برای گپ‌ها^۲ و بلوک‌های^۳ ارائه شده است.

کلمات کلیدی :

آزمونهای آماری، آزمون رن‌ها، دنباله‌های شبه تصادفی، سیستمهای رمزنگاری.

۱- مقدمه

در بسیاری از کاربردها، استفاده از دنباله‌های کاملاً تصادفی از اهمیت قابل توجهی برخوردار می‌باشد. منظور از یک دنباله تصادفی، دنباله‌ای با عناصر مستقل و دارای توزیع یکسان و یکنواخت (*i.i.d*) می‌باشد، لذا مولدهایی مفیدتر می‌باشند که از لحاظ آماری دنباله‌هایی را با ویژگی‌های دنباله‌های تصادفی تولید نمایند. با شناخت خواص آماری دنباله‌های تصادفی، معیارهایی برای دنباله‌های شبه تصادفی بیان شده است [1] و متعاقب آن نیز آزمونهای مختلفی برای بررسی رفتار تصادفی دنباله‌ها و مولدهایشان طرح شده‌اند. آزمونهای آماری بر روی بخش کوچکی از دنباله انجام می‌شود و طی آن، دنباله در آزمون قبول یا رد می‌گردد. برای ارزیابی یک مولد معمولاً تعداد زیادی دنباله مورد آزمایش قرار می‌گیرد و لازم است تعداد معنی داری از دنباله‌ها، از آزمون عبور نمایند. ذکر این نکته ضروری است که قبول یا رد شدن دنباله‌ای در یک آزمون دلیل بر تصادفی بودن یک دنباله نیست، زیرا در این آزمونها رفتار نوعی دنباله‌های تصادفی مد نظر قرار می‌گیرد و بدین جهت دنباله‌هایی که با این رفتار نوعی مطابقت داشته باشند، مورد قبول واقع می‌گردند. آزمون رن‌ها [۲-۳] از جمله آزمونهایی است که در این مقاله مورد بررسی قرار می‌گیرد و پیرو آن آزمونهای دیگری در این خصوص ارائه می‌گردد.

¹ Run.

² Gap.

³ Block.

۲- آزمون رن‌ها

در دنباله‌های کاملاً تصادفی باینری بیت‌های صفر و یک بدون نظم و شکل معینی رخ می‌دهند و بدین لحاظ مدلسازی این دنباله‌ها برای تحقق یک آزمون آماری ضروری است. یکی از مشخصه‌های یک دنباله باینری تعداد و طول رن‌های موجود در یک دنباله می‌باشد. به تعدادی بیت صفر و یا بیت یک که پشت سرهم در دنباله قرار گیرند، یک رن گفته می‌شود. اگر رن شامل تعدادی یک باشد به آن بلوک و اگر رن شامل تعدادی صفر باشد به آن گپ گفته می‌شود. به عنوان نمونه دنبالهٔ 010001100001111110 شامل هفت رن به طول‌های ۱، ۲، ۳، ۴ و ۶ می‌باشد که سه تای آن بلوک به طول‌های ۱ و ۲ و ۶ و چهارتای آن گپ و به طول‌های ۱ و ۳ و ۴ می‌باشد. در یک دنباله باینری کاملاً تصادفی نظیر $X^n = X_1, X_2, \dots, X_n$ ، تعداد رن‌ها در حالت کلی دارای توزیع مشخصی می‌باشند. لذا میانگین و واریانس آنها به عنوان معیاری در آزمون ارایه شده در [۳] در نظر گرفته شده است. به عبارت دیگر ثابت شده است که در یک دنباله n بیتی (هنگامی که n به سمت بی‌نهایت میل کند)، تعداد رن‌های موجود در یک دنباله دارای توزیع نرمال با میانگین و واریانس زیر می‌باشد [۴]:

$$m = 1 + \frac{2n_0n_1}{n} \quad \text{متوسط تعداد رن‌ها} \quad (1)$$

$$\sigma^2 = \frac{(m-1)(m-2)}{n-1} \quad \text{واریانس تعداد رن‌ها}$$

که در رابطهٔ فوق n_0 ، n_1 و n به ترتیب تعداد صفرها، یکها و طول دنباله می‌باشد. در مرجع [۳] نیز به گونه‌ای دیگر رابطه (۱)، مورد استفاده قرار گرفته است. بدین ترتیب پارامتر آزمون مربع کای^۱ برای بررسی توزیع آماری رن‌های موجود در یک دنباله به صورت زیر ارایه شده است.

$$\chi_{run}^2 = \left(\frac{N_R - m}{\sigma} \right)^2 \quad (2)$$

که در عبارت فوق N_R تعداد کل رن‌ها در یک دنبالهٔ نوعی مورد آزمایش با طول n بیت می‌باشد. آزمون فوق بواسطهٔ اینکه تنها متوسط و واریانس تعداد کل رن‌ها را مورد آزمون قرار می‌دهد، نمی‌تواند آزمون کاملی برای بررسی آماری تعداد رن‌ها باشد. در بخش بعد ضمن بیان تئوری توزیع رن‌ها در یک دنبالهٔ کاملاً تصادفی چگونگی اصلاح آزمون رن‌ها را بیان می‌نماییم.

۳- محاسبهٔ تابع احتمال رن‌ها

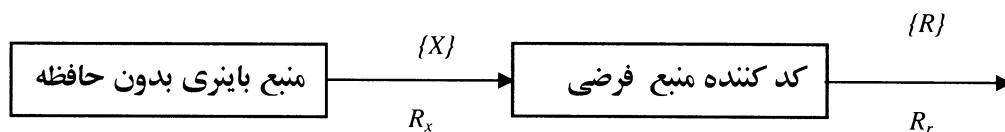
فرض کنید منبع S یک منبع باینری بدون حافظه باشد که بیت‌های صفر و یک را با احتمال $\frac{1}{2}$ تولید کند، بنابراین آنروپی منبع برابر یک بیت بر سمبل بوده و لذا سرعت اطلاعات این منبع برابر بیت^۲ منبع می‌باشد. حال اگر در خروجی این منبع از کد کنندهٔ منبع^۳ فرضی شکل (۱) با مجموعه الفباء $r_1^0, r_1^1, r_2^0, \dots$ به صورت جدول (۱) استفاده نماییم، کد کنندهٔ منبع دارای مجموعه الفبایی با بی‌نهایت عضو می‌باشد. توجه نمایید که علی‌رغم اینکه پیشوند بودن کلمات کد در اینجا مراعات نشده است، منتها فرض بر این است که هیچگاه دو یا چند سمبل یکسان کنار هم کد

¹

² Bit rate.

³ Source encoder.

نشوند. به عنوان نمونه اگر دنباله ورودی برابر 011000 بود، دنباله خروجی تنها به $r_1^0 r_2^1 r_3^0$ کد گردد و هیچگاه دنباله خروجی به یکی از شکل‌های $r_1^0 r_1^1 r_1^0 r_2^0$ یا $r_1^0 r_1^1 r_1^0 r_2^1$ و ... تبدیل نگردد.



شکل (۱)

جدول (۱): نگاشت خروجی منبع فرضی

ورودی	سمبل‌های خروجی
0	r_1^0
1	r_1^1
00	r_2^0
11	r_2^1
000	r_3^0
111	r_3^1
⋮	⋮

به این ترتیب سرعت اطلاعات در خروجی کد کننده منبع برابر R_x بیت بر ثانیه می‌گردد و با توجه به اینکه تعداد بیت در نظر گرفته برای هر سمبل برابر تعداد بیت ورودی است، لذا با توجه به اینکه بیت‌های خروجی منبع باینری یکنواخت و مستقل می‌باشند ($H(X)=1$)، لذا اگر سمبل‌های خروجی کد کننده منبع بخواهد وابسته به هم باشند، بالاجبار سرعت اطلاعات خروجی کمتر از سرعت اطلاعات ورودی می‌گردد ($R_r < R_x$) و این به معنی گم شدن اطلاعات یا ابهام در دکدینگ می‌باشد که با فرض بیان شده، این مسأله متفی است، پس دنباله خروجی کد کننده منبع، دنباله‌ای با عناصر مستقل از هم می‌باشد. حال می‌خواهیم توزیع خروجی کد کننده منبع را محاسبه نماییم. فرض کنید $\{X\}$ دنباله تولید شده توسط منبع باینری باشد.

$$\{X\} = X_1, X_2, X_3, \dots \quad (3)$$

برای تولید سمبل r_j^1 یا r_j^0 در خروجی کد کننده منبع، از ابتدای دنباله شروع می‌کنیم و لذا به سادگی می‌توان احتمال تولید سمبل‌های r_j^0 و r_j^1 را بدست آورد:

$$P(R = r_j^1) = P(X_i = 1, X_{i+1} = 1, \dots, X_{i+j-1} = 1, X_{i+j} = 0) \quad (4)$$

با توجه به استقلال X_i ها از یکدیگر داریم:

$$P(R = r_j^1) = \left(\prod_{k=i}^{i+j-1} P(X_k = 1) \right) \cdot P(X_{i+j} = 0) = \left(\frac{1}{2} \right)^{j+1} \quad (5)$$

به همین ترتیب می توان احتمال تولید سمبل r_j^0 را محاسبه نمود :

$$P(R = r_j^0) = P(X_i = 0, X_{i+1} = 0, \dots, X_{i+j-1} = 0, X_{i+j} = 1) = \left(\frac{1}{2}\right)^{j+1} \quad (6)$$

لازم به ذکر است که در روابط فوق $X_{i+j}=0$ در نظر گرفته شده است، زیرا بنابر فرض هیچگاه دو یا چند سمبل یکسان در خروجی کد کننده مجاز نمی باشد و $X_{i+j}=0$ تضمین کننده این شرط می باشد. با توجه به روابط (5) و (6) می توان نتیجه گرفت، خروجی کد کننده منبع فرضی دنباله ای به صورت $\{R\}$ با سمبل های مستقل از هم می باشد و هر سمبل دارای تابع احتمال به صورت (7) می باشد

$$\{R\} = R_1, R_2, R_3, \dots \quad (7)$$

$$P(R_i = r_j^t) = \left(\frac{1}{2}\right)^{j+1}$$

$$i=1, 2, \dots \quad t=0, 1$$

واضح است که :

$$\sum_{\forall j,t} P(R_i = r_j^t) = \sum_{t=0}^1 \sum_{j=1}^{\infty} \left(\frac{1}{2}\right)^{j+1} = 2 \times \left(\frac{1}{4} + \frac{1}{8} + \dots\right) = 1 \quad (8)$$

با این مدلسازی در واقع r_j^0 گپ به طول z و r_j^1 بلوک به طول z را در دنباله $\{X\}$ مشخص می کند. به عبارت دیگر رابطه (7) تابع احتمال گپها و بلوکها را در یک دنباله باینری کاملاً تصادفی مشخص کرده است. به سادگی می توان تابع احتمال رنها را نیز بدست آورد.

$$P(R = r_j) = P(R = r_j^0 \text{ or } R = r_j^1) = \left(\frac{1}{2}\right)^{j+1} + \left(\frac{1}{2}\right)^{j+1} = \left(\frac{1}{2}\right)^j \quad (9)$$

به عبارت دیگر احتمال وقوع یک رن به طول z (R_j) برابر است با مجموع احتمال رخ دادن یک گپ به طول z یا بلوک به طول z . برای سادگی تابع احتمال برای رنها را به صورت زیر می نویسیم :

$$P(R=R_j) = P(R=j) = P_R(j)$$

$$j=1, 2, \dots$$

بنابر این متوسط طول هر رن را می توان به صورت زیر محاسبه نمود :

$$\bar{R} = \sum_{j=1}^{\infty} jP(R=j) = \sum_{j=1}^{\infty} j\left(\frac{1}{2}\right)^j = 2(\text{bit}) \quad (10)$$

رابطه فوق نشان می دهد که در یک دنباله باینری کاملاً تصادفی به طول n ، به طور متوسط $\frac{n}{2}$ رن وجود دارد که با توجه به روابط 5 الی 7، به طور متوسط نصف رنها $\left(\frac{n}{4}\right)$ گپ و نصف دیگر بلوک می باشد.

۴- چگونگی انجام آزمون رنها

برای انجام آزمون رنها روی یک دنباله باینری به طول n بیت نظیر X^n ، دنباله را به دنباله R^k تبدیل می کنیم.

$$X^n = X_1, X_2, \dots, X_n \quad (11)$$

$$R^k = R_1, R_2, \dots, R_k$$

k با توجه به رابطه 10 به طور متوسط باید حدود $\frac{n}{2}$ باشد. حال فرض کنید آزمون را برای رنهای به طول

حداکثر m بخواهیم انجام دهیم. برای اینکار رنهای به طول بزرگتر از m را به عنوان وقوع یک سمبل فرضی * در نظر می‌گیریم که احتمال وقوع آن با توجه به رابطه (۹) برابر است با:

$$P(R_i = *) = P(R_i > m) = 1 - P(R_i \leq m) \\ = 1 - \left(\frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^m} \right) = \frac{1}{2^m} \quad (12)$$

بنابر این در مجموع می‌توان گفت:

$$P(R_i = j) = \begin{cases} \left(\frac{1}{2}\right)^j & j = 1, 2, \dots, m \\ \frac{1}{2^m} & j = * \end{cases} \quad i = 1, 2, \dots, k \quad (13)$$

بنابراین می‌توان گفت دنباله R^k یک متغیر تصادفی چند جمله‌ای با تابع احتمال رابطه (۱۳) می‌باشد. بنابراین طبق قضیه آزمون مربع کای [۵]، پارامتر آزمون رنها را می‌توان به صورت زیر نوشت:

$$\chi_r^2 = \sum_{i=1}^m \frac{(r_i - k \cdot 2^{-i})^2}{k \cdot 2^{-i}} + \frac{(r^* - k \cdot 2^{-m})^2}{k \cdot 2^{-m}} \quad (14)$$

درجه آزادی متغیر تصادفی χ_r^2 برابر m می‌باشد و شرط انجام آزمون با توجه به تقریب توزیع مربع کای به صورت زیر می‌باشد:

$$\frac{n}{2} \times \frac{1}{2^m} > 5 \quad \text{یا} \quad n > 5 \times 2^{m+1} \quad (15)$$

یا به طور معادل اگر طول دنباله مورد آزمون برابر n باشد، مقدار m باید در شرط زیر صدق نماید:

$$m < \log_2 \left(\frac{n}{10} \right) \quad (16)$$

(در رابطه (۱۴)، r^* تعداد رنهای به طول بیش از m در دنباله R^k می‌باشد)

حال اگر آزمون به این صورت انجام شود که تنها رنهای تا طول حداکثر m مورد آزمون قرار گیرند و بخواهیم سمبل فرضی * را در نظر نگیریم:

$$P(R_i = j / j \leq m) = \frac{P(R_i = j, j \leq m)}{P(j \leq m)} = \frac{P(R_i = j)}{\frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^m}} = \frac{2^{-j}}{1 - 2^{-m}} \quad (17)$$

به این ترتیب دنباله باینری X^n به دنباله R^L تبدیل می‌نماییم که تنها شامل رنهای به طول حداکثر m با تابع احتمال (۱۷) می‌باشد و واضح است که:

$$R^L = R_1, R_2, \dots, R_L \\ \sum_{j=1}^m P(R_i = j / j \leq m) = \sum_{j=1}^m \frac{2^{-j}}{1 - 2^{-m}} = 1 \quad (18)$$

به عبارت دیگر در دنباله R^L هر یک از R_i ها دارای توزیع زیر می‌باشند:

$$P(R_i = j) = \frac{2^{-j}}{1 - 2^{-m}} \quad (19)$$

$$j = 1, 2, \dots, m$$

بنابراین پارامتر آزمون رن‌ها را می‌توان در این وضعیت به صورت زیر نوشت :

$$\chi_r^2 = \sum_{i=1}^m \frac{(r_i - L \frac{2^{-i}}{1 - 2^{-m}})^2}{L \times \frac{2^{-i}}{1 - 2^{-m}}} \quad (20)$$

درجه آزادی متغیر تصادفی χ_r^2 برابر $m-1$ می‌باشد. در حالت ایده‌آل به طور متوسط L نیز برابر است با :

$$E(L) = \frac{n}{2} \times \sum_{j=1}^m P(R_i = j) = \frac{n}{2} \left(\frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^m} \right) = \frac{n}{2} (1 - 2^{-m}) \quad (21)$$

بدین ترتیب شرط انجام آزمون به صورت زیر بدست می‌آید که در واقع همان رابطه (۱۵) می‌باشد. بدین ترتیب

m باید در شرط رابطه (۱۶) نیز صدق نماید.

$$\frac{n}{2} (1 - 2^{-m}) \times \frac{2^{-m}}{(1 - 2^{-5})} > 5 \Rightarrow n > 5 \times 2^{m+1} \quad (22)$$

۵-آزمون‌های گپ و بلوک

اگر بخواهیم آزمونی برای گپها و آزمون دیگری به طور مجزا برای بلوکها ترتیب دهیم. به سادگی از روند ارایه شده در بخش ۴ می‌توان استفاده نمود: به این ترتیب که اگر دنباله n بیتی X^n توسط کد کننده منبع کد گردد و تنها گپ‌های موجود در خروجی مد نظر قرار گیرد (یا به طور معادل بلوکها) :

$$X^n = X_1, X_2, \dots, X_n \quad (23)$$

$$G^k = G_1, G_2, \dots, G_k$$

با توجه به نتیجه رابطه (۱۰)، مقدار k به طور متوسط باید حدود $\frac{n}{4}$ باشد. اگر فرض کنیم منظور از $P(G_i=j)$

یعنی احتمال اینکه گپ i ام به طول j باشد، با توجه به حذف بلوکها در خروجی کد کننده منبع، تابع احتمال گپها به صورت زیر بدست می‌آید :

$$P(G_i = j) = \frac{1}{2^j} \quad i = 1, 2, \dots, k \quad j = 1, 2, 3, \dots \quad (24)$$

به طور مشابه برای بلوکها می‌توان رابطه‌ای مشابه با رابطه (۲۴) نوشت :

$$B^{k'} = B_1, B_2, \dots, B_{k'}$$

$$P(B_i = j) = \frac{1}{2^j} \quad i = 1, 2, \dots, k' \quad j = 1, 2, \dots \quad (25)$$

همانطور که ملاحظه می‌گردد تابع احتمال بیان شده در (۲۴) و (۲۵) کاملاً مشابه رابطه (۹) می‌باشد. بنابراین با

روندی مشابه بخش ۴، می‌توان پارامترهای آزمون گپها و بلوکها را به صورتهای زیر بیان نمود:

$$\chi_g^2 = \sum_{i=1}^m \frac{(g_i - k \cdot 2^{-i})^2}{k \cdot 2^{-i}} + \frac{(g^* - k \cdot 2^{-m})^2}{k \cdot 2^{-m}} \quad (26)$$

$$\chi_b^2 = \sum_{i=1}^m \frac{(b_i - k'.2^{-i})^2}{k'.2^{-i}} + \frac{(b^* - k'.2^{-m})^2}{k'.2^{-m}} \quad (27)$$

در روابط فوق g_i و b_i به ترتیب تعداد گپها و بلوکهای به طول i و g^* و b^* مشخص کننده به ترتیب تعداد گپها و بلوکهایی که طول آنها بیش از m در دنباله های G^k و $B^{k'}$ می باشد. لازم به ذکر است که شرط انجام آزمونهای ارایه شده توسط روابط (26) و (27) عبارت است از :

$$k \times \frac{1}{2^m} > 5 \Rightarrow \frac{n}{4} \times \frac{1}{2^m} > 5 \Rightarrow n > 5 \times 2^{m+1} \quad (28)$$

یا مقدار m انتخابی باید در شرط زیر صدق کند

$$m < \log_2\left(\frac{n}{20}\right) \quad (29)$$

شرط بیان شده در (28) و (29) برای آزمونهای گپ و بلوک به طور یکسان برقرار می باشد. در صورتی که بخواهیم آزمون گپها و بلوکها بر روی گپها و بلوکهای تا طول حداکثر m ، انجام شود و بخواهیم سمبل فرضی g^* و b^* را در نظر بگیریم، پارامترهای آزمون به صورت زیر تبدیل می گردند :

$$\begin{cases} G^L = G_1, G_2, \dots, G_L \\ B^{L'} = B_1, B_2, \dots, B_{L'} \end{cases} \quad (30)$$

$$\begin{cases} P(G_i = j) = \frac{2^{-j}}{1 - 2^{-m}} \\ i = 1, 2, \dots, L \\ P(B_i = j) = \frac{2^{-j}}{1 - 2^{-m}} \\ i = 1, 2, \dots, L' \end{cases} \quad j = 1, 2, \dots, m \quad (31)$$

$$\chi_g^2 = \sum_{i=1}^m \frac{(g_i - L \frac{2^{-i}}{1 - 2^{-m}})^2}{L \times \frac{2^{-i}}{1 - 2^{-m}}} \quad \text{درجه آزادی } m-1 \quad (32)$$

$$\chi_b^2 = \sum_{i=1}^m \frac{(b_i - L' \frac{2^{-i}}{1 - 2^{-m}})^2}{L' \times \frac{2^{-i}}{1 - 2^{-m}}} \quad \text{درجه آزادی } m-1 \quad (33)$$

شرط انجام آزمون و تقریب مناسب توزیع مربع کای برای χ_g^2 و χ_b^2 دقیقاً همان روابط (28) و (29) می باشد.

۶- نتیجه گیری

در یک آزمون آماری استفاده از تابع احتمال و ترتیب دادن آزمون براساس آن بسیار مناسب تر از آزمونی است که براساس متوسط و واریانس تابع احتمال طرح شده است. در این تحقیق به بررسی آمارگان و بیان توزیع رنها و به

تبع آن گپها و بلوکهای یک دنباله کاملاً تصادفی پرداخته شد و سپس براساس آن برای رنھا، گپھا و بلوکھا دو روش آزمون ارایه گردید که یکی دارای درجه آزادی m و دیگری $m-1$ می باشد و بدین ترتیب می توان گپھا و بلوکھای یک دنباله را بطور مجزا مورد آزمون قرارداد.

مراجع:

- [1] S.W.Golomb,*Shift Register Sequences*,Holden-Day, San Fransisco,1982.
- [2] Beker H. and Piper F.,*Cipher Systems: The Protection of Communication*, Northwood Book, London,1982.
- [3] Rukhin,Soto and etc.,"*A Statistical Test suite for Random and Pseudorandon Generator for Cryptographic Application*",NIST Special Publication,2001.
- [4] Kimberely M.,"*Comparison of Two Statistical Tests of Keystream Sequences*",Electronic letter, Vol.23,No.8, PP.365-366, April 1987.
- [5] Larson G.H., *Introduction to Probability and Statistical Inference*, John-Wiley, 1974.