

آمار و احتمالات مهندسی

تمرین هفتم - آمار

حسام و علیرضا

تاریخ تحویل ۱۴۰۱/۱۰/۲۲

سؤال ۱.

فرض کنید برنامه‌ای با $n = ۱۰۰$ فایل کد داریم. بر اثر حمله‌ای که به سرور حاوی فایل‌های برنامه انجام شده، تعدادی باگ به هر فایل اضافه شده است. فرض کنید متغیر X_i تعداد باگ در فایل i ام است که از توزیع پواسون با میانگین ۱ پیروی می‌کند. همچنین تعداد باگ در هر فایل مستقل از فایل‌های دیگر است. احتمال اینکه تعداد کل باگ‌های اضافه شده به برنامه کمتر از ۹۰ باشد را تخمین بزنید.

پاسخ.

با استفاده از قضیه حد مرکزی داریم:

$$\bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

همچنین می‌دانیم در توزیع پواسون داریم:

$$\mu = E[X_i] = \lambda = ۱$$

$$\sigma^2 = Var[X_i] = \lambda$$

بنابراین می‌توان گفت: $\bar{X}_n \sim \mathcal{N}\left(\lambda, \frac{\lambda}{n}\right)$

تعداد کل باگ‌های اضافه شده به برنامه به صورت $Y = \sum_{i=1}^n X_i$ محاسبه می‌شود.

$$\begin{aligned} Y = \sum_{i=1}^n X_i = n\bar{X}_n &\rightarrow P(Y < ۹۰) = P(n\bar{X}_n < ۹۰) = P(\bar{X}_n < ۰٫۹) \\ &= P\left(\frac{\bar{X}_n - \lambda}{\frac{\sqrt{\lambda}}{\sqrt{n}}} < \frac{۰٫۹ - \lambda}{\frac{\sqrt{\lambda}}{\sqrt{n}}}\right) = \Phi\left(\frac{۰٫۹ - \lambda}{\frac{\sqrt{\lambda}}{\sqrt{n}}}\right) = \Phi(-۱) = ۱ - \Phi(۱) = ۰٫۱۵۹ \end{aligned}$$

سؤال ۲.

فرض کنید انتخابات شورای شهر پیش رو است که در آن ۲ کاندیدا با هم به رقابت خواهند پرداخت. یک شرکت نظرسنجی تصمیم گرفته با کمک شما نتیجه انتخابات را پیش‌بینی کند. این شرکت قصد دارد طوری نظرسنجی خود را انجام دهد که با ۹۶ درصد اطمینان، خطای پیش‌بینی کمتر از ۱ درصد باشد.

الف) حداقل از چند نفر باید نظرسنجی کنیم تا دقت خواسته شده توسط شرکت به دست آید؟
 ب) اگر تعداد افراد واجد شرایط رای دادن در شهر کمتر از تعداد افرادی باشد که در قسمت الف به دست آمد، چه اتفاقی خواهد افتاد؟
 این موضوع را تفسیر کنید.

پاسخ.

الف) برای $\hat{\mu}$ آمین نغری که به طور اتفاقی انتخاب شده تعریف می کنیم:

$$X_i = \begin{cases} 0 & \text{برای کاندید A} \\ 1 & \text{برای کاندید B} \end{cases}$$

میانگین جامعه را به دست می آوریم:

$$\bar{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n}$$

خواسته ی شرکت این است که با اطمینان ۹۶ درصد خطای پیش بینی کمتر از ۱ درصد باشد:

$$P(|\bar{X}_n - \mu| \geq 0.01) \leq 0.04$$

با استفاده از نامساوی چیشف داریم:

$$P(\mu - \epsilon < X < \mu + \epsilon) \geq 1 - \frac{\sigma^2}{\epsilon^2} \rightarrow P(|x - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2}$$

$$P(|\bar{X}_n - \mu| \geq 0.01) \leq \frac{\sigma_{\bar{X}_n}^2}{(0.01)^2} = \frac{\sigma_x^2}{n(0.01)^2}$$

متغیر تصادفی X دارای توزیع برنولی است. طبق اصل ناکافی بودن دلیل می دانیم که اگر یک آزمایش تصادفی دارای N نتیجه ممکن باشد و ما هیچ اطلاعی درباره نحوه وقوع آن نداشته باشیم، باید احتمال وقوع هر کدام از آن ها را مساوی فرض کنیم. پس با توجه به اصل ناکافی بودن دلیل می توانیم فرض کنیم که احتمال رای دادن شخص به هر دو کاندید مساوی است (چون صورت سوال هیچ اطلاعاتی در مورد درصد محبوبیت کاندیدها یا احتمال انتخاب هر کاندید توسط اشخاص به ما نداده) پس داریم:

$$X \sim \text{Bernoulli}(0.5) \rightarrow \sigma_X^2 = p(1-p) = \frac{1}{4}$$

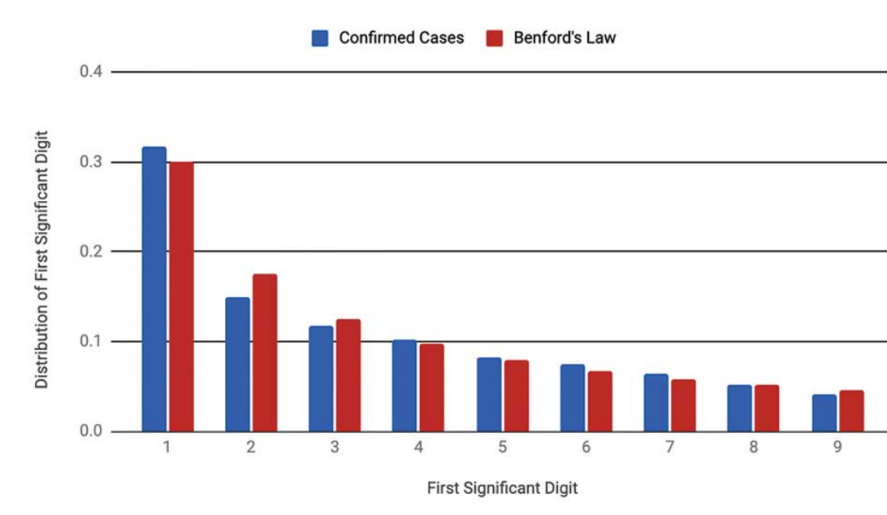
$$\frac{1}{4n(0.01)^2} = 0.04 \rightarrow n = \frac{1}{4 \times 0.04(0.01)^2} = 62500$$

در نتیجه اگر ۶۲۵۰۰ نفر (که به صورت تصادفی و با روش sampling درست انتخاب شده اند) در نظرسنجی شرکت کنند، دقت خواسته شده توسط شرکت به دست خواهد آمد.

ب) اگر تعداد کل افراد واجد شرایط رای دادن، کمتر از ۶۲۵۰۰ نفر باشد، کران به دست آمده توسط نامساوی چیشف بی معنی خواهد بود، چرا که این بدان معنی است که باید از کل جمعیت واجد شرایط رای دادن نظرسنجی کنیم تا با دقت ۹۶ درصد اطمینان داشته باشیم که خطای پیش بینی ما کمتر از ۱ درصد خواهد بود. در حالی که می دانیم نظرسنجی از کل افراد به معنی اطمینان ۱۰۰ درصدی از پیش بینی است. این یکی از ضعف های کران به دست آمده توسط نامساوی چیشف است. دقت کنید که اهمیت نامساوی چیشف در عمومیت آن است. یعنی هیچ فرض و محدودیتی درباره متغیر تصادفی برای برقراری این نامساوی نیاز نیست (البته به جز متناهی بودن واریانس). هرگاه آگاهی بیشتری (افزون بر میانگین و واریانس) درباره ی متغیر تصادفی داشته باشیم (مثلا بیشترین مقدار n در مثال بالا) می توانیم کران های دقیق تری نسبت به کران چیشف به دست آوریم.

سؤال ۳.

نمودار زیر، نشان دهنده توزیع پرارزش ترین رقم میزان مبتلایان روزانه به کووید ۱۹ در آمریکا از ژانویه تا اکتبر سال ۲۰۲۰ (به مدت ۲۹۰ روز) می باشد. مطابق نمودار، پرارزش ترین رقم تعداد افراد مبتلا شده در ماه های مختلف، از قانون بنفورد پیروی می کند! شما به عنوان متخصص آمار،



با فرض صحیح بودن آمار ابتلای کووید در کشور آمریکا قصد دارید تا با استفاده از قانون بنفورد، صحت آمار دیتاست دیگری که به شما داده شده است را بدست آورید.

مطابق با قانون بنفورد، به احتمال 0.47 می‌بایست پرارزش‌ترین رقم تعداد مبتلایان در هر روز برابر با عدد 1 یا 2 باشد. یک دیتاست از میزان مبتلایان در جهان در یک بازه 43 روزه به شما داده شده است. در دیتاست داده شده، در 12 روز از 43 روز، پرارزش‌ترین رقم تعداد مبتلایان برابر با 1 یا 2 می‌باشد. احتمال رخ دادن حالت ذکر شده برای پرارزش‌ترین رقم داده‌ها را در دیتاست داده شده بدست آورید، سپس با توجه به قانون بنفورد، صحت دیتاست داده شده را بررسی کنید. (راهنمایی: از توزیع دوجمله‌ای استفاده کنید. برای بررسی صحت دیتاست داده شده، کافی است تا درصد اختلاف حالت ذکر شده را با پیش‌بینی قانون بنفورد مقایسه کنید).

پاسخ.

متغیر تصادفی X را برابر با تعداد روزهایی در نظر می‌گیریم که پرارزش‌ترین رقم آنها برابر با 1 یا 2 می‌باشد. احتمال آنکه پرارزش‌ترین رقم داده‌ها برابر با 1 یا 2 باشد را میتوان با توزیع دوجمله‌ای با پارامترهای n و p بیان کرد. پس داریم:

$$\binom{43}{12} \times (0.47)^{12} \times (0.53)^{31} = 0.05$$

حال به بررسی صحت دیتای داده شده میپردازیم. آزمون فرض را به صورت مقابل مینویسیم:

H_0 : واقعی بودن دیتاست داده شده

H_A : غیر واقعی بودن دیتاست داده شده

مطابق با قانون بنفورد، پرارزش‌ترین رقم اعداد داده شده با احتمال 0.47 برابر با 1 یا 2 می‌باشد. پس به طور متوسط در یک دیتاست صحیح، تعداد $20 \approx 43 \times 0.47$ عدد از داده‌ها دارای پرارزش‌ترین رقمی برابر با 1 یا 2 می‌باشند. بنابراین برای صحت‌سنجی داده‌های داده شده، اختلاف تعداد روزهایی که دارای رقم پرارزش 1 یا 2 هستند را با 20 مقایسه کرده و صحت آن را می‌سنجیم:

$$\text{Difference Rate} = \frac{|12 - 20|}{20} = 0.4$$

دیتای فوق 40 درصد با میزان بدست آمده توسط قانون بنفورد تفاوت دارد، پس می‌توانیم نتیجه بگیریم که دیتاست داده شده دارای خطا است و واقعی نیست و فرض صفر را رد میکنیم.

سؤال ۴.

فرض کنید X_1, X_2, \dots, X_n یک نمونه تصادفی n تایی از توزیعی با چگالی احتمال زیر باشد.
 الف) برآورد گشتاوری پارامتر θ را به دست آورید. سپس واریانس این توزیع را برحسب پارامتر تخمین زده شده بنویسید.
 ب) همچنین برای نمونه‌ی ۱۰ تایی به دست آمده از توزیع زیر، پارامتر به دست آمده توسط تخمین گر Maximum-Likelihood را محاسبه کنید.

$$f_X(x) = \begin{cases} \theta x^{\theta-1} & 0 < x < 1 \\ 0 & o.w. \end{cases}$$

$$Sample = \{0.0710, 0.2110, 0.4698, 0.3675, 0.5991, 0.9513, 0.6049, 0.9917, 0.1551, 0.2154\}$$

پاسخ.

الف) به دلیل اینکه دو پارامتر مجهول داریم، کافی است $E[X]$ و $E[X^2]$ را محاسبه کنیم:

$$E[X] = \int_0^1 x f_X(x) dx = \int_0^1 \theta x^\theta dx = \left[\frac{\theta x^{\theta+1}}{\theta+1} \right]_0^1 = \frac{\theta}{\theta+1}$$

$$E[X^2] = \int_0^1 x^2 f_X(x) dx = \int_0^1 \theta x^{\theta+1} dx = \left[\frac{\theta x^{\theta+2}}{\theta+2} \right]_0^1 = \frac{\theta}{\theta+2}$$

همچنین داریم:

$$Var[X] = E[X^2] - E[X]^2 = \frac{\theta}{\theta+2} - \left(\frac{\theta}{\theta+1} \right)^2 = \frac{\theta}{(\theta+2)(\theta+1)^2}$$

با استفاده از روش گشتاورها داریم:

$$\hat{M}_1 = E[X] \implies \frac{\sum_{i=1}^n x_i}{n} = \frac{\theta}{\theta+1} \implies \theta = \frac{\bar{X}}{1-\bar{X}}$$

بنابراین برآوردکننده گشتاوری میانگین و واریانس به صورت زیر خواهد بود:

$$\mu = \frac{\theta}{\theta+1} \rightarrow \mu = \bar{X}$$

$$Var[X] = \frac{\theta}{(\theta+2)(\theta+1)^2} = \frac{\frac{\bar{X}}{1-\bar{X}}}{\left(\frac{2-\bar{X}}{1-\bar{X}}\right)\left(\frac{1}{1-\bar{X}}\right)^2} = \frac{\bar{X}(1-\bar{X})^2}{2-\bar{X}}$$

ب) تابع likelihood به صورت زیر تعریف می‌شود:

$$f(x_1, x_2, \dots, x_n; \theta) = \prod_{i=1}^n \theta x_i^{\theta-1} = \theta^n \left(\prod_{i=1}^n x_i \right)^{\theta-1}$$

حال برای تابع log-likelihood داریم:

$$LL(\theta) = n \ln \theta + (\theta - 1) \sum_{i=1}^n \ln x_i$$

برای یافتن تخمین Maximum-Likelihood کافی ست مشتق تابع log-likelihood را گرفته و برابر صفر قرار دهیم:

$$\frac{dLL(\theta)}{d\theta} = \frac{n}{\theta} + \sum_{i=1}^n \ln x_i = 0$$

$$\hat{\theta}_{MLE} = -\frac{n}{\sum_{i=1}^n \ln x_i}$$

برای نمونه داده شده خواهیم داشت:

$$\begin{cases} n = 10 \\ \sum_{i=1}^n \ln x_i = -10.43 \end{cases} \implies \hat{\theta}_{MLE} = -\frac{10}{-10.43} \approx 0.96$$

سؤال ۵.

توزیع مقابل را در نظر بگیرید.

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{\sigma}} e^{-\frac{|x-\mu|}{\sigma}}$$

با استفاده از روش گشتاور، امید ریاضی و واریانس توزیع فوق را بدست آورید.

پاسخ.

توزیع فوق با نام توزیع لاپلاس معروف است و از آن در تحلیل ریسک سرمایه گذاری و تحلیل آمارهای اقتصادی استفاده می شود. از تغییر متغیر $u = x - \mu$ استفاده می کنیم:

$$\begin{aligned} E[X] = m_1 &= \frac{1}{\sqrt{\sigma}} \int_{-\infty}^{\infty} u e^{-\frac{|u|}{\sigma}} du + \mu = \frac{1}{\sqrt{\sigma}} \underbrace{\int_{-\infty}^0 u e^{\frac{u}{\sigma}} du}_{\text{set } u=-u} + \frac{1}{\sqrt{\sigma}} \int_0^{\infty} u e^{-\frac{u}{\sigma}} du + \mu \\ &= -\frac{1}{\sqrt{\sigma}} \int_0^{\infty} u e^{-\frac{u}{\sigma}} du + \frac{1}{\sqrt{\sigma}} \int_0^{\infty} u e^{-\frac{u}{\sigma}} du + \mu = \mu \implies \text{ممان مرتبه اول- امید ریاضی} \end{aligned}$$

$$\begin{aligned} E[X^2] = m_2 &= E[(U + \mu)^2] = E[U^2] + 2\mu E[U] + \mu^2 \\ &= E[U^2] + 2\mu E[X - \mu] + \mu^2 = \frac{1}{\sqrt{\sigma}} \int_{-\infty}^{\infty} u^2 e^{-\frac{|u|}{\sigma}} du + \mu^2 \\ &= \frac{1}{\sqrt{\sigma}} \underbrace{\int_{-\infty}^0 u^2 e^{\frac{u}{\sigma}} du}_{\text{set } u=-u} + \frac{1}{\sqrt{\sigma}} \int_0^{\infty} u^2 e^{-\frac{u}{\sigma}} du + \mu^2 \\ &= \frac{1}{\sqrt{\sigma}} \int_0^{\infty} u^2 e^{-\frac{u}{\sigma}} du + \frac{1}{\sqrt{\sigma}} \int_0^{\infty} u^2 e^{-\frac{u}{\sigma}} du + \mu^2 \\ &= \frac{1}{\sigma} \underbrace{\int_0^{\infty} u^2 e^{-\frac{u}{\sigma}} du}_{\text{set } v=\frac{u}{\sigma}} + \mu^2 = \sigma^2 \int_0^{\infty} v^2 e^{-v} dv + \mu^2 = 2\sigma^2 + \mu^2 \implies \text{ممان مرتبه دوم} \end{aligned}$$

$$\text{Var}[X] = E[X^2] - (E[X])^2 = 2\sigma^2$$

سؤال ۶.

یک تحلیل گر مسابقات بوکس، قصد دارد تا احتمال برخورد موفق ضربات یک بوکسور را بسنجد. او ابتدا با بررسی مسابقات آن بوکسور، نمونه‌ای ۲۰۰ تایی از ضربات او در مسابقات مختلف را بدست می‌آورد. در این نمونه، ۱۳۰ ضربه از ۲۰۰ ضربه به صورت کامل با حریف برخورد کرده است و انحراف معیار نمونه برابر با $\sqrt{2}$ می‌باشد. (ضربات بوکسور نسبت به هم مستقل هستند) مری او ادعا دارد که بوکسور در هر ۲۰۰ ضربه، حداقل ۱۵۰ ضربه را با موفقیت به حریف وارد می‌کند، درستی ادعای وی را بررسی کنید.

پاسخ.

ابتدا فرضیات آزمون فرض را مشخص می‌کنیم:

H_0 : ادعای مری صحیح است و بوکسور بیش از ۱۵۰ ضربه را با موفقیت وارد می‌کند.

H_A : ادعای مری صحیح نیست و نمونه مشاهده شده نمی‌تواند تصادفی باشد.

حال تخمین نقطه‌ای \hat{p} را بدست می‌آوریم:

$$H_0: p = \frac{150}{200} = 0.75$$

$$H_A: \hat{p} = \frac{130}{200} = 0.65$$

حال شرایط لازم برای استفاده از تقریب نرمال را بررسی می‌کنیم. ضربات بوکسور طبق فرض به نسبت به یکدیگر مستقل هستند.

$$np = 200 \times 0.65 = 130 > 10 \quad \checkmark$$

$$n(1-p) = 200 \times 0.35 = 70 > 10 \quad \checkmark$$

حال به محاسبه P-value می‌پردازیم: متغیر تصادفی X را برابر تعداد ضربات موفق قرار می‌دهیم.

$$X \sim N(\mu = \hat{p} = 0.65, \sigma_{\hat{p}} = \frac{\sigma}{\sqrt{n}} = \frac{\sqrt{2}}{\sqrt{200 \times 10}} = 0.1)$$

$$\text{P-value} = P(X > 0.75) = P(Z > \frac{0.75 - 0.65}{0.1}) = P(Z > 1)$$

$$= 1 - P(Z < 1) = 1 - 0.84 = 0.16 > 0.05$$

از آنجایی که P-value بزرگ‌تر از ۰/۰۵ است، با توجه به دیتای داده شده نمی‌توان صحت ادعای مری را رد کرد.

سؤال ۷.

علی در تابستان گذشته در زمینه بلاکچین مطالعاتی را انجام داده است و هم‌اکنون قصد دارد تا با استفاده از شبکه بیت‌کوین، مقداری پول را به حساب یکی از دوستانش بریزد. او هنگام ثبت تراکنش با هزینه‌ای به نام کارمزد تراکنش روبرو شده است. کارمزد تراکنش یک مقدار نسبی است که هر کاربر بابت ثبت شدن تراکنش خود می‌بایست پرداخت کند. اگر مقدار کارمزد تراکنش بسیار کم باشد، تراکنش ثبت نشده و اگر بیشتر از حد معمول باشد، با وجود ثبت شدن تراکنش کاربر پول اضافه‌ای را از دست می‌دهد و ضرر می‌کند. از آنجایی که کارمزد تراکنش یک مبلغ نسبی است، علی با جمع‌آوری آخرین کارمزدهای پرداخت شده در شبکه‌ی بیت‌کوین، سعی دارد تا هزینه تقریبی کارمزد مورد نیاز را تخمین بزند.

علی ۱۰۰ کارمزد را بررسی کرده و میانگین آن‌ها برابر با ۱/۰۸ بدست آمد. انحراف معیار تمامی کارمزدهای شبکه طبق آمارها برابر با ۲/۳ می‌باشد.

الف) بازه اطمینان ۹۵ و ۹۹ درصد را برای میزان کارمزد بهینه‌ای که علی باید پردازد را بدست آورید.

ب) عرض بازه‌های بدست آمده در قسمت الف را با یکدیگر مقایسه کنید و ارتباط میان میزان دقت و عرض بازه را مشخص کنید. شرح دهید که آیا همواره دقت بیشتر در بازه اطمینان برای ما مناسب است یا خیر.

پاسخ.

الف)

$$\bar{X} - \frac{Z \times \sigma}{\sqrt{100}} < \mu < \bar{X} + \frac{Z \times \sigma}{\sqrt{100}}$$

بازه اطمینان ۹۵% :

$$\begin{aligned} \bar{X} - \frac{1.96 \times 2/3}{\sqrt{100}} < \mu < \bar{X} + \frac{1.96 \times 2/3}{\sqrt{100}} \\ \Rightarrow 1/08 - 0/4058 < \mu < 1/08 + 0/4058 \end{aligned}$$

بازه اطمینان ۹۹% :

$$\begin{aligned} \bar{X} - \frac{2.57 \times 2/3}{\sqrt{100}} < \mu < \bar{X} + \frac{2.57 \times 2/3}{\sqrt{100}} \\ \Rightarrow 1/08 - 0/5911 < \mu < 1/08 + 0/5911 \end{aligned}$$

ب)

بازه اطمینان ۹۵٪ کوتاه‌تر (به بیان دیگر دقیق‌تر) از بازه اطمینان ۹۹٪ است. در واقع، بازه‌ای که سطح اطمینان بالاتری دارد، خطای برآورد بیشتری دارد. با توجه به این نکته، استفاده از دقت بسیار بالا می‌تواند منجر به خطای برآورد زیادی شود.