



جلسه‌ی ۲۲: قضیه‌ی حد مرکزی و توزیع نرمال دو متغیره

نگارنده: سیده‌بهناز حسینی سنو

مدیرس: دکتر شهرام خزائی

۱ مقدمه

در جلسات قبل، راجع به متغیرهای تصادفی مستقل با توزیع یکسان^۱ صحبت کردیم. گفتیم چنانچه تعداد n متغیر i.i.d را جمع و بر تعدادشان تقسیم کنیم، حد این عبارت به مقدار متوسط متغیرها میل می‌کند. از این اتفاق دو تفسیر ارائه کردیم:

تفسیر اول: همگرایی در احتمال

$$\forall \varepsilon > 0 : \lim_{n \rightarrow +\infty} \Pr\left\{\left|\frac{X_1 + X_2 + \dots + X_n}{n} - E[X]\right| > \varepsilon\right\} = 0$$

تفسیر دوم: همگرایی almost sure

$$\lim_{n \rightarrow +\infty} \frac{X_1 + X_2 + \dots + X_n}{n} = E[X]$$

(توجه شود که این نوع همگرایی، قوی‌تر از همگرایی در احتمال است؛ اما در این درس نیاز به آن نمی‌پردازیم.)

نوع دیگری از همگرایی وجود دارد که همگرایی در توزیع^۲ نامیده می‌شود و مبنای قضیه‌ی بسیار مهم حد مرکزی^۳ است که در این جلسه به آن خواهیم پرداخت.

۲ همگرایی در توزیع و قضیه‌ی حد مرکزی

قضیه ۱ (قضیه حد مرکزی) اگر X_1, X_2, \dots دنباله‌ای از متغیرهای تصادفی مستقل با توزیع یکسان با مقدار متوسط μ و واریانس σ^2 باشند، توزیع

$$Z_n = \frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}}$$

^۱Independent and identically distributed random variables (i.i.d)

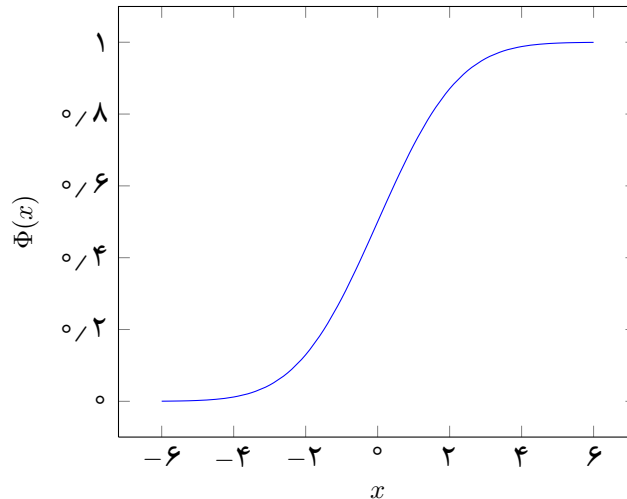
^۲Convergence in distribution

^۳Central limit theorem

به توزیع متغیر تصادفی نرمال استاندارد میل می‌کند. به بیان دیگر:

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pr(Z_n \leq t) &= \lim_{n \rightarrow \infty} \Pr\left(\frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}} \leq t\right) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx = \Phi(t). \end{aligned}$$

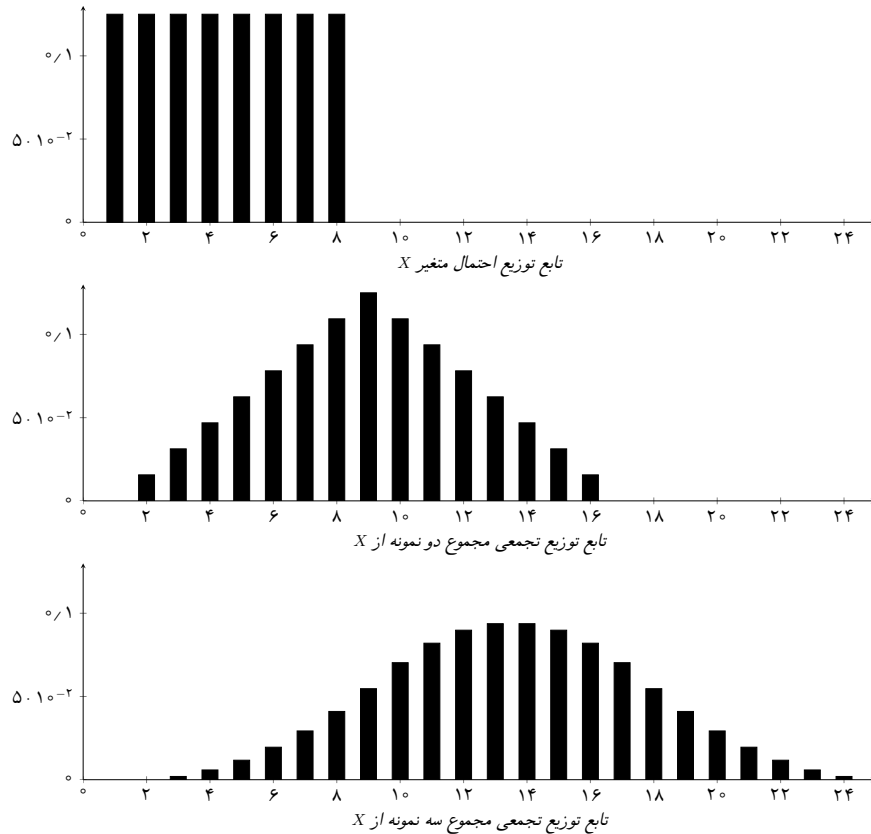
همان‌طور که در جلسات قبل گفته شد، تابع $\Phi(t)$ ، تابع توزیع تجمعی متغیر تصادفی نرمال (گوسی) استاندارد است که به شکل زیر است:



نکته ۱ به طور کلی، توزیع‌های نامتقارن (مانند توزیع پواسون یا توزیع نمایی)، به شرط موجود و محدود بودن مقدار متوسط و واریانس، زمان بیشتری برای همگرایی می‌برند. ♦

قضیه‌ی حد مرکزی نشان می‌دهد که بسیاری از پدیده‌های طبیعی، با تقریب خوبی یک توزیع نرمال استاندارد هستند. تغییرات حاصل از فرایندهای طبیعی، اغلب نتیجه‌ی جمع شدن تعداد زیادی عامل کوچک و ناچیز است. اگرچه تأثیر هر یک از این عوامل به تنهایی ممکن است نادیده گرفته شود، اما تأثیری که مجموع آن‌ها بر جای می‌گذارد قابل چشم‌پوشی نیست. از این رو، مطالعه‌ی توزیع مجموع تعداد زیادی از متغیرهای تصادفی مستقل ضروری است. به طور مثال، نویز موجود در مدارهای مخابراتی حاصل از چندین اثر است و جمع شدن تعداد زیادی نویز مختلف، در نهایت توزیع گاوسی به دست می‌دهد.

مثال ۱ فرض کنید X یک متغیر تصادفی گسسته است که مقادیر ۱ تا ۸ را با احتمال $\frac{1}{8}$ می‌گیرد. تابع جرم احتمال مجموع دو نمونه از متغیر تصادفی X با محاسبه‌ی کانولوشن X با خودش به دست می‌آید. اگر این روند را برای دو، سه، چهار و یا تعداد بیشتری نمونه از X ادامه دهیم، با افزایش تعداد نمونه‌ها مشاهده می‌شود که فرم تابع جرم احتمال مجموع متغیرها به فرم تابع چگالی احتمال متغیر تصادفی نرمال نزدیک می‌شود. اما این به معنای میل کردن تابع جرم احتمال مجموع متغیرها به سمت تابع چگالی احتمال متغیر تصادفی نرمال نیست. در واقع، قضیه حد مرکزی به صورتی که بیان شد، صرفاً همگرایی تابع توزیع تجمعی مجموع متغیرهای تصادفی (که با کم کردن میانگین و تقسیم بر انحراف معیار نرمال نیز شده‌اند) را به سمت تابع توزیع تجمعی احتمال متغیر تصادفی نرمال استاندارد، تضمین می‌کند.



نکته ۲ قضیه حد مرکزی درباره همگرایی تابع توزیع تجمعی دنباله متغیرهای تصادفی $\frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}}$ به سمت تابع توزیع تجمعی متغیر تصادفی نرمال استاندارد، مطرح می‌شود. با این وجود، در عمل برای n های نسبتاً بزرگ (مثلاً $n > 10$) تابع توزیع تجمعی متغیر تصادفی $X_1 + X_2 + \dots + X_n$ با تابع توزیع تجمعی متغیر تصادفی نرمال $N(n\mu, n\sigma^2)$ تقریب زده می‌شود.

مثال زیر کاربردی از این تقریب را نشان می‌دهد.

مثال ۲ کسر f از افراد جامعه نوشابه را به دلستر ترجیح می‌دهند. مقدار f را تخمین بزنید. مطلوب هر مسئله‌ی تخمین، نزدیک بودن جواب به دست آمده به جواب واقعی است. از n نفر به طور تصادفی نمونه‌گیری می‌کنیم. می‌خواهیم با احتمال بیش از ۹۵ درصد، خطای محاسبه‌ی f کمتر از ۰/۰۱ شود. به بیان دیگر می‌خواهیم:

$$\Pr(|M_n - f| > 0.01) \leq 0.05 \quad (1)$$

که

$$M_n = \frac{X_1 + X_2 + \dots + X_n}{n},$$

و متغیر تصادفی X_i را به این صورت تعریف می‌کنیم:

$$X_i = \begin{cases} 1 & \text{اگر شخص } i \text{ ام بگوید بله} \\ 0 & \text{oth.} \end{cases}$$

که X_i دارای توزیع برنولی با پارامتر f بوده و متوسط و انحراف معیار آن به صورت زیر است:

$$E[X_i] = f, \text{Var} = \sqrt{f(1-f)} \leq \frac{1}{4} \quad (2)$$

در جلسه ی قبل، بدون آگاهی از توزیع M_n و تنها با استفاده از قضیه ی چبیشف، تعداد نمونه ها را 50000 تخمین زدیم. با استفاده از قضیه ی حد مرکزی نشان می دهیم تعداد نمونه هایی که در عمل بدان نیاز داریم کمتر از کرانی است که قضیه ی چبیشف به دست می دهد. داریم:

$$\begin{aligned} \Pr(|M_n - f| > 0.01) &= \Pr\left(\frac{X_1 + \dots + X_n - nf}{n} > 0.01\right) \\ &= \Pr\left(\frac{X_1 + \dots + X_n - nf}{\sigma\sqrt{n}} > \frac{0.01\sqrt{n}}{\sigma}\right) \\ &\approx \Pr(|Z| > \frac{0.01\sqrt{n}}{\sigma}) \end{aligned}$$

که Z یک متغیر تصادفی با توزیع نرمال استاندارد است و به عنوان تخمینی از متغیر تصادفی اولیه در نظر گرفته شده است. از (1) و (2) داریم:

$$\Pr(|Z| > \frac{0.01\sqrt{n}}{\sigma}) \leq \Pr(|Z| > \frac{0.01\sqrt{n}}{\frac{1}{2}}) = 0.05 \quad (3)$$

می دانیم:

$$\begin{aligned} \Pr(|Z| \geq a) &= 2 \Pr(Z > a) \\ &= 2(1 - \Pr(Z \leq a)) \\ &= 2(1 - \Phi(a)) \end{aligned} \quad (4)$$

پس از (3) و (4) می توان نوشت:

$$2(1 - \Phi(0.02\sqrt{n})) = 0.05$$

$$\Phi(0.02\sqrt{n}) \cong 0.975 \implies n \cong 9600. \quad \blacklozenge$$

مثال 3 (Random walks) شخص مستی در مبدأ مکان قرار دارد و در هر مرحله به سمت چپ یا راست برمی دارد. مطلوب است متوسط فاصله ی شخص از مبدأ پس از N مرحله. متغیر تصادفی X_i را به صورت زیر تعریف می کنیم:

$$X_i = \begin{cases} 1 & \text{با احتمال } \frac{1}{3} \\ -1 & \text{با احتمال } \frac{2}{3} \end{cases}$$

اگر موقعیت شخص را در لحظه ی N ام با S_N نشان دهیم، داریم:

$$S_N = X_1 + X_2 + \dots + X_N$$

در جلسات قبل، برای سنجش متوسط فاصله‌ی شخص پس از N مرحله نشان دادیم:

$$E[S_N^2] = N, E[|S_N|] \leq \sqrt{N}$$

حال می‌خواهیم ثابت کنیم که جواب مسئله ضریبی از \sqrt{N} است. در اینجا تعدادی متغیر تصادفی *i.i.d* با $E[X_i] = 0$ و $\text{Var}(X_i) = 1$ داریم. باید یک متغیر تصادفی بیابیم که در حد به سمت توزیع نرمال میل کند. پس:

$$\begin{aligned} E[|S_N|] &= E[|X_1 + X_2 + \dots + X_n|] \\ &= \sqrt{N} E\left[\left|\frac{X_1 + \dots + X_N}{\sqrt{N}}\right|\right] \\ &\simeq \sqrt{N} E[|Z|] \\ &= \sqrt{N} \int_{-\infty}^{\infty} |z| \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz \\ &= \sqrt{\frac{2N}{\pi}} \quad \blacklozenge \end{aligned}$$

۳ توزیع نرمال دو متغیره

در جلسات قبل راجع به همبستگی صحبت کردیم. گفتیم متغیرهای تصادفی مستقل، ناهمبسته‌اند. یعنی کوواریانس و در نتیجه ضریب همبستگی آنها برابر صفر است که ضریب همبستگی را چنین تعریف می‌کنیم:

$$\rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_X \sigma_Y} = \frac{E[XY] - \mu_X \mu_Y}{\sigma_X \sigma_Y}$$

توجه شود که اگر دو متغیر تصادفی همبسته نباشند، لزوماً مستقل نیستند؛ یعنی ربط سیستماتیک بین این دو متغیر وجود ندارد، اما ممکن است ربط نهانی داشته باشند. می‌توان متغیرهای نرمالی تعریف کرد که از ناهمبستگی آنها، بتوان مستقل بودنشان را نتیجه گرفت.

متغیرهای تصادفی X و Y را که دارای تابع چگالی مشترک $f_{X,Y}(x,y)$ هستند مشترکاً نرمال گوئیم، اگر دارای ویژگی‌های زیر باشند:

- توزیع کناری X نرمال باشد (پس خود X باید نرمال باشد).

- توزیع Y به شرط $X = x$ به ازای هر x حقیقی نرمال باشد.

- متوسط Y به شرط $X = x$ تابعی خطی از x باشد؛ یعنی $E[Y|X = x] = ax + b$.

- واریانس Y به شرط $X = x$ مقداری ثابت باشد.

می‌توان نشان داد تنها تابع توزیع مشترکی که در شرایط بالا صدق می‌کند بدین شکل است:

$$f_{X,Y}(x,y) = \frac{1}{\sqrt{2\pi\sigma_X\sigma_Y\sqrt{1-\rho}}} \exp\left\{-\frac{1}{2(1-\rho^2)} Q(x,y)\right\},$$

که

$$Q(x,y) = \left(\frac{x-\mu_X}{\sigma_X}\right)^2 - 2\rho\left(\frac{x-\mu_X}{\sigma_X}\right)\left(\frac{y-\mu_Y}{\sigma_Y}\right) + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2$$

و $-1 \leq \rho \leq 1$ - ضریب همبستگی X و Y است. با قرار دادن $\rho = 0$ در این تابع، $f_{X,Y}(x,y)$ حاصل ضرب دو توزیع نرمال X و Y خواهد شد. پس می‌توان نتیجه گرفت که متغیرهای تصادفی مشترکاً نرمال ناهمبسته، مستقل‌اند.

مثال ۴ متغیر تصادفی X را معدل سال اول یک دانشجو و Y را معدل سال آخر وی در نظر بگیرید. معمولاً هر دو متغیر دارای توزیع نرمال هستند و توزیعی مشترکاً نرمال دارند.

مثال ۵ اگر قد پدری را X و قد دخترش را Y بگیریم، این دو متغیر نیز توزیع مشترکاً نرمال خواهند داشت.

توزیع‌های نرمال چند متغیره نیز تعریف می‌شوند که در درس آمار به آن پرداخته خواهد شد.