

Bài 2 - Ôn tập Biến ngẫu nhiên và Giới thiệu Thống kê Tính toán (Review of Random Variables and Introduction to Computational Statistics)

Thống kê máy tính và ứng dụng (CLC)

Vũ Quốc Hoàng (vqhoang@fit.hcmus.edu.vn)

FIT - HCMUS

Ngày 17 tháng 1 năm 2022

Nội dung

1. Bài toán thu thập phiếu thưởng
2. Ôn tập biến ngẫu nhiên rời rạc
3. Luật Zipf và Truyện Kiều - Nguyễn Du
4. Ôn tập biến ngẫu nhiên liên tục
5. Phương pháp xấp xỉ phân phối bằng mô phỏng

Bài toán thu thập phiếu thưởng

Bài toán thu thập phiếu thưởng (coupon collector's problem). Một cửa hàng phát hành n loại phiếu thưởng khác nhau. Mỗi lần khách mua hàng sẽ được cửa hàng phát một tờ phiếu thưởng ngẫu nhiên trong n loại. Khi khách thu thập được đủ n loại phiếu thưởng thì được cửa hàng tặng quà.

Hỏi: khách cần mua hàng bao nhiêu lần từ cửa hàng để được nhận quà?

Ôn tập biến ngẫu nhiên rời rạc

Biến ngẫu nhiên

Nếu giá trị của một đại lượng/tính chất X được xác định hoàn toàn khi biết kết quả ω của một thí nghiệm T thì X được gọi là một đại lượng/biến ngẫu nhiên (**liên quan đến T**)

- Trước khi biết kết quả, ta chỉ biết X có thể nhận một giá trị nào đó trong tập A ,
- Sau khi biết kết quả ω , ta biết X nhận một giá trị cụ thể $x \in A$, ta kí hiệu $X(\omega) = x$.

Biến ngẫu nhiên (random variable) là **hàm** trên không gian mẫu Ω

- $X : \Omega \rightarrow A$, gán mỗi kết quả $\omega \in \Omega$ một giá trị $X(\omega) \in A$,
- A được gọi là **tập/miền giá trị** của X và thường là tập con của \mathbb{R} .

Biến ngẫu nhiên là phương tiện hay được dùng để mô tả các biến cố. Xét biến (số) ngẫu nhiên X liên quan đến thí nghiệm T có không gian mẫu là Ω . Cho $C \subset \mathbb{R}$, ta kí hiệu biến cố “ X nhận giá trị trong C ” là

$$(X \in C) = \{\omega \in \Omega : X(\omega) \in C\}.$$

Ôn tập biến ngẫu nhiên rời rạc

Phân phối của biến ngẫu nhiên

Xét biến ngẫu nhiên X liên quan đến thí nghiệm T có không gian mẫu là Ω . Tập các xác suất $\{P(X \in C) : C \subset \mathbb{R}\}$ xác định một độ đo xác suất trên (không gian mẫu mới) \mathbb{R} và được gọi là **phân phối** (distribution) của X .

- Phân phối của X cho thấy khả năng X nhận các giá trị khác nhau.
- Với phân phối của X , ta khảo sát X mà không cần đề ý đến T hay Ω nữa.
- Nói chung, tập $\{P(X \in C) : C \subset \mathbb{R}\}$ là “rất khó tính toán”. Ta cần cách nào đó giúp **xác định** phân phối của X để “dễ tính toán hơn”.

Ôn tập biến ngẫu nhiên rời rạc

Biến ngẫu nhiên rời rạc và hàm xác suất

- X được gọi là **biến ngẫu nhiên rời rạc** (discrete random variable) nếu tập giá trị của nó là **rời rạc** (**hữu hạn** (finite) hoặc **vô hạn đếm được** (countably infinite)).
- Với X là biến ngẫu nhiên rời rạc, **hàm xác suất** (probability function, probability mass function) của X là hàm $f : \mathbb{R} \rightarrow \mathbb{R}$, được xác định bởi

$$f(x) = f_X(x) = P(X = x), x \in \mathbb{R}.$$

- Hàm xác suất f cho biết khả năng X nhận một giá trị cụ thể.
- Tập số thực $\{x \in \mathbb{R} : f(x) > 0\}$ được gọi là **tập hỗ trợ** (support) của X , kí hiệu $\text{Sup}(X)$.
- Hàm xác suất có tính chất: $f(x) \geq 0, \forall x \in \mathbb{R}$ và $\sum_{x \in \text{Sup}(X)} f(x) = 1$.
- **Hàm xác suất xác định phân phối của biến ngẫu nhiên rời rạc**

$$P(X \in C) = \sum_{x \in C} f(x), C \subset \mathbb{R}.$$

Ôn tập biến ngẫu nhiên rời rạc

Các biến ngẫu nhiên độc lập

Hai biến ngẫu nhiên X, Y được gọi là **độc lập** (independent) nếu với mọi $A, B \subset \mathbb{R}$ ta có

$$P((X \in A) \cap (Y \in B)) = P(X \in A)P(Y \in B).$$

Nghĩa là việc X nhận giá trị nào cũng không ảnh hưởng đến khả năng nhận giá trị nào đó của Y (và ngược lại).

Mệnh đề. Hai biến ngẫu nhiên rời rạc X, Y độc lập khi và chỉ khi

$$P((X = x) \cap (Y = y)) = P(X = x)P(Y = y)$$

với mọi $x, y \in \mathbb{R}$.

Ôn tập biến ngẫu nhiên rời rạc

Kì vọng của biến ngẫu nhiên

Cho biến ngẫu nhiên rời rạc X với hàm xác suất f , **kì vọng** (mean) của X , kí hiệu $E(X)$, là số thực được tính bởi (“**nếu tính được**”)

$$E(X) = \mu_X = \mu = \sum_x xP(X = x) = \sum_x xf(x).$$

Kì vọng của X là giá trị trung bình của các giá trị mà X có thể nhận với **trọng số là xác suất** để X nhận các giá trị tương ứng đó.

Cho biến ngẫu nhiên $X : \Omega \rightarrow \mathbb{R}$ và hàm số $r : \mathbb{R} \rightarrow \mathbb{R}$, ta nói $Y : \Omega \rightarrow \mathbb{R}$ là biến ngẫu nhiên **phái sinh** từ X qua hàm số r , kí hiệu $Y = r(X)$, nếu Y được xác định bởi

$$Y(\omega) = r(X(\omega)), \omega \in \Omega.$$

Khi đó ta có

$$E(Y) = E(r(X)) = \sum_x r(x)f(x).$$

Ôn tập biến ngẫu nhiên rời rạc

Phương sai và độ lệch chuẩn của biến ngẫu nhiên

Cho biến ngẫu nhiên rời rạc X với hàm xác suất f và kì vọng $\mu = E(X)$, **phương sai** (variance) của X , kí hiệu $Var(X)$, là số thực được tính bởi (“**nếu tính được**”)

$$Var(X) = \sigma_X^2 = \sigma^2 = E((X - \mu)^2) = \sum_x (x - \mu)^2 P(X = x) = \sum_x (x - \mu)^2 f(x).$$

Khi đó ta cũng nói $\sigma = \sigma_X = \sqrt{\sigma_X^2} = \sqrt{Var(X)}$ là **độ lệch chuẩn** (standard deviation) của X . Lưu ý: độ lệch chuẩn có **cùng đơn vị** với X nhưng phương sai thì không.

Phương sai (và độ lệch chuẩn) phản ánh sự **phân tán** của phân phối của biến ngẫu nhiên.

Mệnh đề. Cho X là biến ngẫu nhiên (có phương sai), ta có

$$Var(X) = E(X^2) - (E(X))^2.$$

Ôn tập biến ngẫu nhiên rời rạc

Các tính chất quan trọng của kì vọng và phương sai

Cho X_1, X_2, \dots, X_n là các biến ngẫu nhiên (có kỳ vọng), ta có

$$E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i) \quad (\text{linearity of expectation})$$

Cho X là biến ngẫu nhiên và a, b là các hằng số thực, ta có

1. $E(aX + b) = aE(X) + b$,
2. $\text{Var}(aX + b) = a^2 \text{Var}(X)$.

Cho X, Y là hai biến ngẫu nhiên **độc lập**, ta có

1. $E(XY) = E(X)E(Y)$,
2. $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.

Ôn tập biến ngẫu nhiên rời rạc

Hàm đặc trưng của biến cố

Cho biến cố A liên quan đến thí nghiệm T với không gian mẫu Ω , ta gọi **hàm đặc trưng** (characteristic function, indicator function) của A là hàm $\mathbb{I}_A : \Omega \rightarrow \mathbb{R}$ được xác định bởi

$$\mathbb{I}_A(\omega) = \begin{cases} 1 & \text{nếu } \omega \in A, \\ 0 & \text{nếu } \omega \notin A. \end{cases}$$

Hàm đặc trưng giúp khảo sát biến cố như là một biến ngẫu nhiên.

Mệnh đề. Với mọi biến cố A ta có

$$E(\mathbb{I}_A) = P(A).$$

Ôn tập biến ngẫu nhiên rời rạc

Phân phối Bernoulli

Biến ngẫu nhiên rời rạc X được gọi là có **phân phối Bernoulli** (Bernoulli distribution) với tham số p ($0 \leq p \leq 1$), kí hiệu $X \sim \text{Bernoulli}(p)$, nếu X có tập giá trị là $\{0, 1\}$ và

$$f(x) = P(X = x) = \begin{cases} p & \text{nếu } x = 1, \\ 1 - p & \text{nếu } x = 0. \end{cases}$$

Khi đó, X có kì vọng $E(X) = p$ và phương sai $\text{Var}(X) = p(1 - p)$.

Xét thí nghiệm tung một đồng xu với xác suất ra ngửa p , gọi X là “số lần được ngửa” thì $X \sim \text{Bernoulli}(p)$. Trường hợp đồng xu đồng chất thì $X \sim \text{Bernoulli}(0.5)$.

Xét thí nghiệm T với biến cố A có $P(A) = p$, khi đó $\mathbb{I}_A \sim \text{Bernoulli}(p)$.

Ôn tập biến ngẫu nhiên rời rạc

Phân phối nhị thức

Biến ngẫu nhiên rời rạc X được gọi là có **phân phối nhị thức** (binomial distribution) với tham số n ($n \in \mathbb{N}$), p ($0 \leq p \leq 1$), kí hiệu $X \sim \mathcal{B}(n, p)$, nếu X có tập giá trị là $\{0, 1, \dots, n\}$ và

$$f(x) = P(X = x) = C_n^x p^x (1 - p)^{n-x}, x \in \{0, 1, \dots, n\}.$$

Khi đó, X có kì vọng $E(X) = np$ và phương sai $\text{Var}(X) = np(1 - p)$.

Cho thí nghiệm T với biến cố A có $P(A) = p$. Xét thí nghiệm R “thực hiện T lặp lại n lần **độc lập**”, gọi X là “số lần A xảy ra” thì $X \sim \mathcal{B}(n, p)$.

Mệnh đề. Nếu X_1, X_2, \dots, X_n là các biến ngẫu nhiên **độc lập và cùng phân phối** (independent and identically distributed - iid) Bernoulli với tham số p , thường kí hiệu $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$, và $X = \sum_{i=1}^n X_i$ thì $X \sim \mathcal{B}(n, p)$.

Ôn tập biến ngẫu nhiên rời rạc

Phân phối hình học

Biến ngẫu nhiên rời rạc X được gọi là có **phân phối hình học** (geometric distribution) với tham số p ($0 < p \leq 1$), kí hiệu $X \sim \text{Geometric}(p)$, nếu X có tập giá trị là $\{1, 2, \dots\}$ và

$$f(x) = P(X = x) = (1 - p)^{x-1}p, x \in \{1, 2, \dots\}.$$

Khi đó, X có kì vọng $E(X) = \frac{1}{p}$ và phương sai $\text{Var}(X) = \frac{1-p}{p^2}$.

Cho thí nghiệm T với biến cố A có $P(A) = p$. Xét thí nghiệm R “thực hiện T lặp lại nhiều lần **độc lập** cho đến khi A xảy ra thì dừng”, gọi X là “số lần thực hiện” thì $X \sim \text{Geometric}(p)$.

Mệnh đề (tính không nhớ - memoryless). Cho $X \sim \text{Geometric}(p)$, với mọi $n = 1, 2, \dots$ và mọi $k = 0, 1, \dots$ ta có

$$P(X = k + n | X > k) = P(X = n).$$

Ôn tập biến ngẫu nhiên rời rạc

Phương pháp xấp xỉ kỳ vọng bằng mô phỏng

Để xấp xỉ kỳ vọng $E(X)$ của một biến ngẫu nhiên X liên quan đến thí nghiệm T , ta có thể dùng phương pháp thống kê như sau

- Thực hiện lặp lại N lần (độc lập) thí nghiệm T , ghi nhận các giá trị mà X nhận x_1, x_2, \dots, x_N (còn gọi là **mẫu dữ liệu** - sample), và tính **trung bình mẫu**

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_N}{N}.$$

- Khi N đủ lớn, ta có $\bar{x} \approx E(X)$.
- Việc thực hiện lặp lại nhiều lần thí nghiệm T có thể được **mô phỏng** (simulate) trên máy tính.

Bài toán thu thập phiếu thưởng

Tính toán chính xác

Gọi X là số lần khách cần mua hàng để được nhận quà, tức là số lần cần mua hàng từ đầu cho đến khi thu thập **vừa đủ** n loại phiếu thưởng.

Gọi X_i là số lần cần mua hàng từ lúc **vừa đã có** $i - 1$ loại phiếu thưởng cho đến khi thu thập thêm được một loại mới để có **vừa đúng** i loại phiếu thưởng ($i = 1, 2, \dots, n$).

Ta có $X = \sum_{i=1}^n X_i$ và X_i có phân phối hình học với tham số

$$p_i = \frac{n - (i - 1)}{n} = \frac{n - i + 1}{n} \quad (i = 1, 2, \dots, n).$$

Từ đó ta có

$$E(X) = E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i) = \sum_{i=1}^n \frac{n}{n - i + 1} = n \sum_{i=1}^n \frac{1}{i} = nH_n,$$

với $H_n = \sum_{i=1}^n \frac{1}{i}$ được gọi là **số điều hòa** (harmonic number) thứ n .

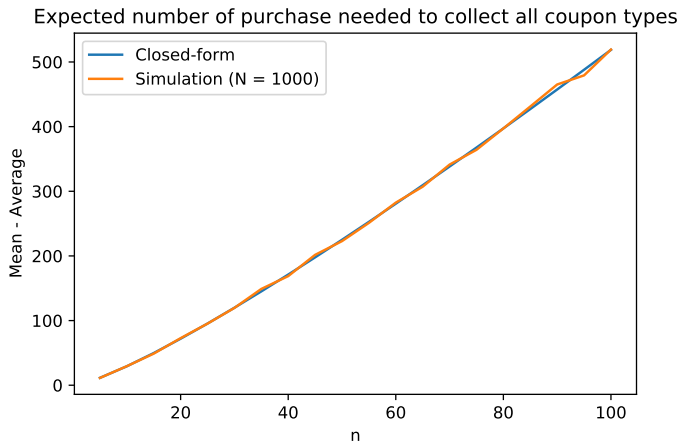
Bài toán thu thập phiếu thưởng

Mô phỏng

```
def num_buy_to_win(n):  
    coupons = []  
    while len(set(coupons)) < n:  
        coupons.append(random.randint(1, n))  
    return len(coupons)  
  
def average(n, N, X):  
    m = sum(X(n) for _ in range(N))  
    return m/N  
  
average(10, 1000, num_buy_to_win)  
#29.175
```

Bài toán thu thập phiếu thưởng

Kết quả



Luật Zipf và Truyện Kiều - Nguyễn Du

Luật Zipf

Luật Zipf (Zipf's law) trong ngôn ngữ học định lượng: **tần số** (frequency) của từ tỉ lệ nghịch với **hạng** (rank) của nó (trong nhiều kho ngữ liệu)

$$f(r) = c \times \frac{1}{r^s} \text{ hay } \log f(r) = \log c - s \log r$$

trong đó hằng số c là hệ số tỉ lệ, hằng số $s \approx 1$ là số mũ, $f(r)$ là tần số của từ có hạng r ($r = 1, 2, \dots$).

Luật Zipf-Mandelbrot mở rộng luật Zipf

$$f(r) = c \times \frac{1}{(r + q)^s} \text{ hay } \log f(r) = \log c - s \log(r + q).$$

(https://en.wikipedia.org/wiki/Zipf%27s_law.)

Luật Zipf và Truyện Kiều - Nguyễn Du

Truyện Kiều - Nguyễn Du

*“ ... Dưới cầu nước chảy trong veo
Bên cầu tơ liễu bóng chiều thướt tha ...”*

Truyện Kiều của đại thi hào Nguyễn Du là một tuyệt tác trong kho tàng văn học Việt Nam

- Thơ lục bát,
- 3,254 câu,
- 22,778 từ,
- 2,383 từ khác nhau.

(https://vi.wikipedia.org/wiki/Truy%E1%BB%87n_Ki%E1%BB%81u.)

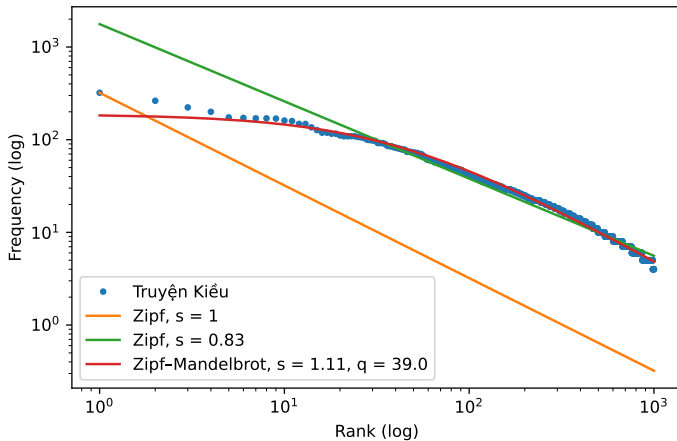
Luật Zipf và Truyện Kiều - Nguyễn Du

Truyện Kiều - Nguyễn Du (tt)

Rank	Word	Frequency	Rank	Word	Frequency
1	một	321	11	rằng	159
2	đã	263	12	lại	148
3	người	223	13	ra	148
4	nàng	200	14	hoa	136
5	lòng	174	15	tình	127
6	lời	172	16	còn	119
7	là	170	17	mới	119
8	cho	170	18	ai	116
9	cũng	169	19	đâu	116
10	có	161	20	chẳng	111

Luật Zipf và Truyện Kiều - Nguyễn Du

Truyện Kiều - Nguyễn Du (tt)



Luật Zipf

Hỏi: tại sao tần số với hạng của từ lại có quan hệ theo **luật lũy thừa** (power law) như vậy?

Ôn tập biến ngẫu nhiên liên tục

Biến ngẫu nhiên liên tục và hàm mật độ xác suất

X được gọi là **biến ngẫu nhiên liên tục** (continuous random variable) nếu có hàm số không âm $f : \mathbb{R} \rightarrow \mathbb{R}$ sao cho với mọi khoảng $[a, b]$ trong \mathbb{R} ta có

$$P(a \leq X \leq b) = \int_a^b f(x)dx.$$

- f được gọi là **hàm mật độ xác suất** (probability density function) của X vì nó cho biết khả năng X nhận giá trị trong các khoảng rất nhỏ của trục số thực \mathbb{R}

$$P(a \leq X \leq a + \epsilon) = \int_a^{a+\epsilon} f(x)dx \approx \epsilon f(a) \text{ khi } \epsilon \text{ rất nhỏ.}$$

- Tập số thực $\{x \in \mathbb{R} : f(x) > 0\}$ được gọi là **tập hỗ trợ** (support) của X , kí hiệu $\text{Sup}(X)$.
- Hàm mật độ xác suất có tính chất: $f(x) \geq 0, \forall x \in \mathbb{R}$ và $\int_{-\infty}^{\infty} f(x)dx = 1$.

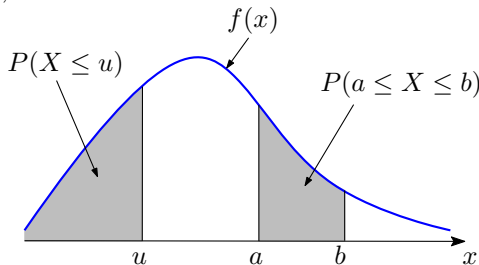
Ôn tập biến ngẫu nhiên liên tục

Hàm mật độ xác suất (tt)

Hàm mật độ xác suất xác định phân phối của biến ngẫu nhiên liên tục

$$P(X \in C) = \int_C f(x)dx, C \subset \mathbb{R}.$$

- $P(X = u) = \int_u^u f(x)dx = 0,$
- $P(X < u) = P(X \leq u) = \int_{-\infty}^u f(x)dx,$
- $P(X > u) = P(X \geq u) = \int_u^{\infty} f(x)dx,$
- $P(a \leq X \leq b) = \int_a^b f(x)dx.$



Từ nhận xét $P(X = u) = 0$, ta thấy tồn tại các biến cố dù có xác suất là 0 nhưng vẫn có thể xảy ra (có E với $P(E) = 0$ nhưng $E \neq \emptyset$).

Ôn tập biến ngẫu nhiên liên tục

Hàm phân phối

Hàm phân phối (tích lũy) (distribution function, cumulative distribution function) của một biến ngẫu nhiên X là hàm số $F : \mathbb{R} \rightarrow \mathbb{R}$ được xác định bởi

$$F(x) = P(X \leq x) = \begin{cases} \sum_{t \leq x} f(t) & \text{nếu } X \text{ rời rạc với hàm xác suất } f, \\ \int_{-\infty}^x f(t) dt & \text{nếu } X \text{ liên tục với hàm mật độ xác suất } f. \end{cases}$$

F xác định phân phối của X .

Hàm phân phối F có các tính chất

1. Tăng: nếu $x_1 \leq x_2$ thì $F(x_1) \leq F(x_2)$,
2. Chuẩn hóa: $\lim_{x \rightarrow -\infty} F(x) = 0$ và $\lim_{x \rightarrow \infty} F(x) = 1$,
3. Liên tục phải: $F(x) = F(x^+) = \lim_{t \rightarrow x, t > x} F(t)$.
4. Nếu X là biến ngẫu nhiên liên tục thì F là hàm liên tục và nếu F có đạo hàm tại x thì $F'(x) = f(x)$.

Ôn tập biến ngẫu nhiên liên tục

Hàm phân phối đồng thời

Hàm phân phối đồng thời (joint distribution function) của hai biến ngẫu nhiên X, Y là hàm số $F_{XY} : \mathbb{R}^2 \rightarrow \mathbb{R}$ được xác định bởi

$$F_{XY}(x, y) = P(X \leq x, Y \leq y) \quad (x, y \in \mathbb{R}).$$

Mệnh đề. Hai biến ngẫu nhiên X, Y độc lập khi và chỉ khi $F_{XY}(x, y) = F_X(x)F_Y(y)$ với mọi $x, y \in \mathbb{R}$.

Hai biến ngẫu nhiên X, Y được gọi là **liên tục đồng thời** (jointly continuous) nếu có hàm số không âm $f_{XY} : \mathbb{R}^2 \rightarrow \mathbb{R}$ sao cho với mọi $C \in \mathbb{R}^2$ ta có

$$P((X, Y) \in C) = \iint_C f_{XY}(x, y) dx dy.$$

Mệnh đề. Hai biến ngẫu nhiên liên tục đồng thời X, Y độc lập khi và chỉ khi $f_{XY}(x, y) = f_X(x)f_Y(y)$ với mọi $x, y \in \mathbb{R}$.

Ôn tập biến ngẫu nhiên liên tục

Kì vọng và phương sai

Cho biến ngẫu nhiên liên tục X với hàm mật độ xác suất f

- **Kì vọng** (mean) của X được tính bởi

$$\mu = E(X) = \int_{-\infty}^{\infty} xf(x)dx,$$

- **Phương sai** (variance) của X được tính bởi

$$\sigma^2 = Var(X) = E((X - \mu)^2) = \int_{-\infty}^{\infty} (x - \mu)^2 f(x)dx,$$

- Với hàm số $r : \mathbb{R} \rightarrow \mathbb{R}$ và $Y = r(X)$

$$E(Y) = E(r(X)) = \int_{-\infty}^{\infty} r(x)f(x)dx.$$

Ôn tập biến ngẫu nhiên liên tục

Phân phối đều

Biến ngẫu nhiên liên tục X được gọi là có **phân phối đều** (uniform distribution) trên $[a, b]$ với $a < b$, kí hiệu $X \sim \mathcal{U}(a, b)$, nếu X có tập giá trị là $[a, b]$ và

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{nếu } a \leq x \leq b, \\ 0 & \text{khác.} \end{cases}$$

Khi đó, X có kì vọng $E(X) = \frac{a+b}{2}$ và phương sai $Var(X) = \frac{(b-a)^2}{12}$.

Gọi X là kết quả của thí nghiệm “chọn **ngẫu nhiên** một điểm trong khoảng $[a, b]$ ” thì $X \sim \mathcal{U}(a, b)$.

Mệnh đề. Cho $X \sim \mathcal{U}(a, b)$ và $d \in (a, b)$, phân phối của X khi biết $X \leq d$ là phân phối đều trên $[a, d]$, thường kí hiệu $(X|X \leq d) \sim \mathcal{U}(a, d)$.

Ôn tập biến ngẫu nhiên liên tục

Phân phối mũ

Biến ngẫu nhiên liên tục X được gọi là có **phân phối mũ** (exponential distribution) với tham số λ ($\lambda > 0$), kí hiệu $X \sim \text{Exp}(\lambda)$, nếu X có tập giá trị là $[0, \infty)$ và

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{nếu } x \geq 0, \\ 0 & \text{khác.} \end{cases}$$

Khi đó, X có kì vọng $E(X) = \frac{1}{\lambda}$, phương sai $\text{Var}(X) = \frac{1}{\lambda^2}$ và hàm phân phối

$$F(x) = 1 - e^{-\lambda x}, x \geq 0.$$

Phân phối mũ có thể được xem như là “phiên bản liên tục” của phân phối hình học.

Mệnh đề (tính không nhớ - memoryless). Cho $X \sim \text{Exp}(\lambda)$, với mọi $t, s \geq 0$ ta có

$$P(X > t + s | X > s) = P(X > t).$$

Ôn tập biến ngẫu nhiên liên tục

Phân phối chuẩn

Biến ngẫu nhiên liên tục X được gọi là có **phân phối chuẩn** (normal distribution) với trung bình μ và phương sai σ^2 ($\sigma > 0$), kí hiệu $X \sim \mathcal{N}(\mu, \sigma^2)$, nếu X có hàm mật độ xác suất

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, x \in \mathbb{R}.$$

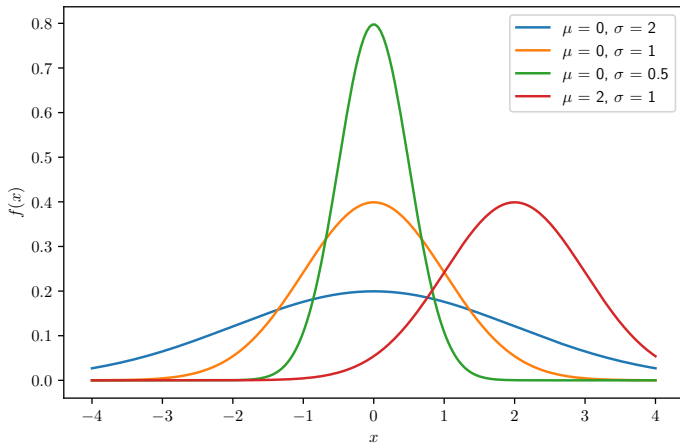
Khi đó, X có kì vọng $E(X) = \mu$ và phương sai $Var(X) = \sigma^2$.

Trường hợp $Z \sim \mathcal{N}(0, 1)$ thì Z được gọi là có **phân phối chuẩn tắc** (standard normal distribution). Hàm mật độ xác suất và hàm xác suất của Z thường được kí hiệu lần lượt là ϕ, Φ , tức là

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, \quad \Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-t^2/2} dt.$$

Ôn tập biến ngẫu nhiên liên tục

Phân phối chuẩn (tt)



Ôn tập biến ngẫu nhiên liên tục

Phân phối chuẩn (tt)

Các tính chất quan trọng của phân phối chuẩn

1. Nếu $X \sim \mathcal{N}(\mu, \sigma^2)$ và $Y = aX + b$ ($a \neq 0$) thì $X \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$,
2. Nếu $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$, $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ và X_1, X_2 độc lập thì $X_1 + X_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$,
3. Nếu $Z \sim \mathcal{N}(0, 1)$ và $X = \sigma Z + \mu$ thì $X \sim \mathcal{N}(\mu, \sigma^2)$,
4. Nếu $X \sim \mathcal{N}(\mu, \sigma^2)$ và $Z = \frac{X - \mu}{\sigma}$ thì $Z \sim \mathcal{N}(0, 1)$, từ đó

$$F_X(x) = P(X \leq x) = P\left(Z \leq \frac{x - \mu}{\sigma}\right) = F_Z\left(\frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right),$$

Luật Zipf và “mô hình chú khỉ”

Hỏi: tại sao tần số với hạng của từ lại có quan hệ theo **luật lũy thừa** (power law) như vậy?

Trả lời: có thể chỉ **do ngẫu nhiên!**

- Li, W. (1992). Random texts exhibit Zipf's-law-like word frequency distribution. *IEEE Transactions on Information Theory*, 38(6), 1842–1845.
- Mô hình: có một chú khỉ ngồi gõ bàn phím từ bộ kí tự có M chữ cái và phím cách. Từ là dãy kí tự phân cách bởi phím cách, chẳng hạn dãy phím gõ `a_mdf__pwell_` tạo ra 3 từ là `a`, `mdf`, `pwell`. Giả sử chú khỉ chưa đi học (nên gõ đại) và rất rảnh rỗi (nên gõ được văn bản rất dài). Trên văn bản mà chú khỉ tạo ra ta cũng thấy luật Zipf giữa tần số với hạng của từ!

Luật Zipf và “mô hình chú khỉ”

Mô phỏng

```
def monkey(N, k, alphabet, space=" "):  
    alphabet += space; words = []; curWord = ""  
    while len(words) < N:  
        letter = random.choice(alphabet)  
        if letter == space:  
            if curWord == "":  
                continue  
            words.append(curWord)  
            curWord = ""  
        else:  
            curWord += letter  
  
    # ...
```

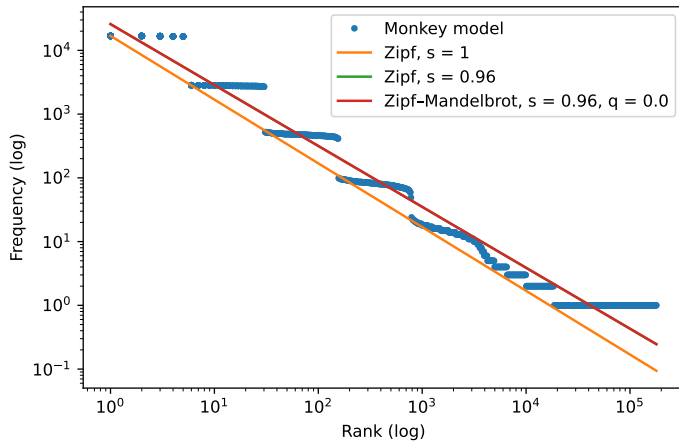
Luật Zipf và “mô hình chú khỉ”

Mô phỏng (tt)

```
def monkey(N, k, alphabet, space=" "):  
    #...  
    word_freqs = collections.Counter(words).most_common()  
    freq = np.array([f for _, f in word_freqs])  
    rank = np.array([int(r) for r in np.logspace(0,  
                                                np.log10(len(freq)), num=k)])  
    return rank, freq[rank - 1]  
  
M = 5 # alphabet size  
N = 500_000 # number of word for simulation  
rank, freq = monkey(N, 1000,  
                    alphabet=string.ascii_lowercase[:M])
```

Luật Zipf và “mô hình chú khỉ”

Kết quả



Phương pháp xấp xỉ phân phối bằng mô phỏng

Biến ngẫu nhiên rời rạc

Để xấp xỉ hàm xác suất f_X của một biến ngẫu nhiên rời rạc X liên quan đến thí nghiệm T , ta có thể dùng phương pháp thống kê như sau

- Thực hiện lặp lại N lần (độc lập) thí nghiệm T và tính các tần suất p_x của biến cố “ X nhận giá trị x ”.
- Khi N đủ lớn, ta có $p_x \approx P(X = x) = f_X(x)$.
- Việc thực hiện lặp lại nhiều lần thí nghiệm T có thể được **mô phỏng** (simulate) trên máy tính.

Phương pháp xấp xỉ phân phối bằng mô phỏng

Biến ngẫu nhiên liên tục

Để xấp xỉ hàm mật độ xác suất f_X của một biến ngẫu nhiên liên tục X liên quan đến thí nghiệm T , ta có thể dùng phương pháp thống kê như sau

- Thực hiện lặp lại N lần (độc lập) thí nghiệm T , ghi nhận các giá trị mà X nhận x_1, x_2, \dots, x_N (còn gọi là **mẫu dữ liệu** - sample).
- Khi N đủ lớn, ta có thể dùng **histogram** hoặc **ước lượng mật độ nhân** (kernel density estimation, KDE) trên mẫu để xấp xỉ f_X .
- Việc thực hiện lặp lại nhiều lần thí nghiệm T có thể được **mô phỏng** (simulate) trên máy tính.

Phương pháp xấp xỉ phân phối bằng mô phỏng

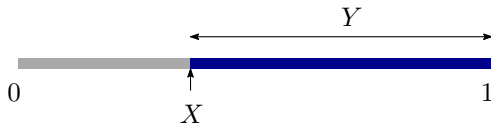
Ví dụ 1

Bài toán. Chọn ngẫu nhiên một điểm trên một thanh có chiều dài 1 đơn vị, cắt tại điểm đó thành hai đoạn và giữ lại đoạn dài hơn. Tính kì vọng và tìm phân phối của chiều dài đoạn giữ lại.

Giải. Gọi X là vị trí ngẫu nhiên chọn trên thanh thì $X \sim \mathcal{U}(0, 1)$. Do đó X là biến ngẫu nhiên liên tục với hàm mật độ xác suất

$$f_X(x) = \begin{cases} 1 & \text{nếu } 0 \leq x \leq 1, \\ 0 & \text{khác.} \end{cases}$$

Gọi Y là chiều dài của đoạn được giữ lại (tức là đoạn dài hơn) thì $Y = \max\{X, 1 - X\}$.



Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 1 (tt)

Ta có kì vọng của chiều dài đoạn giữ lại

$$\begin{aligned} E(Y) &= E(\max\{X, 1 - X\}) = \int_{-\infty}^{\infty} \max\{x, 1 - x\} f_X(x) dx = \int_0^1 \max\{x, 1 - x\} dx \\ &= \int_0^{1/2} \max\{x, 1 - x\} dx + \int_{1/2}^1 \max\{x, 1 - x\} dx = \int_0^{1/2} (1 - x) dx + \int_{1/2}^1 x dx \\ &= \frac{3}{4}. \end{aligned}$$

Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 1 (tt)

Ta tìm hàm phân phối (tích lũy) của biến ngẫu nhiên $Y = \max\{X, 1 - X\}$

$$F_Y(y) = P(Y \leq y) = P(\max\{X, 1 - X\} \leq y), y \in \mathbb{R}.$$

Xét các trường hợp của y

1. $y < 1/2$: $(\max\{X, 1 - X\} \leq y) = \emptyset$ vì $0 \leq X \leq 1$ nên $1/2 \leq \max\{X, 1 - X\}$,

$$P(\max\{X, 1 - X\} \leq y) = P(\emptyset) = 0.$$

2. $1/2 \leq y \leq 1$: $(\max\{X, 1 - X\} \leq y) = (1 - y \leq X \leq y)$,

$$P(\max\{X, 1 - X\} \leq y) = P(1 - y \leq X \leq y) = \int_{1-y}^y f_X(x) dx = \int_{1-y}^y dy = 2y.$$

3. $y > 1$: $(\max\{X, 1 - X\} \leq y) = \Omega$ vì $0 \leq X \leq 1$ nên $\max\{X, 1 - X\} \leq 1$,

$$P(\max\{X, 1 - X\} \leq y) = P(\Omega) = 1.$$

Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 1 (tt)

Từ đó ta có

$$F_Y(y) = \begin{cases} 0 & \text{nếu } y < 1/2, \\ 2y & \text{nếu } 1/2 \leq y \leq 1, \\ 1 & \text{nếu } 1 < y. \end{cases}$$

Lấy đạo hàm của hàm phân phối, ta có hàm mật độ xác suất của Y là

$$f_Y(y) = F'_Y(y) = \begin{cases} 2 & \text{nếu } 1/2 \leq x \leq 1, \\ 0 & \text{khác.} \end{cases}$$

Như vậy Y có phân phối đều trên đoạn $[1/2, 1]$, tức là $Y \sim \mathcal{U}(1/2, 1)$.

Lưu ý, từ phân phối của Y , $Y \sim \mathcal{U}(1/2, 1)$, ta cũng có $E(Y) = \frac{1/2+1}{2} = \frac{3}{4}$.

Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 1 (tt)

```
def greater_len(N):  
    X = np.random.uniform(size=N)  
    Y = np.maximum(X, 1 - X)  
    return Y
```

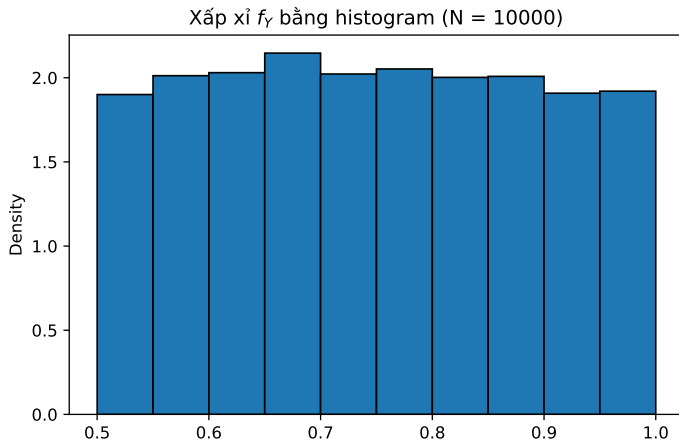
```
N = 10000
```

```
np.mean(greater_len(N))  
#0.7499721269808018
```

```
plt.hist(greater_len(N), density=True, edgecolor="black")
```

Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 1 (tt)



Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 2

Bài toán. Cho X_1, X_2, \dots, X_n là n biến ngẫu nhiên độc lập và cùng phân phối chuẩn $\mathcal{N}(\mu, \sigma^2)$. Đặt

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{và} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

(X_1, \dots, X_n thường được gọi là một mẫu ngẫu nhiên cỡ n , \bar{X} là trung bình mẫu và S^2 là phương sai mẫu.)

Tìm phân phối của các biến ngẫu nhiên $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ và $\frac{\bar{X} - \mu}{S/\sqrt{n}}$.

Trả lời: $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ có phân phối chuẩn tắc $\mathcal{N}(0, 1)$ và $\frac{\bar{X} - \mu}{S/\sqrt{n}}$ có **phân phối Student** (Student's t-distribution) với $n - 1$ bậc tự do.

(https://en.wikipedia.org/wiki/Student%27s_t-distribution.)

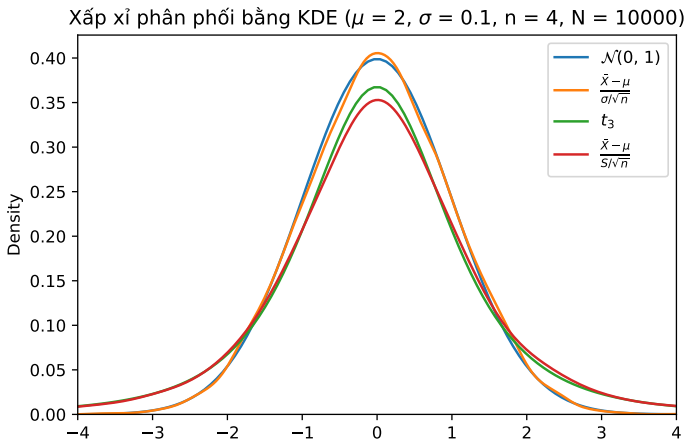
Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 2 (tt)

```
def sample(mu, sigma, n, N):  
    X = np.random.normal(mu, sigma, size=(N, n))  
    X_bar = np.mean(X, axis=1)  
    S2 = np.var(X, axis=1, ddof=1)  
    return X_bar, S2  
  
X_bar, S2 = sample(mu, sigma, n, N)  
  
plt.plot(x, scipy.stats.norm.pdf(x))  
sns.kdeplot((X_bar - mu)/(sigma/np.sqrt(n)))  
plt.plot(x, scipy.stats.t.pdf(x, n - 1))  
sns.kdeplot((X_bar - mu)/(np.sqrt(S2)/np.sqrt(n)))
```

Phương pháp xấp xỉ phân phối bằng mô phỏng

Ví dụ 2 (tt)



Tài liệu tham khảo

Chapter 3-5. Morris H. DeGroot, Mark J. Schervish. *Probability and Statistics*. Addison-Wesley, 2012.

Chapter 3-5. H. Pishro-Nik. *"Introduction to probability, statistics, and random processes"*, available at <https://www.probabilitycourse.com>. Kappa Research LLC, 2014.