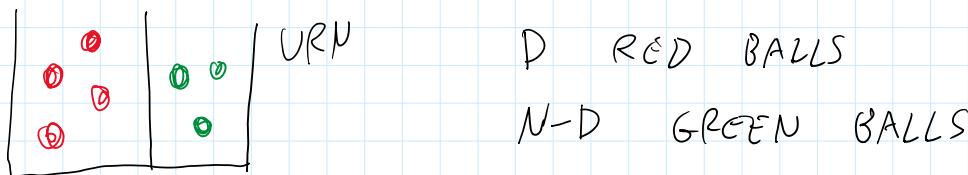


MOODLE - DATA SCIENCE

EXAM - THURSDAY 18 DECEMBER

for DS 16:30 - 19:30 ROOM P300 and LUF1

HYPERGEOMETRIC RANDOM VARIABLES

DRAW n balls

• WITH REPLACEMENT

 $X = \# \text{ of red balls}$

$$X = X_1 + X_2 + \dots + X_n$$

$$X_i = \begin{cases} 1 & \text{if } i\text{-th ball is red} \\ 0 & \text{if } i\text{-th ball is green} \end{cases}$$

$$X_i \sim \text{Ber}(p) \quad p = P(X_i = 1) = \frac{D}{N}$$

$$X_i \perp X_j \text{ for } i \neq j \text{ (independent)}$$

$$X \sim \text{Bin}(n, p) \quad \text{BINOMIAL DISTRIBUTION}$$

• WITHOUT REPLACEMENT

 $Y = \# \text{ of red balls in } n \text{ draws without replacement}$

$$Y \sim \text{Hp}(n, N, D) \quad \text{HYPERGEOMETRIC DISTRIBUTION}$$

Density of Y , $K = 0, \dots, n$

$$p_Y(K) = P(X = K) = \frac{\text{favorable cases}}{\text{possible cases}} = \frac{\binom{P}{K} \binom{N-D}{n-K}}{\binom{N}{n}}$$

$\begin{matrix} K \text{ red} & n-K \text{ green} \\ \textcolor{red}{\bullet} & \textcolor{green}{\bullet} \\ \textcolor{red}{\bullet} & \textcolor{red}{\bullet} \textcolor{green}{\bullet} \textcolor{red}{\bullet} \textcolor{green}{\bullet} \end{matrix} \quad n \text{ balls}$

Expectation of Y ? $E_Y(Y)$?

Consider $Y = Y_1 + Y_2 + \dots + Y_n$

$$\checkmark \quad \checkmark \quad \dots \quad \checkmark \quad \checkmark \quad \checkmark \quad \dots \quad \checkmark \quad \checkmark \quad \dots \quad \checkmark \quad \checkmark \quad \dots \quad \checkmark \quad \checkmark \quad \dots$$

Consider $Y = Y_1 + Y_2 + \dots + Y_n$

Y_i random variable $Y_i = \begin{cases} 1 & \text{if } i\text{-th ball is red} \\ 0 & \text{if } i\text{-th ball is green} \end{cases}$

$Y_i \sim \text{Ber}(p_i)$

$p_i = P(Y_i = 1) = P(\text{i-th ball is red - drawn WITHOUT replacement})$

It can be proved that $p_i = \frac{D}{N} \quad \forall i = 1, \dots, n$

$Y_i \sim \text{Ber}\left(\frac{D}{N}\right)$

but $Y_i \neq Y_j$ for $i \neq j$ (NOT independent)

We can compute $E[Y] = n \frac{D}{N}$

$$E[Y] = E[Y_1] + E[Y_2] + \dots + E[Y_n] = n \cdot \frac{D}{N}$$

$\text{Var}(Y) = ?$

$$\begin{aligned} \text{Var}(Y) &= \text{Var}(Y_1 + \dots + Y_n) \\ &= \text{Var}(Y_1) + \dots + \text{Var}(Y_n) + 2\text{Cov}(Y_1, Y_2) + 2\text{Cov}(Y_1, Y_3) \\ &\quad + \dots + \text{Cov}(Y_1, Y_n) + \dots + \text{Cov}(Y_{n-1}, Y_n) = \dots \end{aligned}$$

$$\text{Var}(Y) = n \frac{D(N-D)}{N^2} \frac{N-n}{N-1} = \text{Var}(X) \frac{n(n-1)}{N-1}$$

recall $\text{Var}(X) = n \text{Var}(X_1) = n \frac{D(N-D)}{N^2}$
if $X \sim \text{Bin}(n, \frac{D}{N})$

$$\text{Note } \frac{N-n}{N-1} < 1 \Rightarrow \text{Var}(Y) < \text{Var}(X)$$

Poisson Distribution

$X_n \sim \text{Bin}(n, p)$

in general, it is difficult to compute density $p_X(k)$ when n is large

Consider $X \sim \text{Poi}(\lambda)$ will be a possible approximation for $\text{Bin}(n, p)$ when n is large, but np is bounded

$$\boxed{\lambda = np}$$

$X: \Omega \rightarrow N$ POISSON r.v. parameter λ $X \sim \text{Poi}(\lambda)$

$$p_X(n) = e^{-\lambda} \frac{\lambda^n}{n!}$$

 $\forall n \in N = \{0, 1, \dots\}$

$$\text{Recall } e^\lambda = \sum_{n=0}^{\infty} \frac{\lambda^n}{n!}$$

$$\text{Hence } \sum_n p_X(n) = 1$$

$$\begin{aligned} n! &= n \underbrace{(n-1)(n-2)\dots 1}_{= (n-1)!} \\ &= n(n-1)! \end{aligned}$$

$$\begin{aligned} E[X] &= \sum_{n=0}^{\infty} n p_X(n) = \sum_{n=0}^{\infty} n e^{-\lambda} \frac{\lambda^n}{n(n-1)!} \\ &= \sum_{n=1}^{\infty} e^{-\lambda} \frac{\lambda \cdot \lambda^{n-1}}{(n-1)!} = \lambda e^{-\lambda} \sum_{n=1}^{\infty} \frac{\lambda^{n-1}}{(n-1)!} \\ &= \lambda e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = \lambda e^{-\lambda} \cdot e^\lambda = \lambda \end{aligned}$$

 $k = n-1$

$E[X] = \lambda$

$V_{\text{var}}(X) = \lambda$

$E[X^2] = \lambda^2 + \lambda$

$= \sum_{n=1}^{\infty} n^2 P(X=n)$

Thm LAW OF SMALL NUMBERS

Let X_1, \dots, X_n independent such that $X_i \sim \text{Bin}(p_i)$ Set $S_n = X_1 + \dots + X_n$ and let $W_n \sim \text{Poi}(p_1 + \dots + p_n)$ Then, for any $A \subseteq N$

$|P(S_n \in A) - P(W_n \in A)| \leq \sum_{i=1}^n p_i^2$

↳ distance between distribution of S_n and W_n

Corollary

If we set $p_i = \frac{\lambda}{n}$ for every $i = 1, \dots, n$

$S_n \sim \text{Bin}(n, \frac{\lambda}{n}) \quad W_n \sim \text{Poi}(\lambda)$

Then

$|P(S_n \in A) - P(W_n \in A)| \leq \sum_{i=1}^n \left(\frac{\lambda}{n}\right)^2 = n \frac{\lambda^2}{n^2} = \frac{\lambda^2}{n}$

$i=1$ (\bar{x}_1)

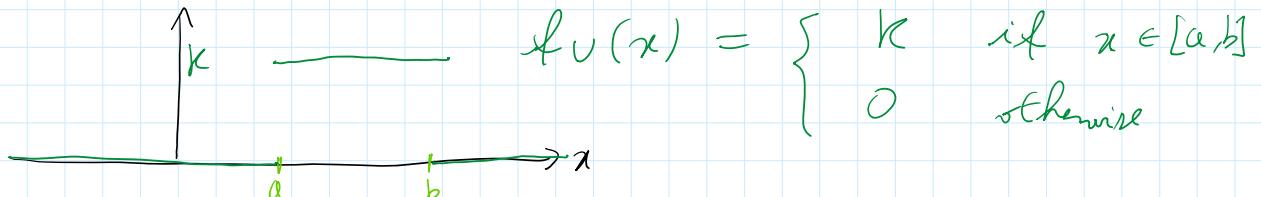
\bar{x}^2

n

ABSOLUTELY CONTINUOUS DISTRIBUTION

• UNIFORM DISTRIBUTION

$$U \sim U(a, b) \quad -\infty < a < b < \infty$$



$$\int_{-\infty}^{+\infty} f_U(x) dx = 1 = \int_a^b K dx = 1 = K(b-a) = 1 \Leftrightarrow K = \frac{1}{b-a}$$

$$f_U(x) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise} \end{cases} \quad U \sim U(a, b)$$

$$f_U(x) = \frac{1}{b-a} \mathbf{1}_{[a,b]}(x) \quad P(a \leq U \leq b) = 1$$

$$\bullet \text{ distribution} \quad P(U \in A) = \int_A f_U(x) dx$$

if $a < c < d < b$

$$P(U \in [c, d]) = P(c \leq U \leq d) = \int_c^d f_U(x) dx$$

$$= \int_c^d \frac{1}{b-a} dx = \frac{d-c}{b-a} = P(c \leq U \leq d)$$

$$\begin{aligned} E[U] &= \int_{-\infty}^{+\infty} x f_U(x) dx = \int_a^b x \frac{1}{b-a} dx \\ &= \frac{1}{b-a} \left[\frac{x^2}{2} \right]_{x=a}^{x=b} = \frac{b^2 - a^2}{2} \cdot \frac{1}{b-a} = \frac{(b-a)(b+a)}{2} \cdot \frac{1}{b-a} \end{aligned}$$

$$E[U] = \frac{a+b}{2} = \text{middle point}$$

EXPONENTIAL n. 5 $X \sim Exp(\lambda)$

$$F_X(x) = P(X \leq x) = \begin{cases} 1 - e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

$$f_X(x) = \lambda e^{-\lambda x} \mathbf{1}_{[0, +\infty)}(x)$$

$$= \Gamma(1) - 1 \quad 1 / \Gamma(1) = 1$$

$$f_X(x) = \lambda e^{-\lambda x} \mathbb{1}_{(0,+\infty)}$$

$$E[X] = \frac{1}{\lambda} \quad \text{Var}(X) = \frac{1}{\lambda^2}$$

- MINIMUM and MAXIMUM of r.v.

Consider X, Y r.v. $X \perp\!\!\!\perp Y$

What is distribution of

$$Z = \max(X, Y) \quad W = \min(X, Y) ?$$

$$[Z(w) = \max(X(w), Y(w))]$$

We compute distribution function of Z and W

- $P(Z \leq z) = F_Z(z) = P(\max(X, Y) \leq z)$

$$\left\{ \max(X, Y) \leq z \right\} = \{X \leq z \text{ and } Y \leq z\}$$

$$= P(X \leq z, Y \leq z) \quad \cap = \{X \leq z\} \cap \{Y \leq z\}$$

INDEPENDENCE $P(X \leq z) P(Y \leq z) = F_X(z) F_Y(z)$

$$\Rightarrow F_{\max(X,Y)}(z) = F_X(z) F_Y(z)$$

- if X and Y were absolutely continuous, the density of $Z = \max(X, Y)$ is $\left[\frac{d}{dz} F_Z(z) = f_Z(z) \right]$

$$f_Z(z) = \frac{d}{dz} F_Z(z) = \frac{1}{\lambda^2} (F_X(z) \cdot F_Y(z))$$

- $P(W > w) = 1 - F_W(w) = P(\min(X, Y) > w)$

$$= P(X > w, Y > w) \stackrel{\text{INDEP.}}{=} P(X > w) P(Y > w)$$

$$F_{\min(X,Y)} = 1 - (1 - F_X(w))(1 - F_Y(w))$$

Example

Let $X \sim \text{Exp}(\lambda)$, $Y \sim \text{Exp}(\mu)$, $X \perp\!\!\!\perp Y$

What is distribution of $\min(X, Y)$

$$\Rightarrow \min(X, Y) \sim \text{Exp}(\lambda + \mu)$$

$$P(X > x) = \begin{cases} e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$P(X > x) = \begin{cases} e^{-\lambda x} & \text{if } x \geq 0 \\ 1 & \text{if } x < 0 \end{cases}$$

$$\bullet P(\min(X, Y) > z) = P(X > z) P(Y > z) \quad \text{if } z \geq 0$$

$$= e^{-\lambda z} \cdot e^{-\mu z} = e^{-(\lambda + \mu)z}$$

$$\bullet P(\min(X, Y) > z) = 1 \quad \text{if } z < 0$$

$$\Rightarrow \min(X, Y) \sim \text{Exp}(\lambda + \mu)$$

because we computed the distribution function

Thm A random variable $T : \Omega \rightarrow [0, +\infty)$ has exponential distribution if and only if it has the following MEMORYLESS PROPERTY

Ⓐ $P(T > s + t \mid T > s) = P(T > t) \quad \forall s, t > 0$
 (and it holds $P(T > t) = e^{-\lambda t}$)

Proof

(\Rightarrow) $T \sim \text{Exp}(\lambda)$. We show that Ⓚ holds.

$$P(T > t) = e^{-\lambda t}$$

$$P(T > s + t \mid T > s) = \frac{P(T > s + t, T > s)}{P(T > s)}$$

$$\{T > s + t\} \subseteq \{T > s\}$$

$$T > s + t > s \Rightarrow \omega \in \{T > s + t\} \Rightarrow \omega \in \{T > s\}$$

$$\text{then } \{T > s + t\} \cap \{T > s\} = \{T > s + t\}$$

$$= \frac{P(T > s + t)}{P(T > s)} = \frac{e^{-\lambda(s+t)}}{e^{-\lambda s}} = e^{-\lambda t}$$

(\Leftarrow) Given T r.v satisfying Ⓚ,

we have to show that $T \sim \text{Exp}(\lambda)$, that is, there exists λ such that $P(T > t) = e^{-\lambda t}$

there exists λ such that $P(T > \epsilon) = e^{-\lambda \epsilon}$

$$P(T > s + \epsilon | T > s) = \frac{P(T > s + \epsilon, T > s)}{P(T > s)} = \frac{P(T > s + \epsilon)}{P(T > s)}$$

(1)

$$= P(T > \epsilon)$$

If we call $g(\epsilon) = P(T > \epsilon)$, function of ϵ , we derive the equality

$$g(s + \epsilon) = g(\epsilon) g(s)$$

CAUCHY FUNCTIONAL INEQUALITY

If g is continuous, the unique solutions to this equation are the exponential

