

## 제 0 강 통계학 Review

### Part I. 확률론 (Probability Theory)

#### I. 확률변수 (Random Variable)와 확률분포

##### A. 확률변수 $X$ 는 표본공간 $\Omega$ 상에서 정의되는 **real valued function** 임

i. 어떤 확률적 실험의 결과로 나올 수 있는 모든 가능한 결과에 대해 어떤 실수값이 대응되어야 함

ii. 하나의 실험에 대해 여러 가지의 확률변수가 정의될 수 있음

1. 주사위 던지는 실험: 던진 결과 나오는 값을 대응시켜주는 확률변수  $X$ , 짝수는 0 홀수는 1 을 대응시켜주는 확률변수  $Y$  등등

##### B. 확률(밀도)함수 (Probability (density) Function)

i. 확률함수 : 확률변수  $X$  가 취하는 값이 이산적(discrete)인 경우,  $X$  가 취하는 각 값에 대해 취하는 확률을 대응시켜주는 함수

ii. 확률밀도함수 : 확률변수  $X$  가 취하는 값이 연속적(continuous)한 경우,  $X$  가 취하는 값이  $I$  라는 구간에 속할 확률이  $P(X \in I) = \int_I f(x)dx$  로 주어질 때,  $f(x)$  를 확률밀도함수라 함

C. (누적) 분포함수 ((Cumulative) Distribution Function): 확률변수  $X$  의 분포함수 :  $F(x) \equiv P(X \leq x)$

D. 결합 확률(밀도)함수 (Joint Probability (density) Function):

i. 두 개의 확률변수  $X, Y$  가 이산적이 경우 결합확률함수는

$$f_{xy}(x_i, y_j) \equiv P(X = x_i, Y = y_j) \text{ 와 같이 주어지며}$$

ii. 두 개의 확률변수  $X, Y$  가 연속적이 경우 결합확률밀도함수

$$f_{xy}(x, y) \text{ 는 } P((X, Y) \in D) = \iint_D f_{xy}(x, y) dx dy \text{ 와 같이 주어진다.}$$

E. 한계 확률(밀도)함수 (Marginal Probability (density) Function)

i.  $f_x(x_i) = \sum_j f_{xy}(x_i, y_j), i = 1, 2, \dots$  :  $X$  의 한계확률함수

ii.  $f_x(x) = \int_{-\infty}^{\infty} f_{xy}(x, y) dy$  :  $X$  의 한계확률밀도함수

F. 조건부 확률(밀도)함수 (Conditional Probability (density) Function)

i. 이산적인 경우:  $Y = y_j$  로 주어졌을 때,  $X$  의 조건부확률함수는

$$f_{x|y}(x_i | y_j) \equiv P(X = x_i | Y = y_j) = \frac{P(X = x_i, Y = y_j)}{P(Y = y_j)} = \frac{f_{xy}(x_i, y_j)}{f_y(y_j)} = \frac{f_{xy}(x_i, y_j)}{\sum_i f_{xy}(x_i, y_j)}$$

와 같이 주어짐

- ii. 연속적인 경우:  $Y=y$  로 주어졌을 때,  $X$  의 조건부확률밀도함수

$$f_{x|y}(x|y) = \frac{f_{xy}(x,y)}{f_y(y)} = \frac{f_{xy}(x,y)}{\int_{-\infty}^{\infty} f_{xy}(x,y) dx} \text{ 와 같이 주어짐}$$

### G. 수학적 기대값 (Mathematical Expectation)

- i.  $g(X)$  가 확률변수  $X$  의 함수라고 할 때,

$$E(g(X)) \equiv \sum_i g(x_i) f(x_i) \quad (\text{이산적인 경우}),$$

$$E(g(X)) \equiv \int_{-\infty}^{\infty} g(x) f(x) dx \quad (\text{연속적인 경우})$$

1.  $g(X) = X$  일 때,  $E(X)$  를 확률변수  $X$  의 평균(보통  $\mu_X$  또는  $\mu$  로 표기)이라고,

2.  $g(X) = (X - E(X))^2$  일 때,  $V(X) \equiv E[(X - E(X))^2]$  을 확률변수  $X$  의 분산 (보통  $\sigma_X^2$  또는  $\sigma^2$  로 표기)이라고 함

a.  $\sigma_X \equiv \sqrt{V(X)}$  :  $X$  의 표준편차(Standard Deviation)

- ii.  $g(X, Y)$  가 확률변수  $X, Y$  의 함수라고 할 때,

$$E(g(X, Y)) \equiv \sum_i \sum_j g(x_i, y_j) f(x_i, y_j) \quad (\text{이산적인 경우}),$$

$$E(g(X,Y)) \equiv \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y)f(x,y)dxdy \quad (\text{연속적인 경우})$$

1.  $g(X,Y) = (X - \mu_X)(Y - \mu_Y)$  일 때,  $E[(X - \mu_X)(Y - \mu_Y)]$  를  $X$  와  $Y$  의 공분산(Covariance) 이라고함 ( $Cov(X,Y)$ 로 표기)

iii. 수학적 기대값과 관련된 사항 몇 가지

$$1. \quad V(X) = E(X^2) - [E(X)]^2, \quad Cov(X,Y) = E(XY) - E(X) \cdot E(Y)$$

$$2. \quad E(c) = c, \quad V(c) = 0, \quad c \text{ 가 상수일 때,}$$

$$3. \quad E(aX + bY + c) = aE(X) + bE(Y) + c, \quad ,$$

$$V(aX + bY + c) = a^2V(X) + b^2V(Y) + 2abCov(X,Y) \quad (a, b, c \text{ 는 상수})$$

$$4. \quad E\left(\sum_i a_i X_i\right) = \sum_i a_i E(X_i), \quad V\left(\sum_i a_i X_i\right) = \sum_i \sum_j a_i a_j Cov(X_i, X_j) \quad a_i \text{ 는 상수, 단 } Cov(X_i, X_i) = V(X_i)$$

$$5. \quad \rho_{XY} \equiv \frac{Cov(X,Y)}{\sqrt{V(X)} \cdot \sqrt{V(Y)}} : \text{ 확률변수 } X, Y \text{ 간의 상관계수(correlation coefficient)}$$

## H. 확률적 독립 (Stochastic Independence)

- i. 일반적으로  $n$  개의 확률변수  $X_1, X_2, \dots, X_n$  가 서로 확률적으로 독립(mutually Stochastically Independent)이라는 것은 각 확률변수의 한계확률(밀도)함수를  $f_1(x_1), f_2(x_2), \dots, f_n(x_n)$  이라고 할 때,
- $X_1, X_2, \dots, X_n$  의 결합확률밀도함수가
- $$f(x_1, x_2, \dots, x_n) = f_1(x_1) \times f_2(x_2) \times \dots \times f_n(x_n) = \prod_i f_i(x_i) \quad \text{와} \quad \text{같이}$$
- 한계확률밀도함수의 곱으로 나타나는 경우로 정의된다.

### ii. 확률적 독립과 관련된 몇 가지 사항

1. 확률변수  $X_1, X_2, \dots, X_n$  가 서로 확률적으로 독립일 경우,

a.  $E\left\{\prod_{i=1}^n u_i(X_i)\right\} = \prod_{i=1}^n E(u_i(X_i)), \text{ ex.) } E(X_1 X_2) = E(X_1)E(X_2)$

- b.  $Cov(X_i, X_j) = 0, \rho_{X_i X_j} = 0$  (그러나 이 경우 逆은 일반적으로 성립하지 않는다)

c.  $Var\left(\sum_i a_i X_i\right) = \sum_i a_i^2 V(X_i), \text{ ex) } V(X_1 \pm X_2) = V(X_1) + V(X_2)$

## II. 이산적 확률분포들 (Discrete Probability Distributions)

A. 베르누이 실험 (Bernouilli Experiment) : 가능한 결과가 오직 두 가지로만 나타나는 실험 (편의상 그 두 가지 결과 중 하나를 성공(S)로 나머지 하나를 실패(F)로 명명.

i. 확률변수  $X$  를 한 번의 베르누이 실험에서 나오는 성공의 횟수라고 하고, 성공의 확률을  $p$  라고 하면  $X$  의 확률분포는 다음과 같다.

$$1. \quad P(X=1)=p, \quad P(X=0)=1-p \equiv q \quad \text{또는} \quad f(i)=p^i q^{1-i}, \quad i=0,1$$

$$2. \quad E(X)=p, \quad V(X)=pq$$

확률분포명 및 표기	확률변수	확률함수	평균, 분산	기타
Binomial 분포 (이항분포)  $\mathcal{B}(n,p)$	N 번의 독립적 베르누이 실험시, 성공의 횟수	$f(i)=\binom{N}{i} p^i q^{N-i}, \quad i=0,1,2,\dots,N$	$E(X)=Np$ $V(X)=Npq$	
Poisson 분포	시간당 $\psi$ 의 비율로	$f(i)=\frac{\lambda^i e^{-\lambda}}{i!}, \quad i=0,1,2,\dots$	$E(X)=\lambda$	$\lambda_1, \lambda_2$ 의 모수를

$\mathcal{P}(\lambda)$ $\lambda = T \times \psi$	<p>물고기가 잡힐 때, T 시간 동안 잡힐 물고기의 수</p> <p>n 이 충분히 크고 p 가 충분히 작을 때, <math>\mathcal{B}(n, p)</math>는 <math>\mathcal{P}(np)</math>로 근사</p>		$V(X) = \lambda$	<p>갖는 두 개의 독립적인 Poisson 분포의 합은 <math>\lambda_1 + \lambda_2</math>의 모수를 갖는 포아송분포를 합</p>
<p>Multinomial 분포</p>	<p>각 실험이 3 개(혹은 그 이상)의 결과(A, B, C)가 가능한 실험을 독립적으로 행할 때, A 의 횟수(X)와 B 의 횟수(Y)</p>	$f(i, j) = \binom{n}{i, j, k} p_X^i p_Y^j p_Z^k,$ $i = 0, 1, \dots, n, j = 0, 1, \dots, n - i$ <p>단, <math>k = n - i - j, p_Z = 1 - p_X - p_Y</math></p>	$Cov(X, Y) = -np_X p_Y$	<p>X 와 Y 의 한계확률분포는 각각 이항분포 <math>\mathcal{B}(n, p_X), \mathcal{B}(n, p_Y)</math>임</p>

### III. 연속적 확률분포들 (Continuous Probability Distributions)

#### A. 기본 분포들

확률분포명 및 표기	확률밀도함수, 분포함수	평균, 분산	기타
Uniform 분포 $\mathcal{U}(a,b)$	$f(x) = \frac{1}{b-a}, a < x < b, \quad \text{zero elsewhere.}$ $F(x) = 0, x \leq a,$ $= \frac{x-a}{b-a}, a < x < b$ $= 1, x \geq b$	$E(X) = \frac{a+b}{2}$ $V(X) = \frac{(b-a)^2}{12}$	
Exponential 분포 $\mathcal{E}(\beta)$ , 단, $\beta$ 는 물고기 한 마리를 낚는 데 걸리는 시간	$f(x) = \frac{1}{\beta} e^{-\frac{x}{\beta}}, x > 0, \text{ zero elsewhere}$ $F(x) = 0, x \leq 0,$ $= 1 - e^{-\frac{x}{\beta}} x > 0$	$E(X) = \beta$ $V(X) = \beta^2$	물고기 한마리를 낚는 데 걸리는 시간이 $\beta$ (또는 시간당 $\psi \equiv 1/\beta$ 의 비율로 물고기가 낚임)일 때 첫 물고기를 낚을 때까지 걸리는 시간이 Exponential 분포
Normal 분포(정규분포) $\mathcal{N}(\mu, \sigma^2)$	$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, -\infty < x < \infty$	$E(X) = \mu$ $V(X) = \sigma^2$	$X_1 \sim N(\mu_1, \sigma_1^2),$ $X_2 \sim N(\mu_2, \sigma_2^2)$ $aX_1 + bX_2$ $\sim N(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2 + 2ab\sigma_{12})$



B. 변환을 통해 도출되는 분포들

확률분포명 및 표기	변환	확률밀도함수	평균, 분산	기타
(표준) Cauchy 분포 $\mathcal{C}(0,1)$		$f(x) = \frac{1}{\pi[1+x^2]}, \quad -\infty < x < \infty$	0 을 중심으로 대칭인 분포이나, 꼬리 부분이 표준정규분포에 비해 두꺼우며, 적률이 존재하지 않는다.	
Chi-square(카이제곱) 분포 $\chi_n^2$ :자유도가 인 카이제곱분포	$X = \sum_{i=1}^n Z_i^2$ $Z_1, \dots, Z_n$ 은 서로 독립인 표준정규분포	$f(x) = \frac{x^{\frac{n}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)}, \quad x > 0, \quad \text{zero elsewhere}$ (굳이 외울 필요는 없음)	$E(X) = n$ $V(X) = 2n$	
(Student) t 분포 $t_n$ :자유도가 인 n 분포	$X = \frac{Z}{\sqrt{V/n}}$ $Z \sim N(0,1),$ $V \sim \chi_n^2,$ $Z, V$ 는 서로 확률적 독립	$f(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\sqrt{n\pi}} \frac{1}{\left(1+\frac{x^2}{n}\right)^{\frac{n+1}{2}}},$ $-\infty < x < \infty$ (굳이 외울 필요는 없음)	$E(X) = 0, \quad (n \geq 2 \text{ 인 경우에만 존재})$ $V(X) = \frac{n}{n-2}, \quad (n \geq 3 \text{ 인 경우에만 존재})$	자유도가 1 인 t 분포는 표준 Cauchy 분포가 됨

<p>F-분포</p> <p><math>F_{n,m}</math>: 분자의 자유도가 n 이고 분모의 자유도가 m 인 F 분포</p>	$X = \frac{V_1/n}{V_2/m}$ <p><math>V_1 \sim \chi_n^2</math>,  <math>V_2 \sim \chi_m^2</math> 는 서로 확률적 독립</p>	$f(x) = \frac{\Gamma(\frac{n+m}{2})}{\Gamma(\frac{n}{2})\Gamma(\frac{m}{2})} \left(\frac{n}{m}\right)^{\frac{n}{2}} \frac{x^{\frac{n}{2}-1}}{(1+\frac{n}{m}x)^{\frac{n+m}{2}}}$ <p>, <math>x &gt; 0</math>, zero elsewhere  (균이 외울 필요는 없음)</p>	$E(X) = \frac{m}{m-2}, (n \geq 3 \text{ 인 경우에만 존재})$ $V(X) = \frac{2m^2(n+m-2)}{(m-2)^2(m-4)},$ <p>, ( <math>n \geq 5</math> 인 경우에만 존재) (균이 외울 필요는 없음)</p>	<p>X 가 <math>F_{n,m}</math> 인 경우 <math>1/X</math> 는 <math>F_{m,n}</math> 임</p> <p><math>t_n^2 = F_{1,n}</math></p>
--	--	--	--	--

C. 기타 확률분포들

확률분포 명 및 표기	변환	확률밀도함수	평균, 분산	기타
Lognormal 분포	$Y = \exp(X)$ , $X \sim N(\mu, \sigma^2)$	$f(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{\ln y - \mu}{\sigma}\right)^2} \frac{1}{y}, y > 0$		
Logistic 분포		$f(x) = \frac{e^{-x}}{(1+e^{-x})^2}, -\infty < x < \infty,$ $F(x) = \frac{1}{1+e^{-x}}$	0 을 중심으로 대칭인 분포이나, 꼬리 부분이 표준정규분포에 비해 두꺼움	