

7.2 기댓값과 확률변수의 변환

문제 7.2.1

확률변수의 기댓값

확률변수의 확률밀도함수 \longrightarrow 확률변수의 이론적 평균값 = 기댓값 (expectation)

- 확률변수 X 의 기댓값을 구하는 연산자 (operator): $E[X]$
- 확률변수 X 의 기댓값: μ_X 또는 μ

* 이산확률변수의 기댓값 = 표본공간의 원소 x_i 의 가중평균
(가중치는 x_i 가 나올 확률, 즉 pmf $p(x_i)$)

$$\mu_X = E[X] = \sum_{x_i \in \Omega} x_i p(x)$$

sample(X).
|번씩만 나오는 경우의 수 (0)

문제 7.2.1

공정한 동전이 있고, 이 동전의 앞면이 나오면 1, 뒷면이 나오면 0인 확률변수 X 가 있다.
 $E[X]$ 를 구하라.

$$1 \cdot \frac{1}{2} + 0 \cdot \frac{1}{2} = \frac{1}{2}$$

(참고: 데이터 공간에서 기댓값에 대응하는 값인 표본평균을 구하는 공식)

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \text{ sample.}$$

기댓값 공식과 표본평균 공식에서 x_i 의 의미가 다른 점이 유의할 것

- 기댓값 공식에서 x_i 는 표본공간의 모든 원소를 뜻하지만,
표본평균 공식에서 x_i 는 선택된 (sampled, realized) 표본만을 뜻한다.

연습문제 7.2.2

기대값을 구하는 공식에서는 확률을 가중치로 곱한다. 2번에 연 공분배를 구하는 공식에서는 확률 가중치가 없나?

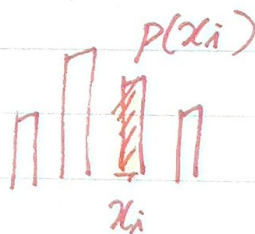
* 연속확률변수의 기대값은 확률밀도함수 $p(x)$ 를 가중치 하여 모든 가능한 표본 x 를 적분한 값이다.

$$\mu_x = E[X] = \int_{-\infty}^{\infty} x p(x) dx$$

기대값은 여러 가능한 x 값을 확률 (또는 확률밀도) 값에 따라 가중치를 한 것이므로 가장 확률 (또는 확률밀도) 이 높은 x 값 근처의 값이 된다. 즉, 확률 (또는 확률밀도) 가 모여있는 곳의 위치를 나타낸다.

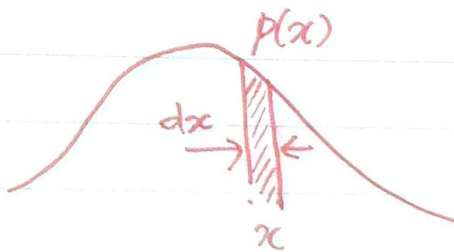
① discrete X

$$E[X] = \sum_{x_i \in \Omega} x_i \underbrace{p(x_i)}_{\text{pmf}}$$



② continuous X

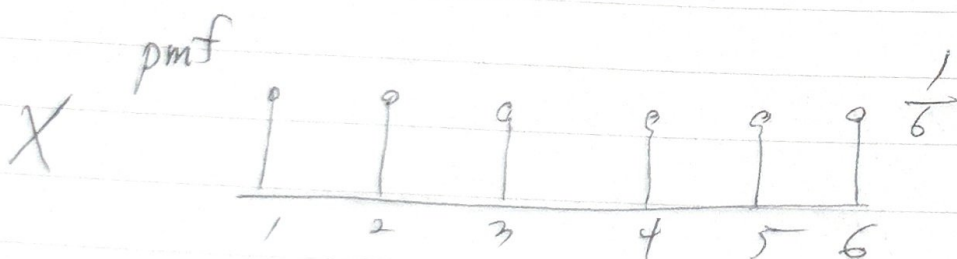
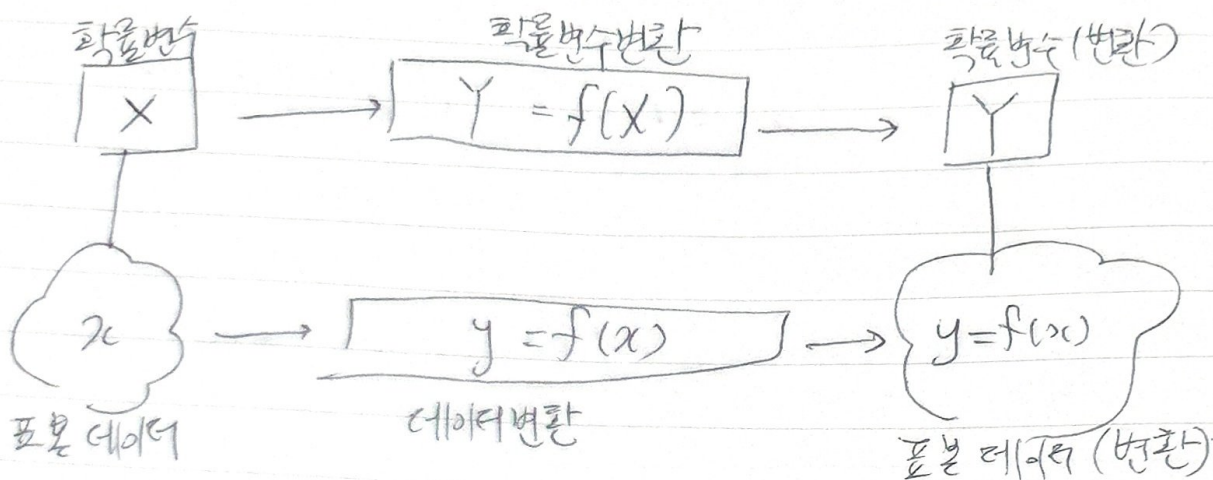
$$E[X] = \int_{-\infty}^{\infty} x \underbrace{p(x)}_{\text{pdf}} dx$$



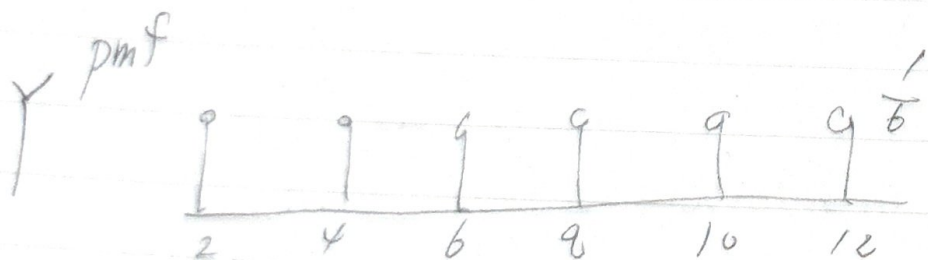
(확률변수의 변환)

우리가 얻은 데이터의 값을 어떤 함수 f 에 넣어서 변환시킨다고 가정하자. 그러면 새로운 데이터 집합이 생긴다.

$$\{x_1, x_2, \dots, x_n\} \rightarrow \{f(x_1), f(x_2), \dots, f(x_n)\}$$



fair dice
(1~6)



fair dice x2.
(2; 4, 6, 8, 10, 12)

확률변수 X 에서 표본을 N 번 뽑아서 2 값을 더하는 경우에는, 다음처럼 원래 확률변수(복사본 X_1, X_2, \dots, X_N)를 만든 다음 이 복사본 확률변수의 표본값을 더한 형태로 변환식을 써야 한다.

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

이렇게 복사본을 만들어 첨가를 붙이는 이유는 X_1 과 X_2 가 같은 확률변수를 가지는 확률변수지만 표본값이 다르기 때문이다.

만약 다음과 같이 쓰면,

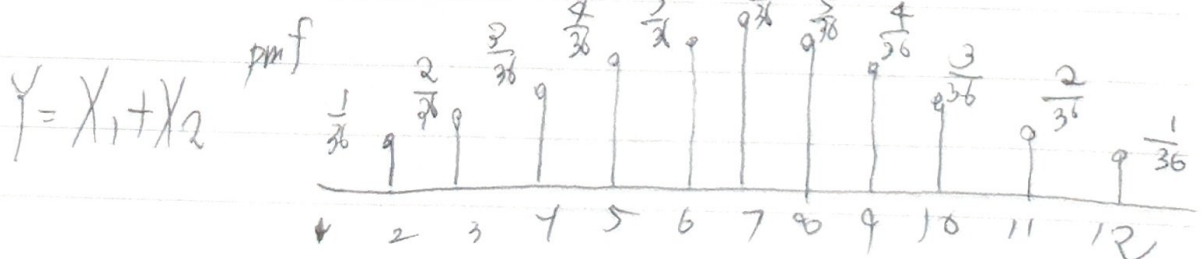
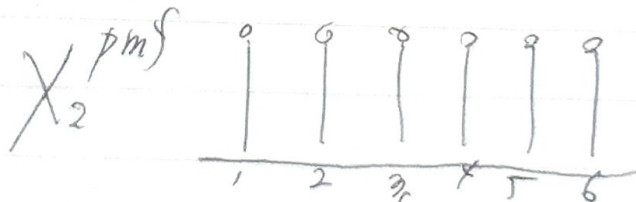
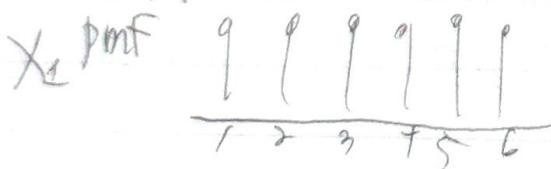
$$Y = X + X + \dots + X$$

이 식은 다음처럼 전혀 다른 확률변수를 가리킨다.

$$Y = N \cdot X \neq (X_1 + X_2 + \dots + X_N)$$

연습문제 7.2.5

확률변수 X_1 과 X_2 는 각각 주사위를 던져 나오는 수를 나타내는 확률변수다. 그리고 Y 는 두 주사위를 동시에 던져 나오는 수의 합을 나타내는 확률변수다. 확률변수 X_1, X_2, Y 의 확률질량함수의 그래프를 각각 그려라.



기대값의 성질

상수 c 에 대해.

$$E[c] = c$$

<선형성>

$$E[cX] = cE[X]$$

$$E[X+Y] = E[X] + E[Y]$$

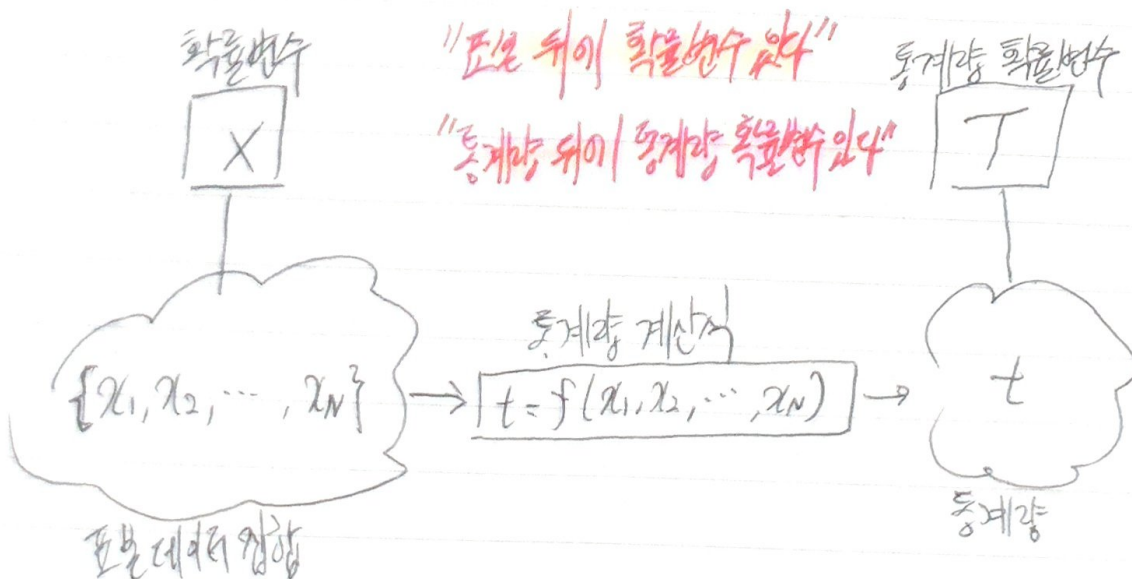
$$E[c_1X + c_2Y] = c_1E[X] + c_2E[Y]$$

통계량 Statistics

확률변수 X 로부터 데이터집합 $\{x_1, x_2, \dots, x_N\}$ 을 얻었다고 하자.

이 데이터 집합 (sample)의 모든 값을 정해진 어떤 공식에 넣어서
해산의 숫자를 구한 것을 통계량 (statistics) 이라고 한다.

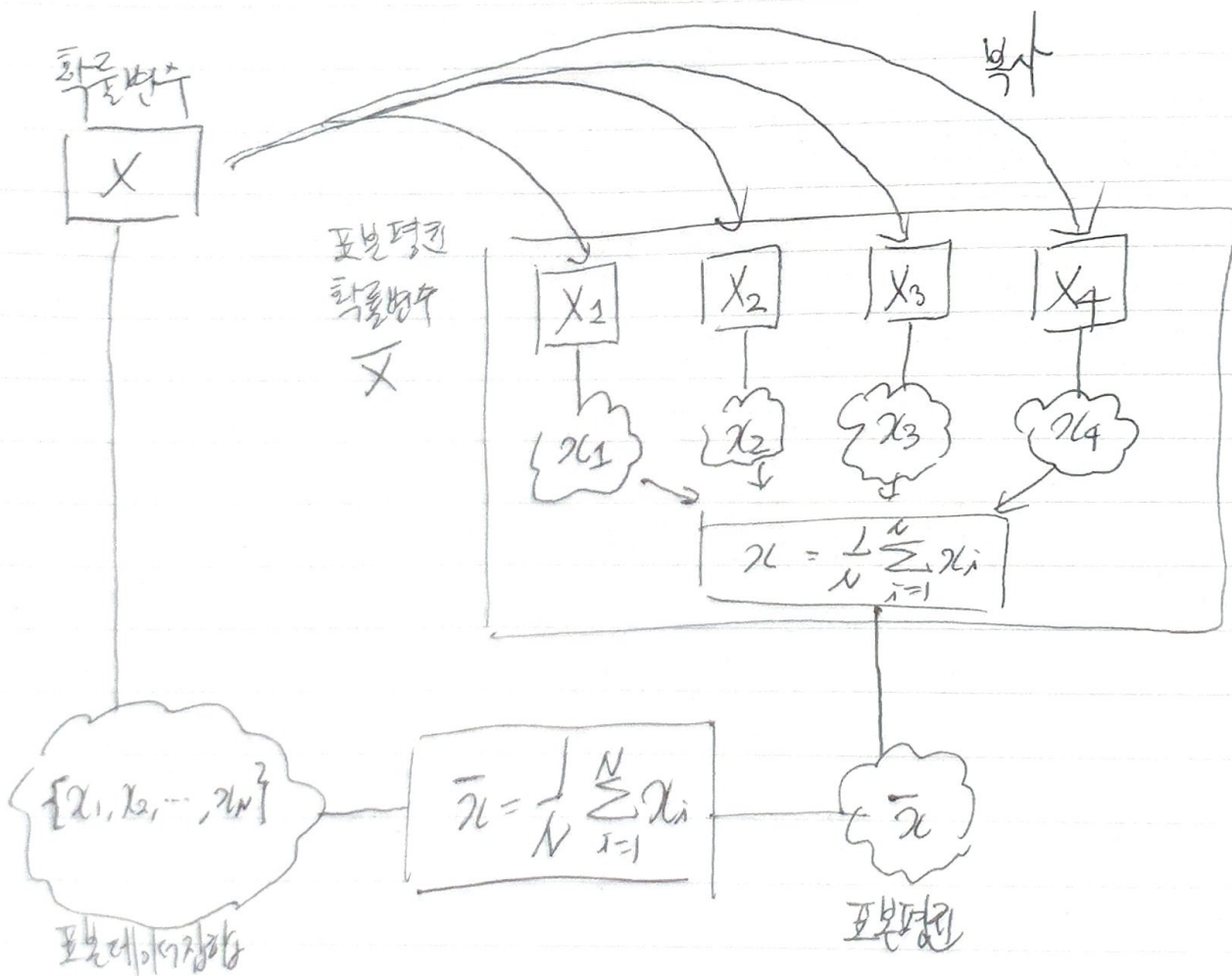
예를 들어 표본평균, 표본분산, 표본중앙값, 표본분산 등은 모두 통계량이다.
통계량도 확률변수의 범위에 포함된다.



표본평균 확률변수

확률변수로부터 N 개의 표본을 만들어 이 표본집합의 표본평균을 구하면 이렇게 구한 표본평균 값도 확률변수가 된다. 표본평균 확률변수의 원래의 확률변수 X 의 값을(bar)을 취하여 \bar{X} 와 같이 표기한다.

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$



(선습문제 7.2.6)

표본평균 \bar{x} 의 값은 확률적인 데이터이고, 이를 생성하는 확률변수 X 는 기타 값이 정의할 수 없었다. 그렇다면 (편향) 샘플분산 S^2 의 값은 확률적인 데이터인가? 만약 그렇다면 이를 생성하는 확률변수 S^2 은 어떻게 정의해야 하는가?

$$\bar{x} = \frac{1}{N} \sum x_i$$

$$S^2 = \frac{1}{N} \sum (x_i - \frac{1}{N} \sum x_i)^2$$

기댓값과 표본평균과의 관계

표본평균도 확률변수이고 기댓값이 존재한다. 표본평균의 기댓값은 원래의 확률변수의 기댓값과 같다는 것을 다음처럼 증명할 수 있다.

$$E[\bar{x}] = E[X]$$

(증명)

표본평균은 표본에서 확률변수들이 모인 값 = 원래의 확률변수들에서 확률변수들이 모인 값

$$E[\bar{x}] = E\left[\frac{1}{N} \sum_{i=1}^N x_i\right]$$

$$= \frac{1}{N} \sum_{i=1}^N E[x_i]$$

$$= \frac{1}{N} \sum_{i=1}^N E[X]$$

$$= \frac{1}{N} \cdot N \cdot E[X]$$

$$= E[X]$$

x_i 는 X 의 copy이고 기댓값은 같다.

표본평균은 확률변수의 기댓값 근처의 값이 된다.

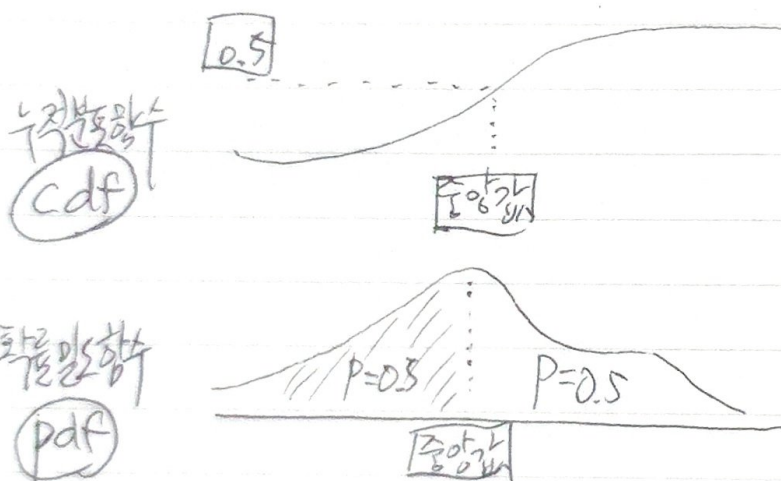
$$E[\text{fair die}] = 3.5 \rightarrow E[\text{fair dice}] \approx 3.5$$

중앙값

확률변수의 중앙값 (median)은 중앙값보다 큰 값이 나올 확률과 작은 값이 나올 확률이 0.5로 같은 값을 뜻한다. 따라서 다음과 같이 누적확률분포 $F(x)$ 에서 중앙값을 계산할 수 있다.

$$0.5 = F(\text{중앙값})$$

$$\text{중앙값} = F^{-1}(0.5)$$



최빈값

most frequent value.

- 이산확률분포에서 — "가장 확률값이 큰 수" = 최빈값 (mode)
- 연속확률분포에서 — "확률밀도함수 pdf $p(x)$ 의 값이 가장 큰 확률변수의 값" (확률밀도함수의 최대값의 위치)

$$\text{최빈값} = \arg \max_x p(x)$$