

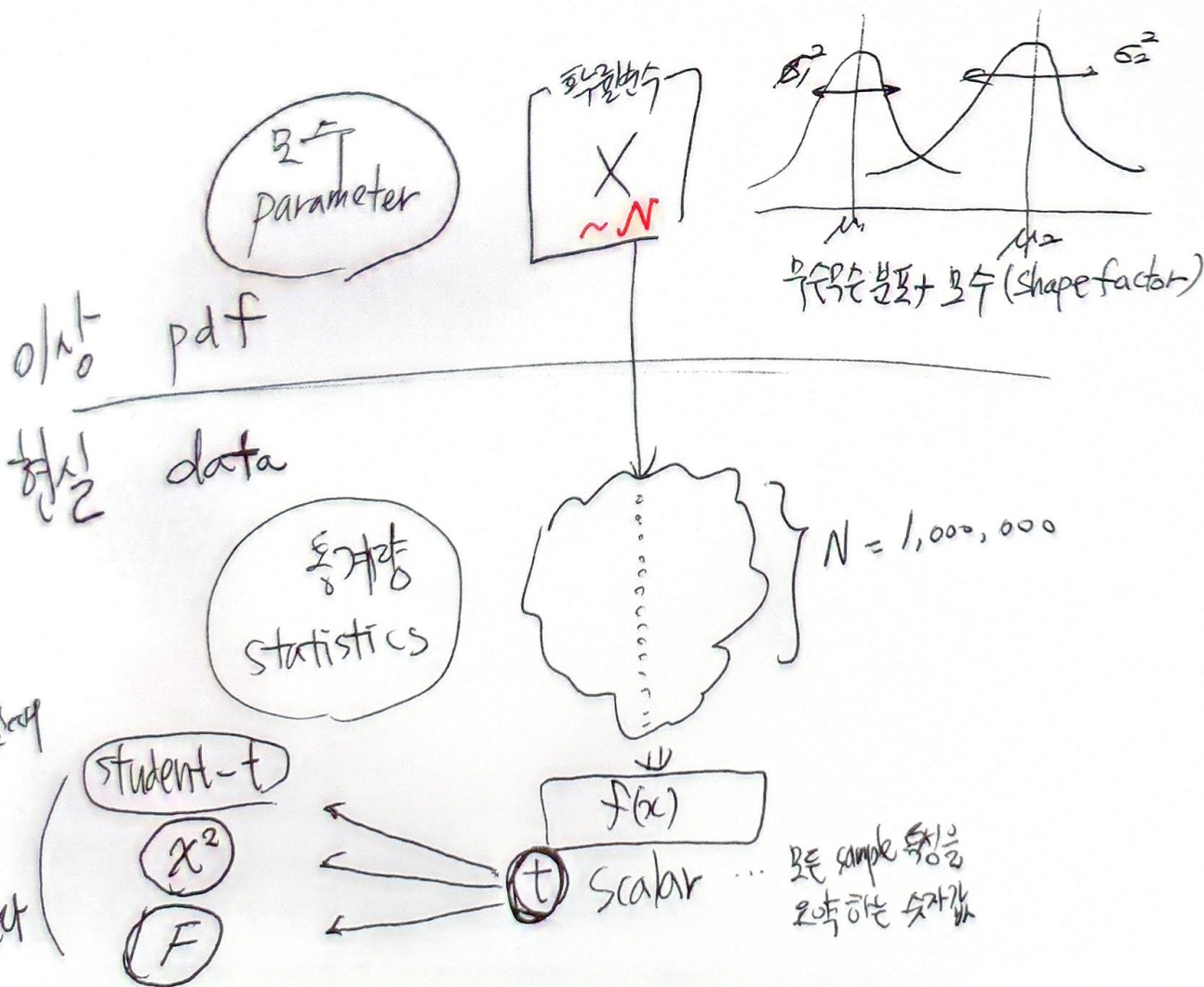
# 8.5 스튜던트 t 분포, 카이제곱 분포, F 분포

- 정규분포에서 파생된 분포들  $\rightarrow$  스튜던트 t, 카이제곱, F
- 위 파생 분포는 통계학 분포라고 부리며, 나중에 공부할 가설 검정이 쓰인다.

## 스튜던트 t 분포

현실 데이터 중.. (정규분포와 상당히 유사하지만, 양 끝단의 비중이)  
 정규분포에 비해 더 큰 데이터  
 = 정규분포보다 극단적 현상이 더 자주 발생했다는 의미  
 = fat tail

EX) 주식 수익률은 보통 정규분포를 따르는 것으로 가정  
 But 정규분포에서는 거의 발생할 수 없는 극단적 사건 발생  
 = black swan



스튜던트 t 분포 (or t 분포) 확률밀도함수

$$t(x; \mu, \lambda, \nu) = \frac{\sqrt{\lambda} \Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left( 1 + \lambda \frac{(x-\mu)^2}{\nu} \right)^{-\frac{\nu+1}{2}}$$

\*  $\lambda = \text{정규분포의 정밀도 } (\sigma^2)^{-1} = (\beta \text{ 이 대응하는 개념})$

\*  $\Gamma(x) = \text{gamma function}$

$$\Gamma(x) = \int_0^{\infty} u^{x-1} e^{-u} du$$

정규분포와 달리, 정수값을 가지는 자유 (degree of freedom) 라는 변수 (parameter)  $\nu$  를 추가적으로 가진다. 스튜던트 t 분포에서는 변수  $\nu$  로 2 이상의 자연수를 사용한다.

변수  $\nu = 1 \rightarrow$  코시 분포 (Cauchy distribution)

코시 분포에서 양수인 부분만 사용  $\rightarrow$  하트코시 분포 (Half-Cauchy distribution)

t 통계량

(previous) 정규분포의 표본을 표본평균으로 나누어 정규화한 t 통계량은 항상 정규분포가 된다.

$\rightarrow$  BUT t 통계량을 구하려면 표본분산의 정확한 표본평균을 알아야 한다.

--- 현실적으로 불가능!

--- 표본에서 측정한 표본표준편차 (sample standard deviation)  
으로 정규화할 수밖에 없음

t 통계량 =  $N$  개 표본  $x_1, \dots, x_N$  에서 계산한 표본평균을  
표본표준편차로 정규화한 값



t-통계량은 자유도가  $N-1$ 인 스튜던트 t 분포를 이룬다.

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{N}}} \sim t(x; 0, 1, N-1)$$

이 식에서  $\bar{x}$ ,  $s$ 는 각각 표본평균, 표본표준차다.

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

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

이 식은 주위 정규분포의 가설검에 관한 각종 상황에서 사용된다.

→  $N$ 이 커질수록  $t$ 의 분포는 정규분포에 가까워진다.

### 카이제곱분포

정규분포를 따른 확률변수  $X$ 의  $N$ 개의 표본  $x_1, \dots, x_N$ 의 합 (또는 제곱)은  
표본분산으로 정규화하면 스튜던트 t 분포를 따를 것을 배웠다.

그러나 이  $N$ 개의 표본들을 단순히 더하는 것이 아니라 제곱하여 더하면 야, 수렴값을 갖는 분포가 된다. 이 분포를 카이제곱 (chi-squared) 분포라고 하며,  
 $\chi^2(x; \nu)$ 으로 표기한다. 카이제곱분포는 스튜던트 t 분포처럼 자유도  $\nu$ 를 갖는다.

$$x_i \sim \mathcal{N}(x)$$

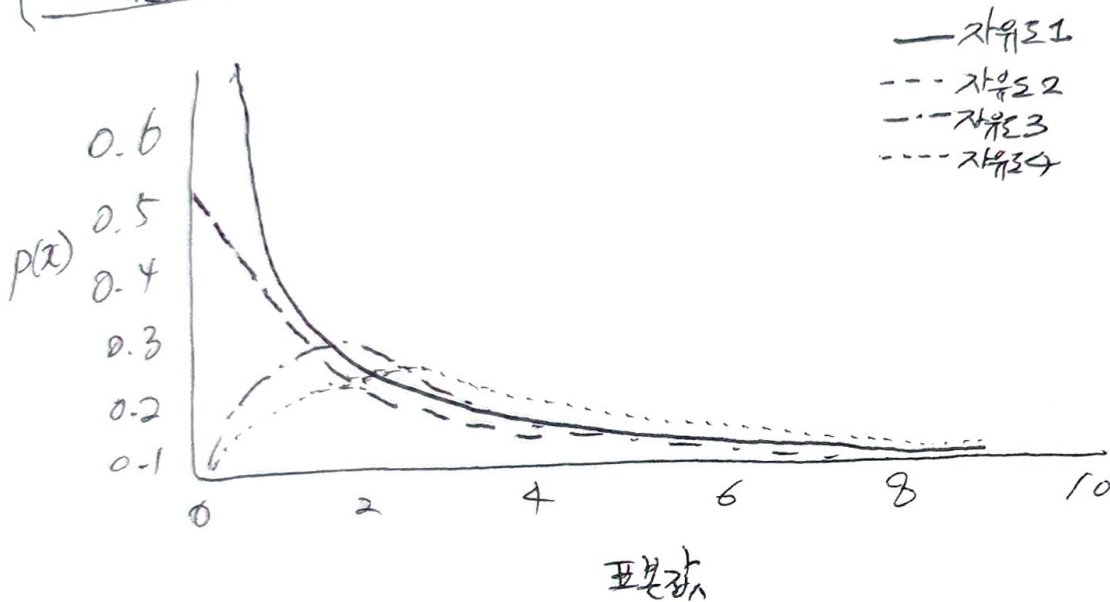
↓

$$\sum_{i=1}^N x_i^2 \sim \chi^2(x; \nu)$$

카이제곱분포의 확률 밀도함수는 다음과 같다.

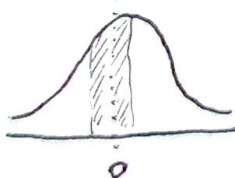
$$\chi^2(x; \nu) = \frac{x^{(\nu/2-1)} e^{-x/2}}{2^{\nu/2} \Gamma(\frac{\nu}{2})}$$

Scipy.stats.chi2()



\* 자유도의 의미 for chi-square

$X \sim N$

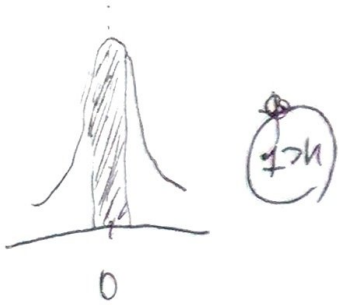
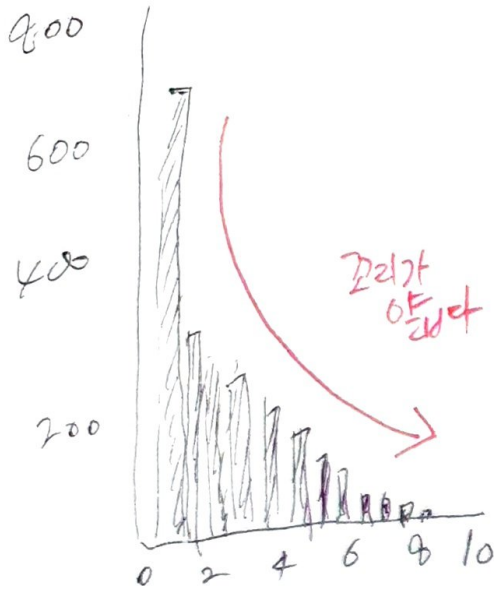


$\chi_1^2 \sim \chi^2(N=1)$

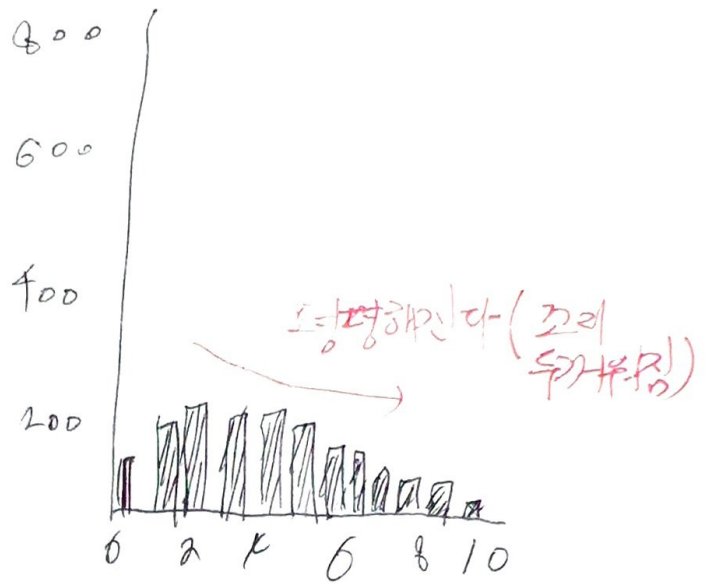
$\chi_1^2 + \chi_2^2 \sim \chi^2(N=2)$

$\chi_1^2 + \chi_2^2 + \chi_3^2 \sim \chi^2(N=3)$

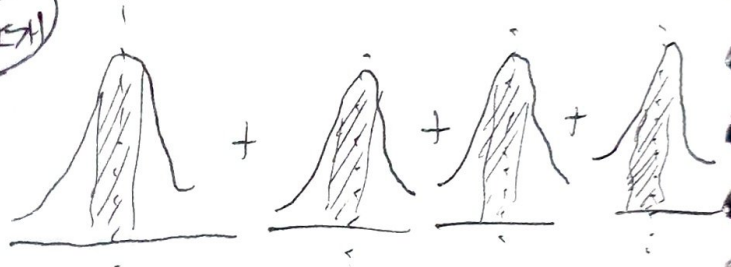
제공량의 분포 ( $N=1$ )



제공량의 분포 ( $N=4$ )



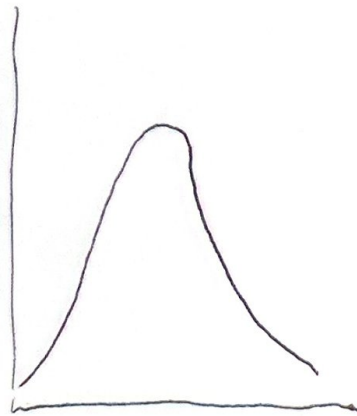
(4개)



$N=6$

(정규분포의  
정점 가져와야)

$N=30$



# F 분포

(스튜던트 t, 카이제곱분포) — 정규분포를 따르는 확률변수  $X$ 로부터 나온  $N$ 개의 표본에서 만들 수 있음

• 카이제곱분포를 따르는 독립적인 두 개의 확률변수  $x_1^2(x; N_1)$ 과  $x_2^2(x; N_2)$ 의 확률변수 표본을 각각  $x_1, x_2$ 라고 할 때, 이들 각각  $N_1, N_2$ 로 나누어 비율을 구하면  $F(x; N_1, N_2)$  분포가 된다.  $N_1, N_2$ 는 F분포의 자유도 모수라고 한다.

$$x_1 \sim \chi^2(N_1), x_2 \sim \chi^2(N_2) \rightarrow \frac{\frac{x_1}{N_1}}{\frac{x_2}{N_2}} \sim F(x; N_1, N_2)$$

F 분포의 확률밀도함수를 다음과 같다.

$$f(x; N_1, N_2) = \frac{\sqrt{\frac{(N_1 x)^{N_1} N_2^{N_2}}{(N_1 x + N_2)^{N_1 + N_2}}}}{x \beta\left(\frac{N_1}{2}, \frac{N_2}{2}\right)}$$

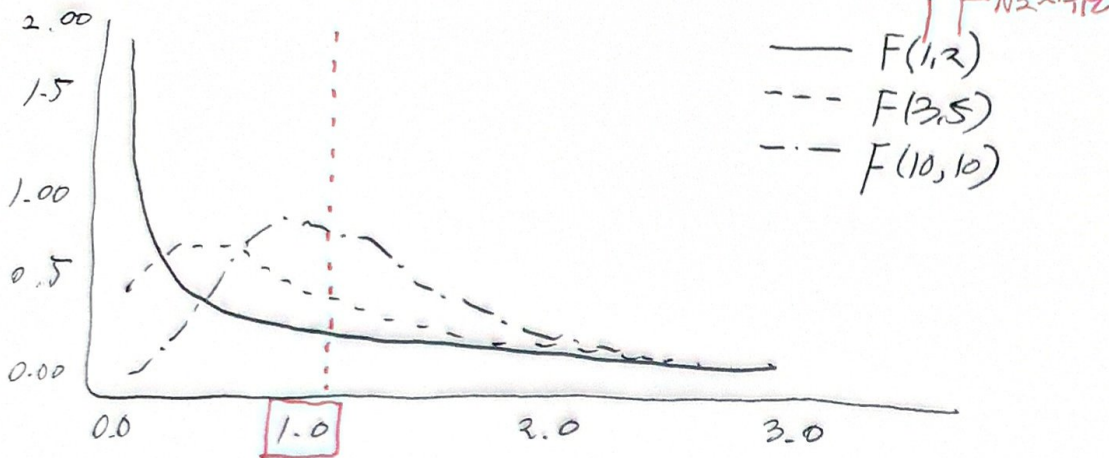
$\beta(x) = \text{베타}, \text{특수함수}$

\* 스튜던트 t 분포의 표본값을 제공한 값은 F분포를 따른다.

$$t(N)^2 = F(1, N)$$



# < 자유도이 다른 F분포의 모양 >



동일한 분포 (카이제곱)

이서 2개의 확률변수 분포를 가져와 나눈다면,  
직관적으로 생각했을 때 1이 많이 나오지 않을까?

→ 자유도가 높아질수록 최빈값이 1에 가까워짐을 확!

## 활용

스튜던트 t 분포, 카이제곱 분포, F 분포는 모두 정규 분포의 통계량 분포 (statistics distribution)의 일종이다. 선형 회귀 분석에서 이 통계량 분포들은 각각 다음 값에 대한 확률 분포로 사용된다.

- 스튜던트 t 분포 : 추정된 가중치에 대한 확률 분포
- 카이제곱 분포 : 2차 제곱합에 대한 확률 분포
- F 분포 : 비교 대상이 되는 선형 모형의 2차 제곱합에 대한 비율의 확률 분포

