

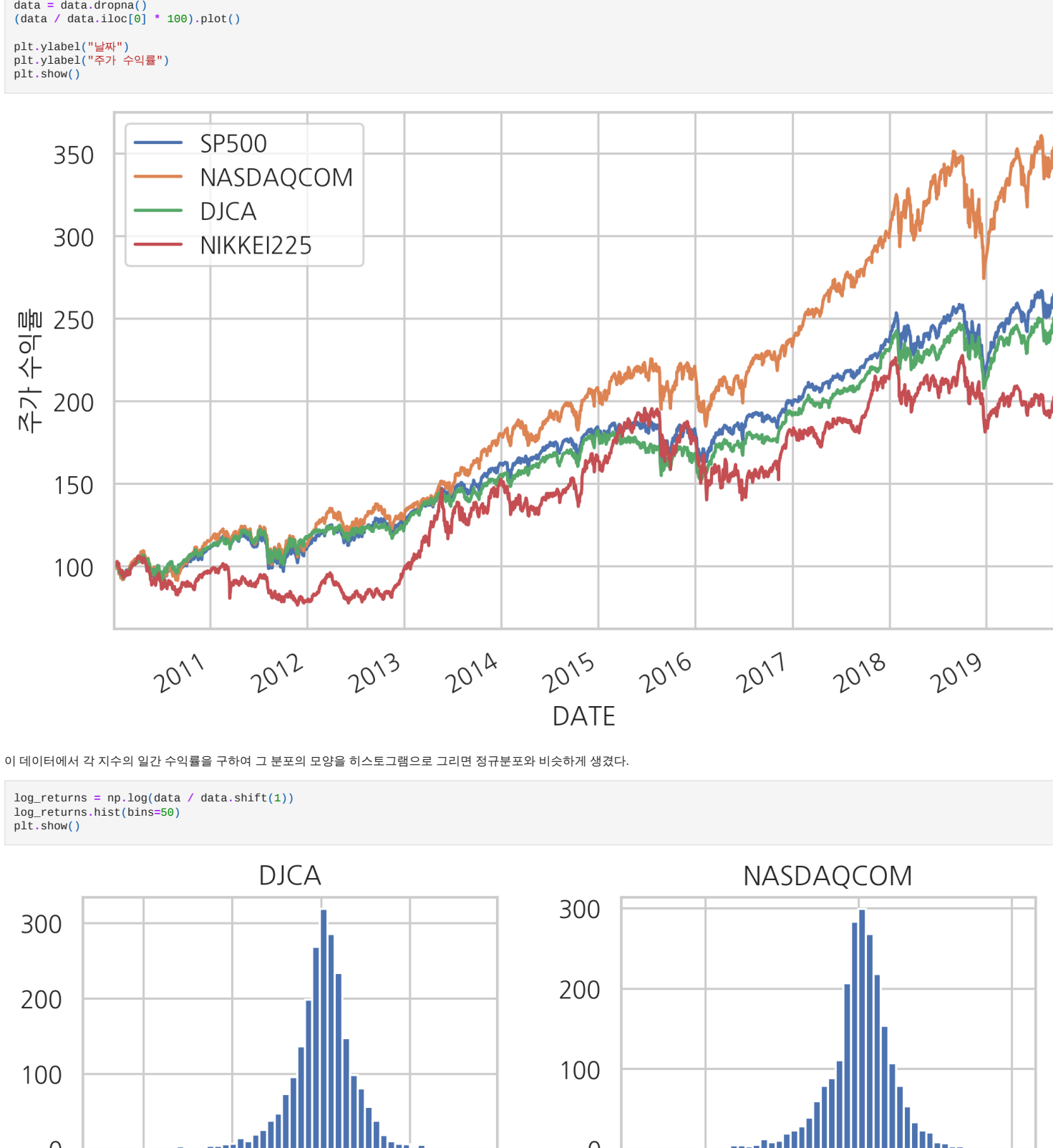
8.5 스튜던트 t분포, 카이제곱분포, F분포

이 절에서는 정규분포에서 파생된 분포를 공부한다. 정규분포에서 생성된 표본 데이터 집합에 여러 수식을 적용하여 값을 변환시키면 데이터 집합의 분포 모양이 달라지는데 적용된 수식에 따라 스튜던트 t분포, 카이제곱분포, F분포가 만들어진다. 이 분포들은 통계량 분포라고도 부르는데 나중에 공부할 것일 것일 것이다.

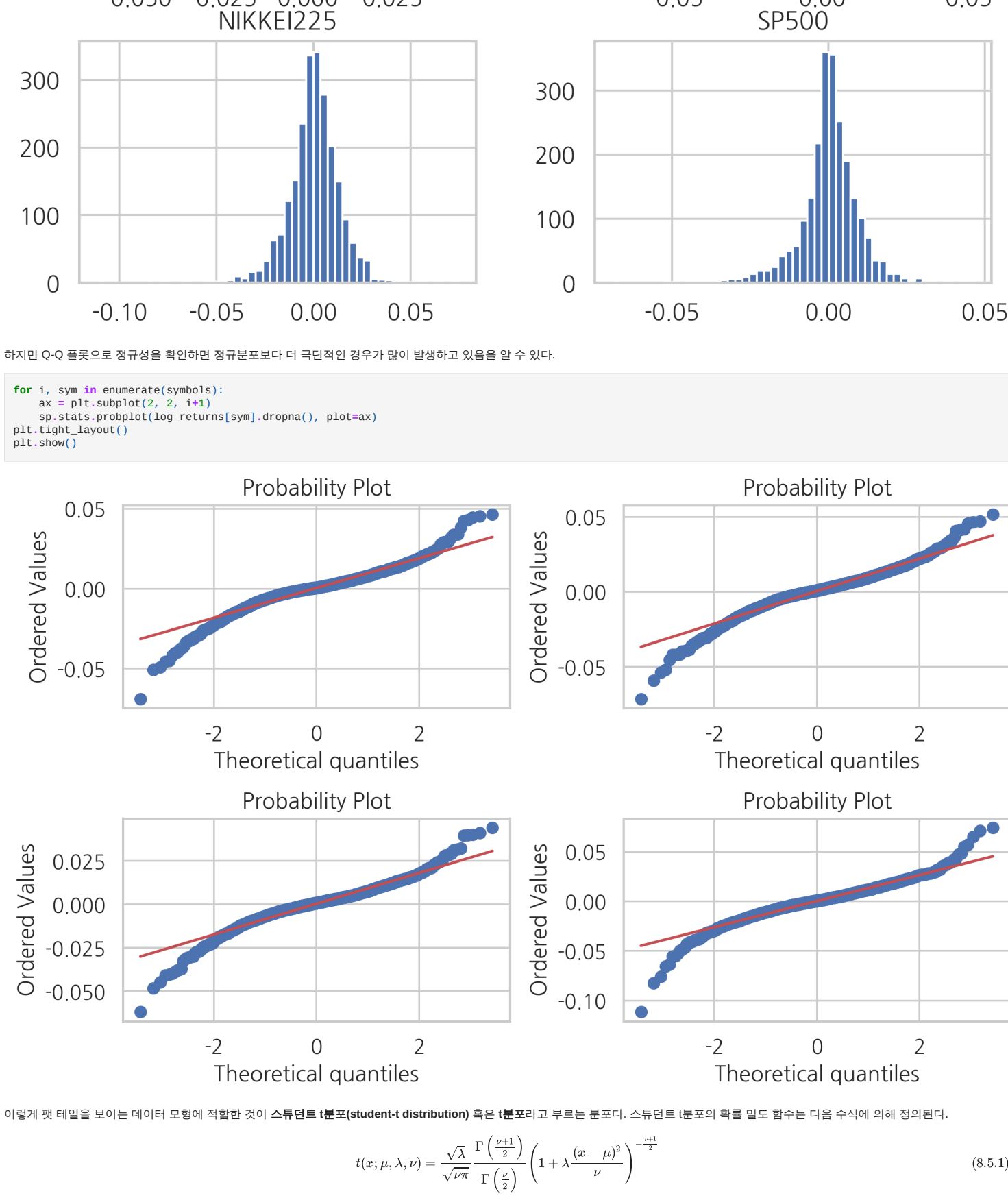
스튜던트 t분포

현실의 데이터를 실험하면 정규분포와 상관이 유사하지만 양 끝단의 비중이 정규분포에 비해 데이터 끝단을 더욱 극단적으로 밀어내어 꼬리를 길게 만든다. 정규분포와 가정했을 때 극단적 현상이 더 자주 발생한다는 뜻이다. 분포의 모양을 매우 양 끝단(고리) 부분이 정규분포보다 두 배 더 넓어지고 해서 이를 커테일(Curt) 현상이라고 한다. 예를 들어 주식의 수익률은 보통 정규분포를 따르는 것으로 가정하는데 실제로는 정규분포에 어는 자주 벗어날 수 있는 극단적인 사건들이 종종 발생하곤 한다. 금융시장에서는 이러한 현상을 블랙 스완(Black swan)이라고도 한다.

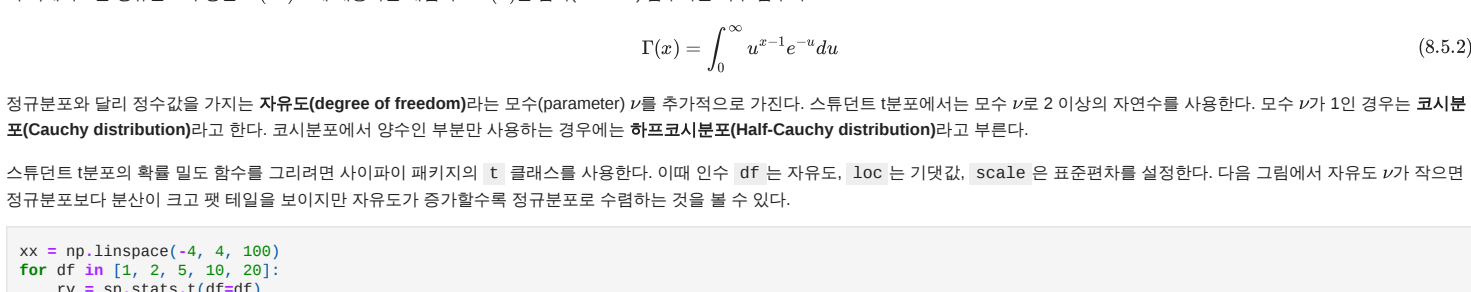
실제 과거의 주가 데이터를 확인해보자. 다음은 S&P 500, 나스닥(Nasdaq), 다우존스(Dow-Jones), 니케이255(Nikkei255) 네 가지의 주가지수 데이터이다. 비교를 위해 2010년의 값을 100으로 통일하였다.



이 데이터에서 각 지수의 일간 수익률을 구하여 그 분포의 모양을 히스토그램으로 그리고 정규분포와 비슷하게 생겼다.



하지만 Q-Q 플롯으로 정규성을 확인하면 정규분포보다 더 극단적인 경우가 많이 발생하고 있음을 알 수 있다.



이렇게 몇몇 데이터를 보면 데이터 모양에 적합한 것이 스튜던트 t분포(student-t distribution) 혹은 t분포라고 부른다. 스튜던트 t분포의 확률 밀도 함수는 다음 수식에 의해 정의된다.

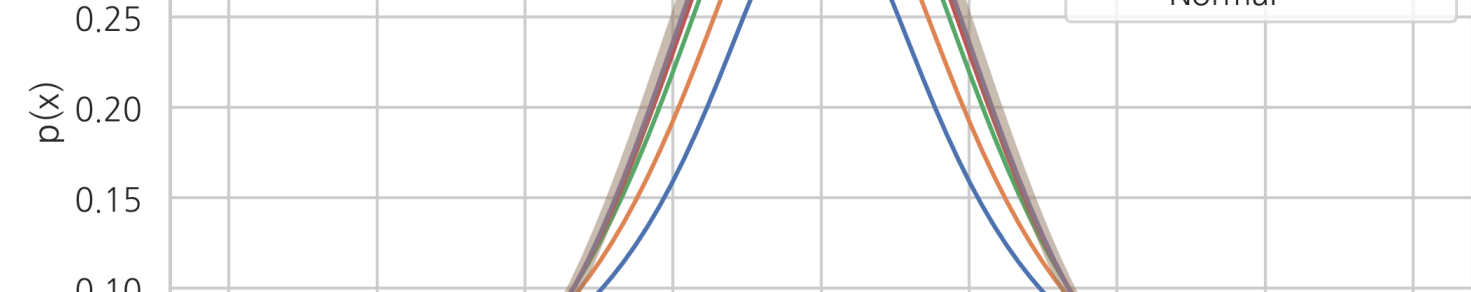
$$f(x; \mu, \lambda, \nu) = \frac{\sqrt{\lambda}}{\sqrt{\nu\pi}} \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})} \left(1 + \lambda \frac{(x-\mu)^2}{\nu} \right)^{-\frac{\nu+1}{2}} \quad (8.5.1)$$

이 식에서 λ 는 정규분포의 정밀도 (σ^2)⁻¹에 대응하는 개념이고 $\Gamma(x)$ 는 감마(Gamma) 함수라는 특수 함수이다.

$$\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du \quad (8.5.2)$$

정규분포와 달리 정수값을 가지는 자유도(degree of freedom)라는 모수(parameter) ν 를 추가적으로 가진다. 스튜던트 t분포에서는 모수 ν 로 2 이상의 자연수를 사용한다. 모수 ν 가 1인 경우는 코시분포(Cauchy distribution)라고 한다. 코시분포에서 양수인 부분만 사용하는 경우에는 하프코시분포(Half-Cauchy distribution)라고 부른다.

스튜던트 t분포의 확률 밀도 함수를 그리려면 사이파이 패키지의 t 클래스를 사용한다. 이때 인수 df는 자유도, loc는 기댓값, scale은 표준편차를 설정한다. 다음 그림에서 자유도 ν 가 작으면 정규분포보다 분산이 크고 몇 테일을 보이지만 자유도가 증가할수록 정규분포로 수렴하는 것을 볼 수 있다.



스튜던트 t분포의 기댓값과 분산은 다음과 같다.

- 기댓값

$$E[X] = \mu \quad (8.5.3)$$

- 분산

$$\text{Var}[X] = \frac{\nu}{\lambda(\nu-2)} \quad (8.5.4)$$

분산의 대한 식은 $\nu > 2$ 인 경우만 적용된다. $\nu = 1, 2$ 일 때는 분산이 무한대가 된다.

t 통계량

정규분포의 분포를 표준편차로 나눠 정규화된 z 통계량인 z 통계량인 정규분포가 된다는 것은 이미 공부하였다. 그런데 z 통계량을 구하려면 확률분포의 정확한 표준편차를 우리가 알고 있어야 한다. 하지만 현실적으로는 표준편차를 정확히 알 수 없기 때문에 표본에서 측정된 표본표준편차(sample standard deviation)로 정규화를 수행해 있다. 정규분포로부터 얻은 N 개의 표본 x_1, \dots, x_N 에 계산한 표본평균을 표본표준편차로 정규화한 값을 t 통계량이라고 한다.

t 통계량은 자유도가 $N - 1$ 인 스튜던트 t분포를 이룬다.

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{N}}} \sim t(x; 0, 1, N - 1) \quad (8.5.5)$$

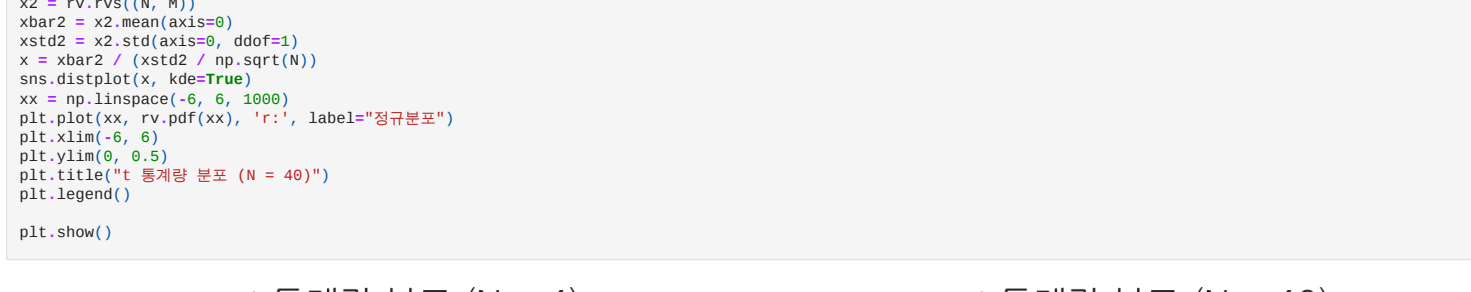
이 식에서 \bar{x} , s 는 각각 표본평균, 표본표준편차이다.

$$\bar{x} = \frac{x_1 + \dots + x_N}{N} \quad (8.5.6)$$

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (8.5.7)$$

이 정리는 주위 정규분포의 기댓값에 관한 각종 검증에서 사용된다.

다음은 시뮬레이션을 사용하여 표본표준편차로 정규화된 표본평균과 정규분포를 비교한 것이다. 왼쪽은 $N = 4$, 오른쪽은 $N = 40$ 인 경우이다.



카이제곱분포

정규분포를 따르는 표본 변수 X_1, \dots, X_N 의 합의 제곱(또는 평균)은 표본 분산으로 정규화하면 스튜던트 t분포를 따른다는 것을 배웠다.

그런데 이 제곱의 표준편차를 단순히 더하는 것이 아니라 제곱을 하여 더하면 양수만을 가지는 분포가 된다. 이 분포를 카이제곱(chi-squared)분포라고 하며 $\chi^2(x; \nu)$ 으로 표기한다. 카이제곱분포는 스튜던트 t분포와 마찬가지로 자유도 모수를 가진다.

$$x_i \sim \mathcal{N}(x) \quad (8.5.8)$$

$$\downarrow$$

$$\sum_{i=1}^N x_i^2 \sim \chi^2(x; \nu = N) \quad (8.5.10)$$

카이제곱분포의 확률 밀도 함수는 다음과 같다.

$$\chi^2(x; \nu) = \frac{1}{2^{\nu/2} \Gamma(\frac{\nu}{2})} x^{\nu/2-1} e^{-x/2} \quad (8.5.11)$$

사이파이 stats 서브패키지의 chi2 클래스를 사용하여 확률 밀도 함수의 모양을 살펴보면 다음과 같다.

