

8.4 정규분포와 중심극한정리

정규분포 (normal distribution) 또는 가우스 정규분포 (Gaussian normal distribution)라는 분포는 자연 현상에서 나타나는 숫자를 확률 모형으로 모형화할 때 많이 사용한다.

* 정규분포의 확률밀도함수 (pdf)

$$N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

분산의 역수를 정밀도 (precision) 라고 부기도 한다.

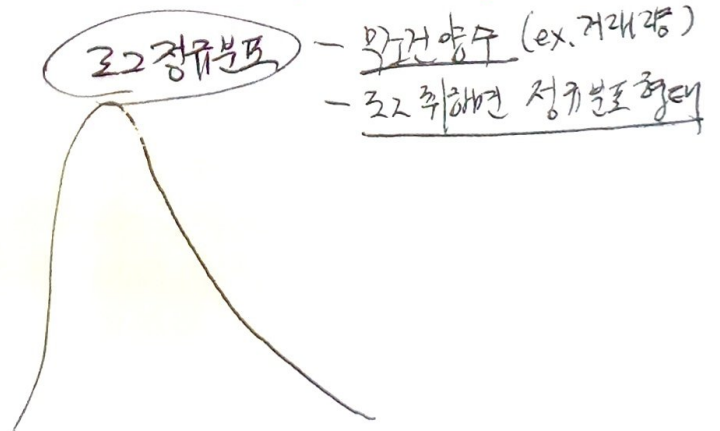
$$\beta = \frac{1}{\sigma^2}$$

→ 역수를 취하는 이유 = 편이 편함

정규분포 중에서도 평균이 0이고 분산이 1인 ($\mu=0, \sigma^2=1$) 정규분포를 표준정규분포 (standard normal distribution) 라고 한다.

정규분포의 확률밀도함수는 다음과 같은 성질을 가진다.

- $x = \mu$ 일 때 확률밀도가 최대가 된다.
- $x = \infty$ 또는 $x = -\infty$ 로 갈수록 확률밀도가 줄어든다.



Q-Q 플롯

* 정규분포는 여러 연속확률분포 중에서도 가장 널리 사용되는 확률분포.

→ 어떤 확률변수의 분포가 **정규분포인지 아닌지 확인**하는 것은 중요한 통계적 분석 중 하나

Q-Q (Quantile-Quantile) 플롯

- 분석할 표본 데이터의 분포와 정규분포의 분포 형태 비교
- 표본 데이터가 정규분포를 따르는지 검사하는 간단한 시각적 도구
- 동일 분위수에 해당하는 **정상 분포의 값과 주어진 데이터값을 pairing 하여** 그린 scatter plot

Q-Q process (대략적 방법론)

① 표본 데이터 정렬

② 하나하나의 표본 데이터가 전체 데이터 중의 몇 %에 해당하는지 위치값 계산
 • 위치값이란 특정 순위 (order) 값이 나타낼 가능성이 높은 값을 뜻하는
 순서통계량 (order statistics) 이라는 값 사용 (**몇 번째 데이터일 가능성이 가장 높은 값**)

③ 각 표본 데이터의 위치값이 정규분포의 누적확률함수 (cdf) 값이 되는
 표준 정규분포의 표준값 계산

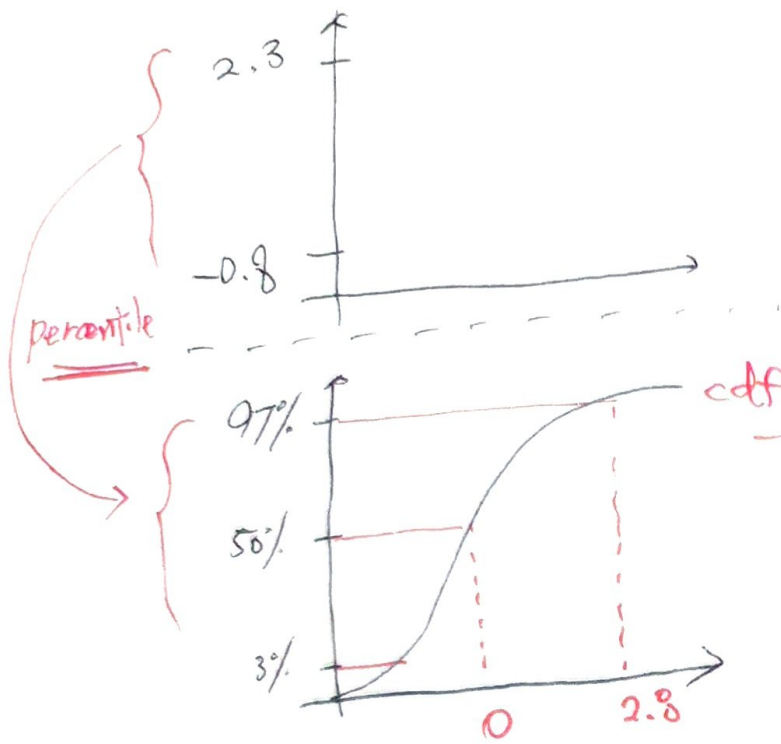
즉, 위치값에 대한 누적확률함수의 역함수 값 구하기

= 표준 정규분포의 분위함수 (quantile function)

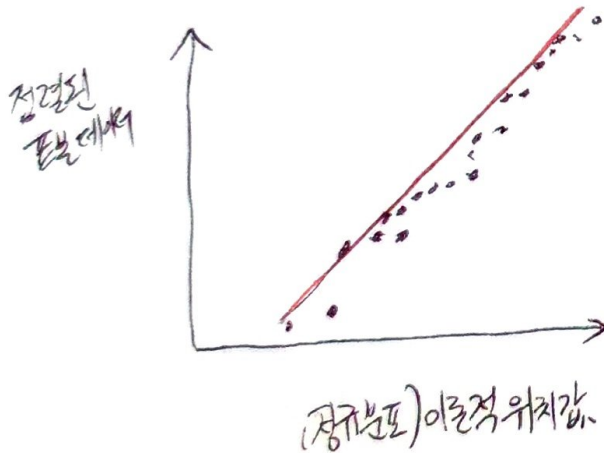
④ 표준 정규분포의 1% 분위함수값은 $F^{-1}(0.01)$, 약 -2.326

⑤ 정렬된 표본 데이터 (ordered values) 와 그에 대응하는 분위수 (theoretical quantiles) 를 하나의 쌍으로 맞추어 2차원 공간에 하나의 점 (point) 으로 그린다.

⑥ 모든 표본에 대해 ② - ⑤ 반복하여 스캐터플롯 완성

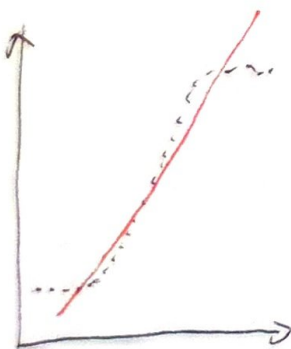


- 정규분포라고
가정하고
cdf matching

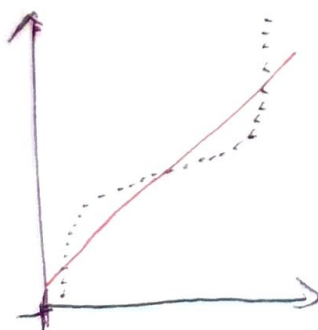
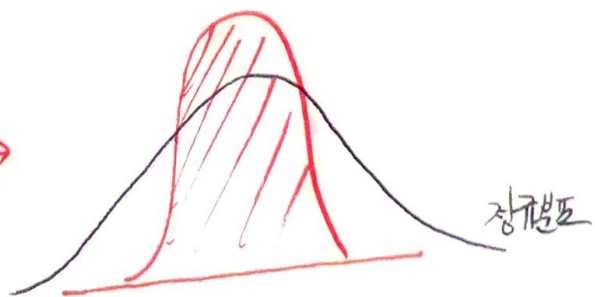


일직선이 형성되면
정규분포에 가깝다는 의미.

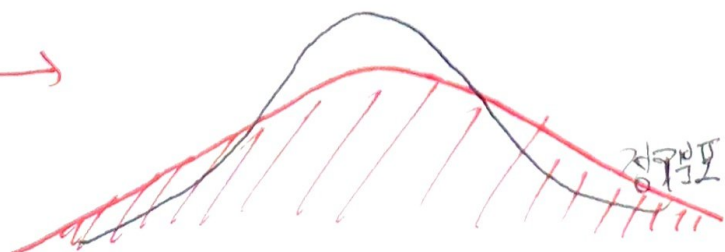
qq plot



정규분포 아님



정규분포 아님



중심극한정리 Central Limit Theorem

실세계에서 발생하는 현상 중 많은 것들을 정규분포로 모델링가능.

→ 그 이유 중 하나: 중심극한정리 (Central Limit Theorem)

→ "어려 확률변수" 합이 정규분포와 비슷한 분포를 이루는 현상"

분포 X_1, X_2, \dots, X_N 이 기대값이 μ 이고 분산이 σ^2 로 동일한 분포 (가치값과 분산의 값이 동일한 것이며, 분포의 모양은 달라도 됨)이며 서로 독립인 확률변수들이라 하자.

X_1, X_2, \dots, X_N 에서 뽑은 각각의 표본 데이터 x_1, x_2, \dots, x_N 의 표본평균

$$\bar{x}_N = \frac{1}{N} (x_1 + \dots + x_N)$$

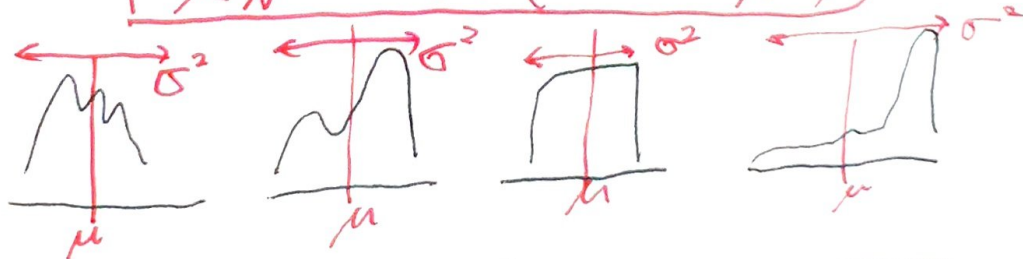
→ 마찬가지로 예측할 수 없는 확률변수다. 이 확률변수를 \bar{x}_N 이라곤 하자.
중심극한정리는 다음과 같다.

N 개의 양의 분포로부터 얻은 표본의 평균은
 N 이 증가할수록 기대값이 μ , 분산이 $\frac{\sigma^2}{N}$ 인 정규분포로 수렴한다.

" N 이 커질수록 수렴한다"

$$\bar{X}_N \xrightarrow{d} N\left(\mu; \mu, \frac{\sigma^2}{N}\right)$$

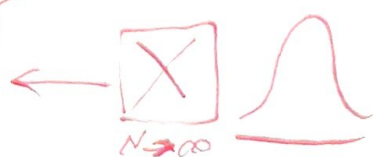
pdf



(N 개의 양)
(μ, σ^2)

X_1 X_2 X_8 X_9

$$\bar{x}_N = \frac{1}{N} (x_1 + x_2 + x_3 + x_4)$$



평균이 0, 분산이 1이 되도록 다음처럼 정규화를 하면
다음과 같이 쓸 수 있다.

{ N 개의 임의의 분포로부터 얻은 표본의 평균을 정규화하면
 N 이 증가할수록 표준정규분포로 수렴한다. }

$$\frac{\bar{X}_N - \mu}{\frac{\sigma}{\sqrt{N}}} \xrightarrow{d} \mathcal{N}(x; 0, 1)$$

* 중심극한정리가 중요한 이유?

현실이 존재하는 매우 복잡한 데이터들은 각각은 서로 매우 다르지만
많이 합쳐수록 정규분포의 특징을 나타낸다

→ 잘 모르겠으면 일단 많이 모아서 정규분포로 분석해보자 *

정규분포의 통계량 분포

임의의 분포가 아닌 N 개의 정규분포로부터 얻은 표본 데이터를 구한 표본평균은
어떤 분포를 가지게 될까?

N 개의 정규분포로부터 얻은 표본의 합은 N 과 상관없이 $N\mu$,
분산이 $N\sigma^2$ 인 정규분포다.

$$x_i \sim \mathcal{N}(\mu, \sigma^2) \rightarrow \sum_{i=1}^N x_i \sim \mathcal{N}(N\mu, N\sigma^2)$$

정규분포의 분포에 상수를 빼거나 곱해도 정규분포다. 이 경우에도 위와 같이
기댓값이 0, 표준편차가 1이 되도록 정규화를 하면 다음과 같이 쓸 수 있다.

$$x_i \sim \mathcal{N}(\mu, \sigma^2) \rightarrow \boxed{z} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{N}}} \sim \mathcal{N}(0, 1)$$

정규분포 표본의 변환을 정규화된 통계량 = Z 통계량

* 중심극한정리 — 무작위 표본, 무한대이어야 정규분포가 된다

* Z 통계량 — 정규분포 표본, N 개수와 상관없이
항상 정확하게 표준정규분포.

선형회귀모형과 정규분포

* 정규분포는 선형회귀모형에서 잡음 (disturbance) 을 모형화하는 데 사용

{ 선형회귀모형은 입력변수 x_1, \dots, x_n 이 종속변수 y 이 선형적으로
영향을 미치는 모형이다.

$$\hat{y} = w_1 x_1 + \dots + w_n x_n \approx y$$

이 모형은 다음과 같이 표현될 수 있다.

$$y = w_1 x_1 + \dots + w_n x_n + \epsilon$$

ϵ 는 잡음 (disturbance) 이라고 하며, 우리가 값을 측정할 수 있는 양을 뜻한다.
예측값과 실제값의 차이를 뜻하는 잔차 (residual) 와는 다르다.

$\epsilon = \text{disturbance}$
≠
residual

잡음은 선형리모델을 만들 때 하나하나의 영향력이 작거나 크거나
특정재|항들에서 무시할 수없는 변수들의 영향을 하나로 합친 것이다.

즉, 원래 입력은 x_1, \dots, x_N, \dots 의 거의 무한한 갯수의 입력변수의
영향을 받는다.

$$y = w_1 x_1 + \dots + w_N x_N + w_{N+1} x_{N+1} + w_{N+2} x_{N+2} + \dots$$

하지만 이 중에서 입력변수 x_1, \dots, x_N 만이 영향력이 크거나 측정이 쉽다면
다른 변수의 영향을 하나의 확률변수로 합쳐서 표현할 수 있다.

↓

$$\epsilon = w_{N+1} x_{N+1} + w_{N+2} x_{N+2} + \dots$$

(측정X, 구할수없는 어떤 값.)

중심정리(정리)에 의해,

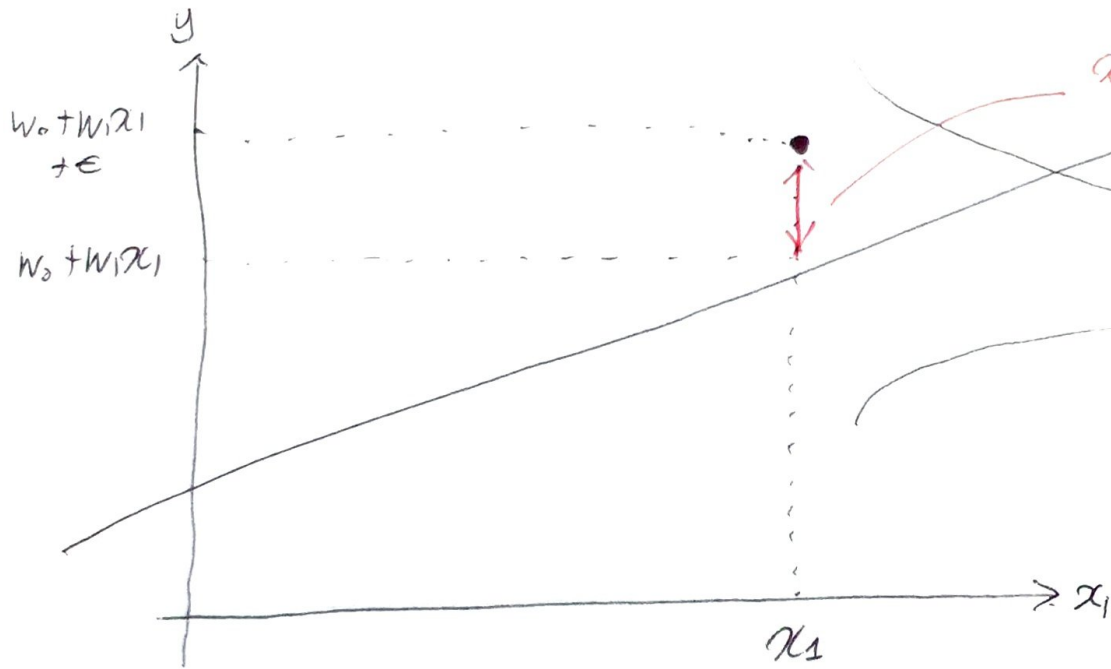
입력의 확률변수의 합은 정규분포와 비슷한 형태가 된다. 또한 ϵ 의 기대값이
0이 아니라면 다음처럼 상수인 $w_0 = E[\epsilon]$ 을 추가하는 대신 ϵ 의 기대값이
0이 되도록 할 수 있으므로,

$$y = w_0 + w_1 x_1 + \dots + w_N x_N + \epsilon$$

잡음 ϵ 이 기대값인 0인 정규분포라고 가정하는 것은 합리적이다.

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y = w_0 + w_1 x_1 + \overbrace{w_2 x_2 + w_3 x_3 + \dots}^{\epsilon}$$



x_1 이외의 실제 값을
설명하는 다른
영향들이 모여서

우리가 분석하고
못한 또 다른
설명변수들이
누적되어 미친
영향이크다!!
||
 ϵ
(epsilon)

<선형회귀모델의 오차 범위>

- ① ϵ 를 배제
- ② ϵ 가 정규분포에서 비롯된다고 가정하고,
 ϵ 의 범위를 이용