

011 - Normal Curves

EPIB 607

Sahir Rai Bhatnagar
Department of Epidemiology, Biostatistics, and Occupational Health
McGill University

`sahir.bhatnagar@mcgill.ca`

slides compiled on September 26, 2021



The Normal (Gaussian) distribution

What is it?

- A distribution that describes continuous (numerical) data
- Can also be used to approximate discrete data distributions
- Range is (technically) infinite, though the probability of seeing very large or very small values is extremely tiny
- Fully described by only two parameters, the mean and variance (μ and σ^2)
- **NOTE:** R use the short-hand: $X \sim \mathcal{N}(\mu, \sigma)$, denoting the normal distribution as a function of the mean and *standard deviation*. This is not standard; many texts instead write $X \sim \mathcal{N}(\mu, \sigma^2)$. Be careful of this!

The Normal (Gaussian) distribution

Carl Gauss was a German mathematician who developed a number of important advances in statistics such as the method of least squares.



Figure: The Deutsche Bundesbank issued Deutsche Mark banknotes in 15 different denominations, including this 10 Deutsche Marks banknote featuring Carl Friedrich Gauss.

The Normal distribution

Where do Normal data come from?

- Natural processes
 - ▶ Blood pressure
 - ▶ Height
 - ▶ Weight
- “Man-made” (or derived)
 - ▶ Binomial (proportion) and Poisson (count) data are approximately Normal under certain conditions
 - ▶ Sums and means of random variables (Central Limit Theorem)
 - ▶ Data can sometimes be made to look Normal via transformations (squares, logs, etc)

The Normal distribution

For Normal data, we can use the ~~Gaussian tables~~ **R** to answer the questions:

- What is the probability that a single observation X is
 - ▶ greater than X^* ?
 - ▶ less than X^* ?
 - ▶ between X_L^* and X_U^* ?
- That is, we can find out information about the percent distribution of X as a function of thresholds X^* , or X_L^* and X_U^* .
- We can also use the ~~Normal tables~~ **R** to find out information about thresholds X^* that will contain particular percentages of the data. I.e., we can find what threshold values will
 - ▶ Exclude the lower ω^* % of a population
 - ▶ Exclude the upper ω^* % of a population
 - ▶ Contain the middle ω^* % of a population

The Normal distribution

We can use ~~the Gaussian tables~~ R to answer these questions **no matter what the values of** μ and σ^2 .

That is, the % of the Normal distribution falling between $X_L^* = \mu - m_1\sigma$ and $X_U^* = \mu + m_2\sigma$ where m_1, m_2 are any multiples **remains the same** for any μ and σ .

How so??

Because we can **standardize** any $X \sim \mathcal{N}(\mu, \sigma)$ to find $Z \sim \mathcal{N}(0, 1)$

The Normal distribution

An illustration using IQ scores, which we presume have a $\mathcal{N}(100, 13)$ distribution of scores.

Q1: What percentage of scores are **above** 130?

Two steps:

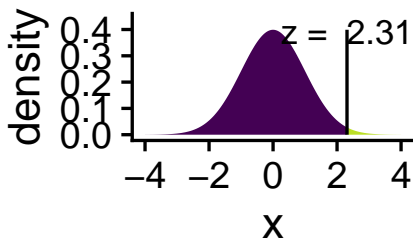
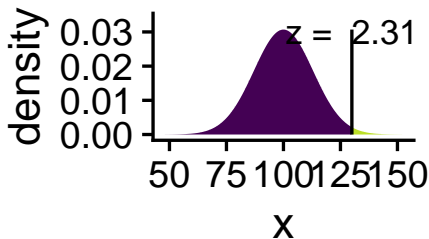
1. Change of location from $\mu_X = 100$ to $\mu_Z = 0$
2. Change of scale from $\sigma_X = 13$ to $\sigma_Z = 1$

Together, this gives us

$$Z = \frac{X - \mu_X}{\sigma_X} = \frac{130 - 100}{13} = 2.31$$

The Normal distribution

The position of $X=130$ in a $\mathcal{N}(100, 13)$ distribution is the same as the place of $Z = 2.31$ on the $\mathcal{N}(0, 1)$, which we call the **standardized** Normal distribution (or Z-distribution).



The Normal distribution

How are the values in the Normal tables found?

Normal density:

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp \frac{-(x - \mu)^2}{2\sigma^2}$$

Probabilities found by integration (area under the Normal curve):

$$P(a \leq x \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}\sigma} \exp \frac{-(x - \mu)^2}{2\sigma^2} dx$$

The Normal distribution

(The percent above $X = 130$) = (% above $Z = 2.31$) = 1.04%

How do we know this? We look at the lower tail probability of 2.31 [i.e., the % below 2.31], and then subtract it from 1:

1. $P(X < 130) = P(Z < 2.31) = 0.9896$
2. $P(X > 130) = 1 - P(X < 130) = 0.0104$

So 130 is the 98.96th percentile of a $\mathcal{N}(100,13)$ distribution.

Reminder about percentiles and quantiles

- **Quantile**

- ▶ Any set of data, arranged in ascending or descending order, can be divided into various parts, also known as partitions or subsets, regulated by quantiles.
- ▶ Quantile is a generic term for those values that divide the set into partitions of size n , so that each part represents $1/n$ of the set.
- ▶ Quantiles are not the partition itself. They are the numbers that define the partition.
- ▶ You can think of them as a sort of numeric boundary.

- **Percentile**

- ▶ Percentiles are quite similar to quantiles: they split your set, but only into two partitions.
- ▶ For a generic k th percentile, the lower partition contains $k\%$ of the data, and the upper partition contains the rest of the data, which amounts to $100 - k \%$, because the total amount of data is 100% .
- ▶ Of course k can be any number between 0 and 100.

More about percentiles and quantiles

- In class, we will find ourselves asking for the quantiles of a distribution.
- Percentiles go from 0 to 100
- Quantiles go from any number to any number
- Percentiles are examples of quantiles and you might find some people use them interchangeably (though this may not always be correct since quantiles can take on any value, positive or negative).
- **In particular**, R uses the term quantiles.
- **In the previous example**, we saw that $P(Z < 2.31) = 0.9896$. In R, 2.31 is called the quantile .

The Normal distribution

(The percent above $X = 130$) = (% above $Z = 2.31$) = 1.04%

But wait!! The standard Normal is symmetric about 0, so we can do this another way... The % **above** 2.31 is equal to the % **below** -2.31:

$$\begin{aligned}P(X > 130) &= P(Z > 2.31) \\&\Rightarrow P(Z > 2.31) = P(Z < -2.31) \\&\Rightarrow P(X > 130) = P(Z < -2.31) = 0.0104\end{aligned}$$

So 130 is the 98.96th percentile of a $\mathcal{N}(130, 13)$ distribution. What is the 1.04th percentile?

Transform from $Z = -2.31$ back to X :

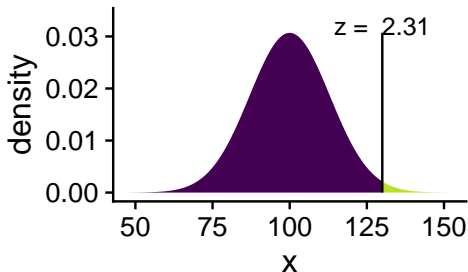
$$X = \sigma Z + \mu = 13(-2.31) + 100 = 69.97.$$

For probabilities we use *pnorm*

```
stats::pnorm(q = 130, mean = 100, sd = 13)
```

```
## [1] 0.9894919
```

```
mosaic::xpnorm(q = 130, mean = 100, sd = 13)
```



```
## [1] 0.9894919
```

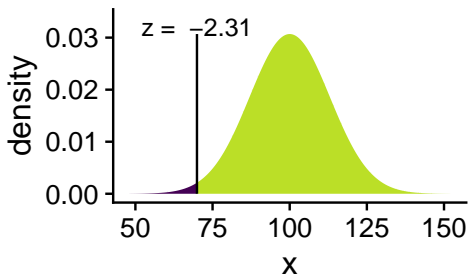
- `pnorm` returns the integral from $-\infty$ to q for a $\mathcal{N}(\mu, \sigma)$
- `pnorm` goes from *quantiles* (think *Z* scores) to probabilities

For quantiles we use *qnorm*

```
stats::qnorm(p = 0.0104, mean = 100, sd = 13)
```

```
## [1] 69.94926
```

```
mosaic::xqnorm(p = 0.0104, mean = 100, sd = 13)
```



```
## [1] 69.94926
```

- `qnorm` answers the question: What is the Z-score of the p th percentile of the normal distribution?
- `qnorm` goes from *probabilities* to quantiles

The Normal distribution

Q2: What is the probability of seeing an IQ score **as extreme as** (think highly unusual) 130?

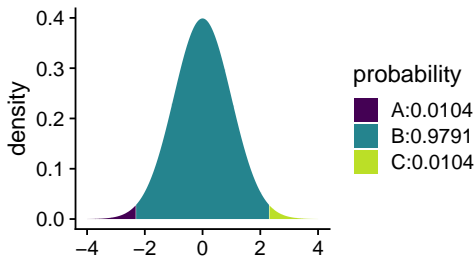
1. Again, we find that $X = 130$ is the same percentile of the IQ Normal distribution as $Z = 2.31$ is of the standard Normal.
2. To see what scores are as extreme, we want to know the probability that $Z > 2.31$ or that $Z < -2.31$.
3. As we saw previously, $P(Z > 2.31) = P(Z < -2.31) = 0.0104$, so the probability of seeing an IQ as extreme or more so than 130 is $2 \times 0.0104 = 0.0208$.

Finding tail probabilities

```
# lower.tail = TRUE is the default
stats::pnorm(q = -2.31, mean = 0, sd = 1, lower.tail = TRUE) +
stats::pnorm(q = 2.31, mean = 0, sd = 1, lower.tail = FALSE)

## [1] 0.02088815
```

```
mosaic::xpnorm(q = c(-2.31, 2.31), mean = 0, sd = 1)
```



```
## [1] 0.01044408 0.98955592
```

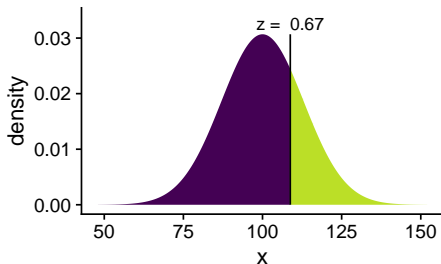
The Normal distribution

Q3: What is the 75th percentile of the IQ scores distribution?

We now have to reverse the sequence of steps:

- **Ask yourself:** What Z value corresponds to a probability of 0.75? Should you use `pnorm` or `qnorm`?

```
mosaic::xqnorm(p = 0.75, mean = 100, sd = 13)
```



```
## [1] 108.7684
```

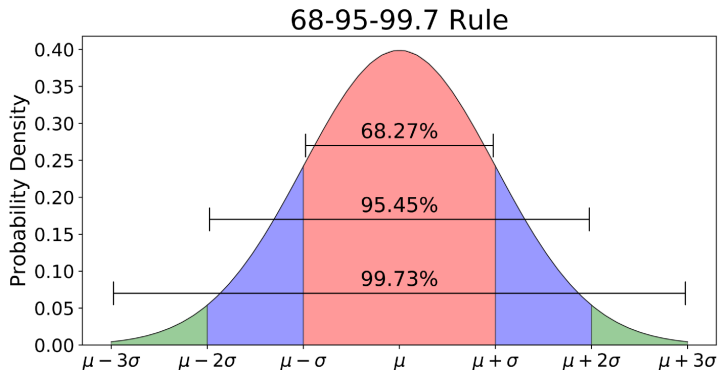
This tells us that 75% of the IQ scores fall below 108.8.

Empirical Rule or 68-95-99.7% Rule

In any normal distribution with mean μ and standard deviation σ :

- Approximately 68% of the data fall within one standard deviation of the mean.
- Approximately 95% of the data fall within two standard deviations of the mean.
- Approximately 99.7% of the data fall within three standard deviations of the mean.

Empirical Rule or 68-95-99.7% Rule



Properties of Normal random variables

Special properties of the Normal distribution:

- If Y is a Normal random variable, then so is $a + bY$.
- If X and Y are two Normal random variables, then $X + Y$ is a Normal random variable. What is the mean and variance of this new random variable?
- If X and Y are two Normal random variables and $\rho_{XY} = 0$ (correlation between X and Y), then X and Y are independent.

Properties of Normal random variables

Let $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma)$, and let each Y_i be independent of the others.
(think simple random sample)

Then $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ has what distribution?

- The sum of Normal random variables is Normal, so \bar{Y} is a Normal random variable.
- $E(\bar{Y}) = \frac{1}{n} \sum_{i=1}^n E(Y_i) = \frac{1}{n} \sum_{i=1}^n \mu = \mu.$
- $Var(\bar{Y}) = Var(\frac{1}{n} \sum_{i=1}^n Y_i) = \frac{1}{n^2} \sum_{i=1}^n Var(Y_i) = \sigma^2/n.$
- Standard Error of $\bar{Y} = \sqrt{Var(\bar{Y})} = \sigma/\sqrt{n}$

Session Info

```
R version 4.0.4 (2021-02-15)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Pop!_OS 21.04

Matrix products: default
BLAS:   /usr/lib/x86_64-linux-gnu/openblas-pthread/libblas.so.3
LAPACK: /usr/lib/x86_64-linux-gnu/openblas-pthread/libopenblas-p-r0.3.13.so

attached base packages:
[1] tools      stats      graphics  grDevices  utils      datasets  methods
[8] base

other attached packages:
[1] DT_0.16 mosaic_1.7.0 Matrix_1.3-2 mosaicData_0.20.1
[5] ggformula_0.9.4 ggstance_0.3.4 lattice_0.20-41 kableExtra_1.2.1
[9] socviz_1.2 gapminder_0.3.0 here_0.1 NCStats_0.4.7
[13] FSA_0.8.30 forcats_0.5.1 stringr_1.4.0 dplyr_1.0.7
[17] purrr_0.3.4 readr_1.4.0 tidyr_1.1.3 tibble_3.1.3
[21] ggplot2_3.3.5 tidyverse_1.3.0 knitr_1.33

loaded via a namespace (and not attached):
[1] fs_1.5.0 lubridate_1.7.9 webshot_0.5.2 httr_1.4.2
[5] rprojroot_2.0.2 backports_1.2.1 utf8_1.2.2 R6_2.5.1
[9] DBI_1.1.1 colorspace_2.0-2 withr_2.4.2 tidyselect_1.1.1
[13] gridExtra_2.3 leaflet_2.0.3 curl_4.3.2 compiler_4.0.4
[17] cli_3.0.1 rvest_1.0.0 pacman_0.5.1 xml2_1.3.2
[21] ggdendro_0.1.22 labeling_0.4.2 mosaicCore_0.8.0 scales_1.1.1
[25] digest_0.6.27 foreign_0.8-81 rmarkdown_2.9.7 rio_0.5.16
[29] pkgconfig_2.0.3 htmltools_0.5.2 highr_0.9 dbplyr_1.4.4
[33] fastmap_1.1.0 htmlwidgets_1.5.3 rlang_0.4.11 readxl_1.3.1
[37] rstudioapi_0.13 farver_2.1.0 generics_0.1.0 jsonlite_1.7.2
[41] crosstalk_1.1.1 zip_2.2.0 car_3.0-9 magrittr_2.0.1
[45] Rcpp_1.0.7 munsell_0.5.0 fansi_0.5.0 abind_1.4-5
[49] lifecycle_1.0.0 stringi_1.7.3 carData_3.0-4 MASS_7.3-53.1
[53] plyr_1.8.6 grid_4.0.4 blob_1.2.1 ggrepel_0.8.2
[57] crayon_1.4.1 cowplot_1.1.0 haven_2.3.1 splines_4.0.4
[61] hms_1.0.0 pillar_1.6.2 reprex_0.3.0 glue_1.4.2
[65] evaluate_0.14 data.table_1.14.0 modelr_0.1.8 vctrs_0.3.8
[69] tweenr_1.0.1 cellranger_1.1.0 gtable_0.3.0 polyclip_1.10-0
[73] assertthat_0.2.1 TeachingDemos_2.12 xfun_0.25 ggforce_0.3.2
[77] rnorm_1.4.1 broom_0.7.2 viridisLite_0.4.0 ellipsis_0.3.2
```