

Lecture - Day 3 (part 2)

Normal distribution

MF9130E V24

2024.04.10

Chi Zhang

Oslo Center for Biostatistics and Epidemiology

chi.zhang@medisin.uio.no

Random variables

Discrete and continuous

Random variable: a quantity that can take random values, with a certain probability

Discrete variables

- coin tossing H,T
- birth boy, girl

Continuous variables

- weight and height
- age

Properties of probability:

Non-negative (0 or above), less than 1, prob of all outcomes sum up to 1.

Histogram and bar plot

X-axis is usually what **values** the variable can be - either a single value, or a range

Y-axis is the **frequency**; or **proportions** (probability) corresponding to that variable value (or range)

Histogram

Continuous variable

Probability distribution: relationship between outcome and probability

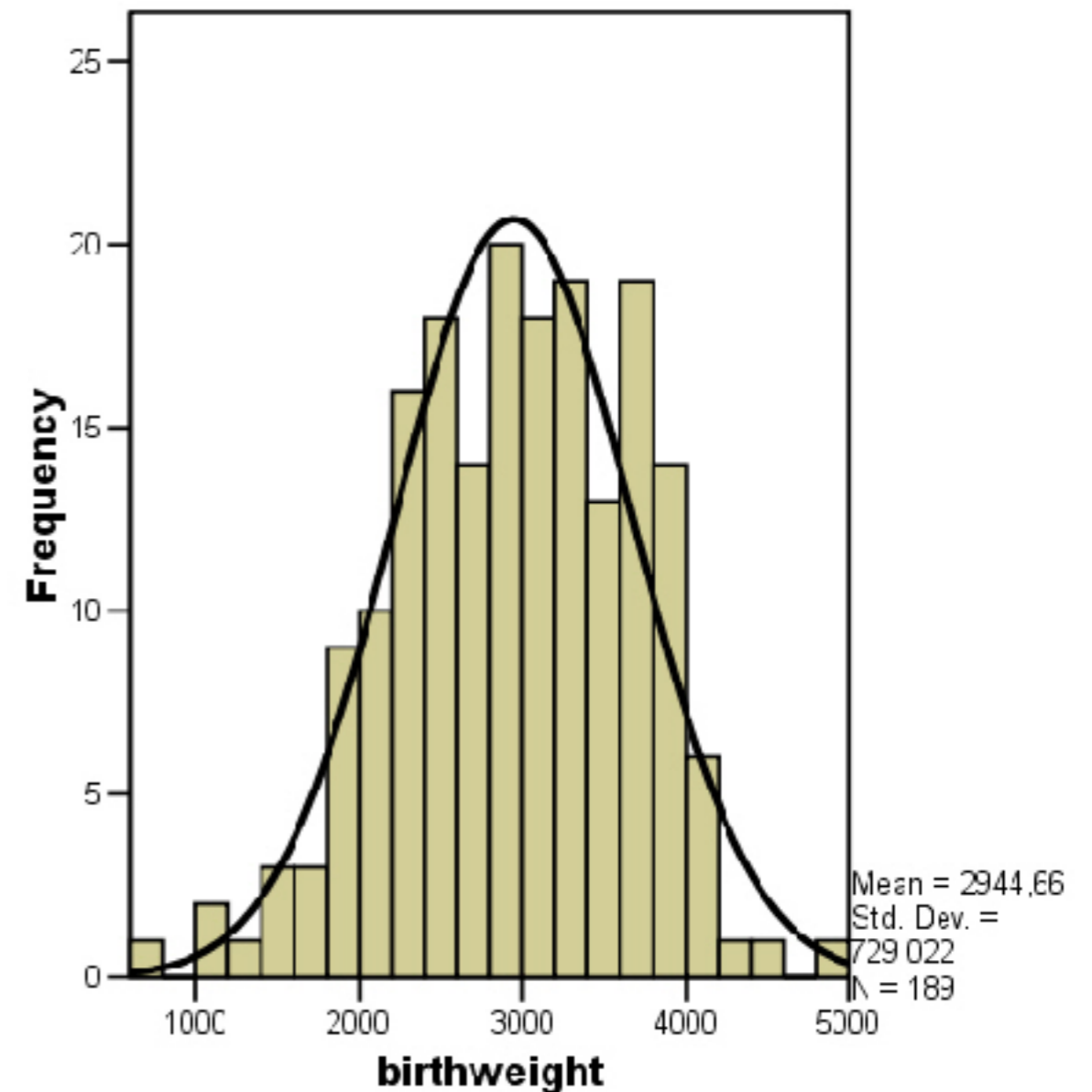
Outcome here is the **measurement (on a continuous scale)**, instead of categories

E.g. Birth weight of 189 newborns
3001, 2918, 3000, 3001,.....

How to find the probability? Try **counting elements in an interval**

- 20 measurement between 2800 and 3000;
- 30 between 3000 and 3200; for example.
- the interval (bin) can be bigger, or smaller

You can get the proportion for each interval by dividing the counts over N



Probability distribution

Probability density

A continuous probability distribution is defined by the **probability density function** $f(x)$ (a relationship between the measurement x and its **density** at this point)

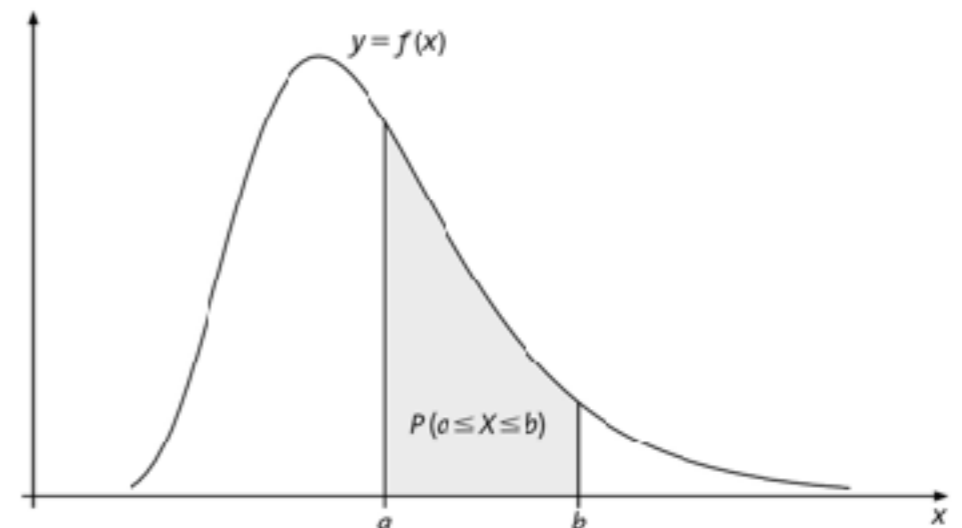
NOT a probability itself! $P(x) = 0$. Need integration

The following properties are important

$$f(x) \geq 0$$

The area under the curve in total is 1

$P(a \leq X \leq b)$ is the **area under the curve** from a to b



Figur 5.1 Sammenligningstetthet for en målevariabel

We are usually more interested in the probability of X **greater or smaller** than a certain value; or **between** two values.

e.g. we care about probability of birth weight between 3000 to 3010, but not equal to 3000 exactly.

Probability distribution

Probability density

Birthweight example: $N=189$

The area under the curve in total is 1

- if you sum all the bars (2 + 3 + 1 + ...)

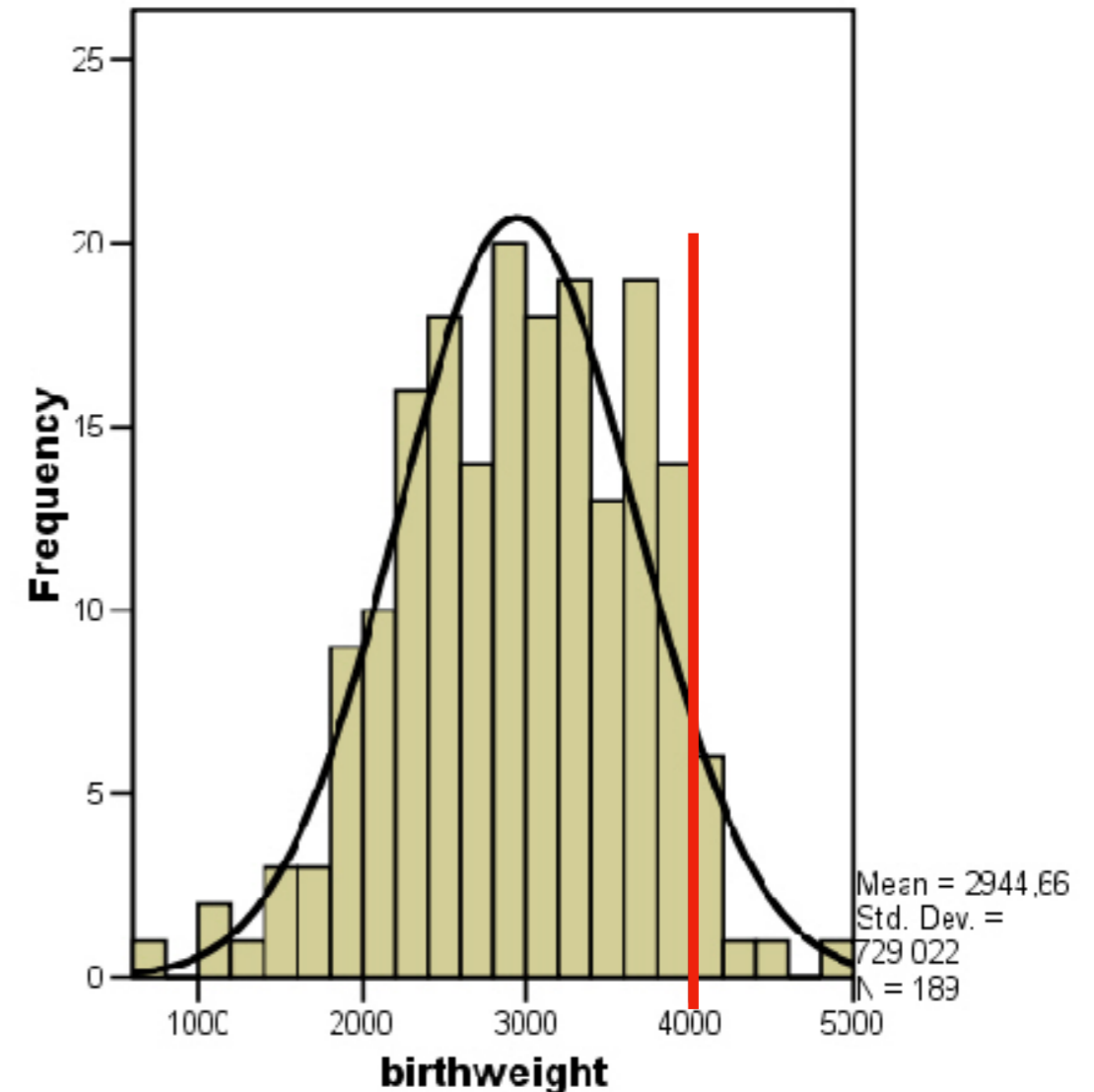
the total is 189

- (the **height** of each bar)

$P(a \leq X \leq b)$ is the area under the curve from a to b

- $P(2800 \leq X \leq 3000) = 20/189$

- $P(X > 4000) = 9/189$

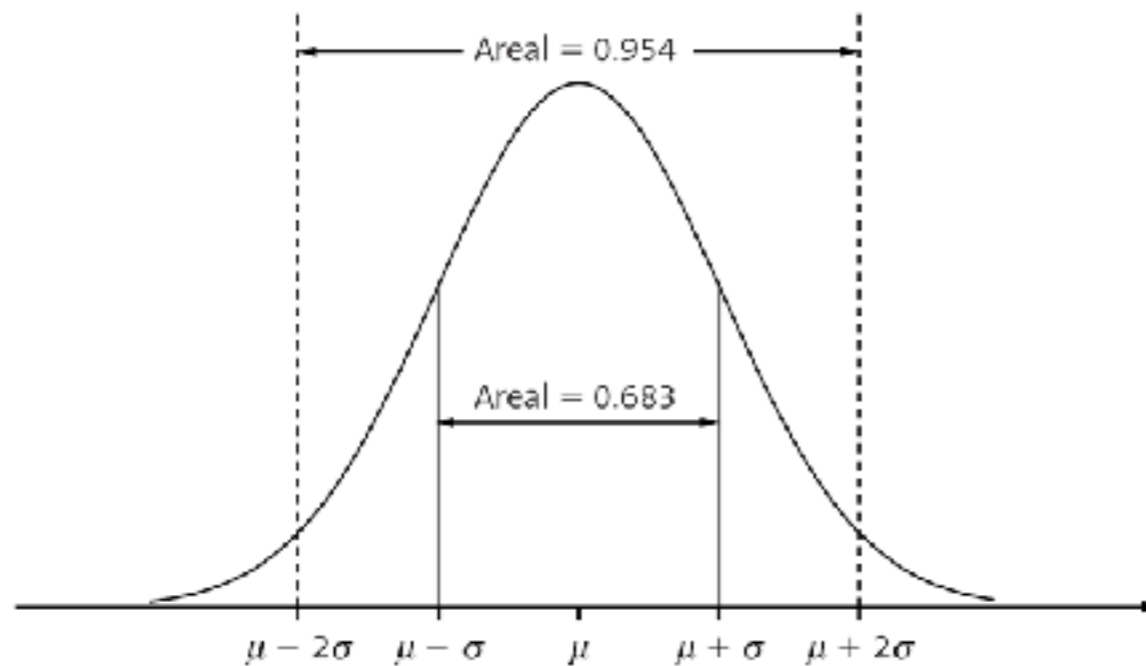


Normal distribution

Probability density function of the normal distribution is

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

where μ is the mean, σ the standard deviation and $\exp(a) = e^a$



Figur 5.4 Tegning av en normalfordeling. Det er avmerket at $\mu \pm \sigma$ dekker 68 %, mens $\mu \pm 2\sigma$ dekker 95 %.

$$X \sim N(\mu, \sigma^2)$$

Read: mu, sigma

Normal distribution

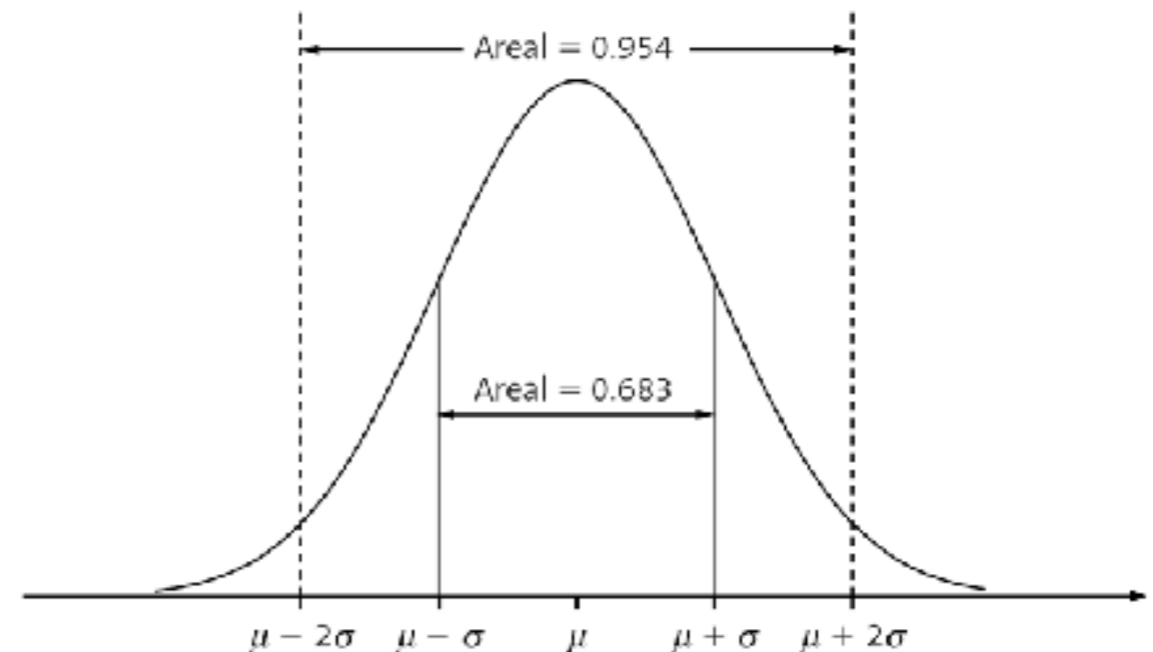
Symmetric, bell shape

Mean defines the location (where it is centered)

SD (sigma) defines the variation, i.e. spread

An interval with center at the mean value, going 2 standard deviation each way covers **approximately 95%** of the distribution

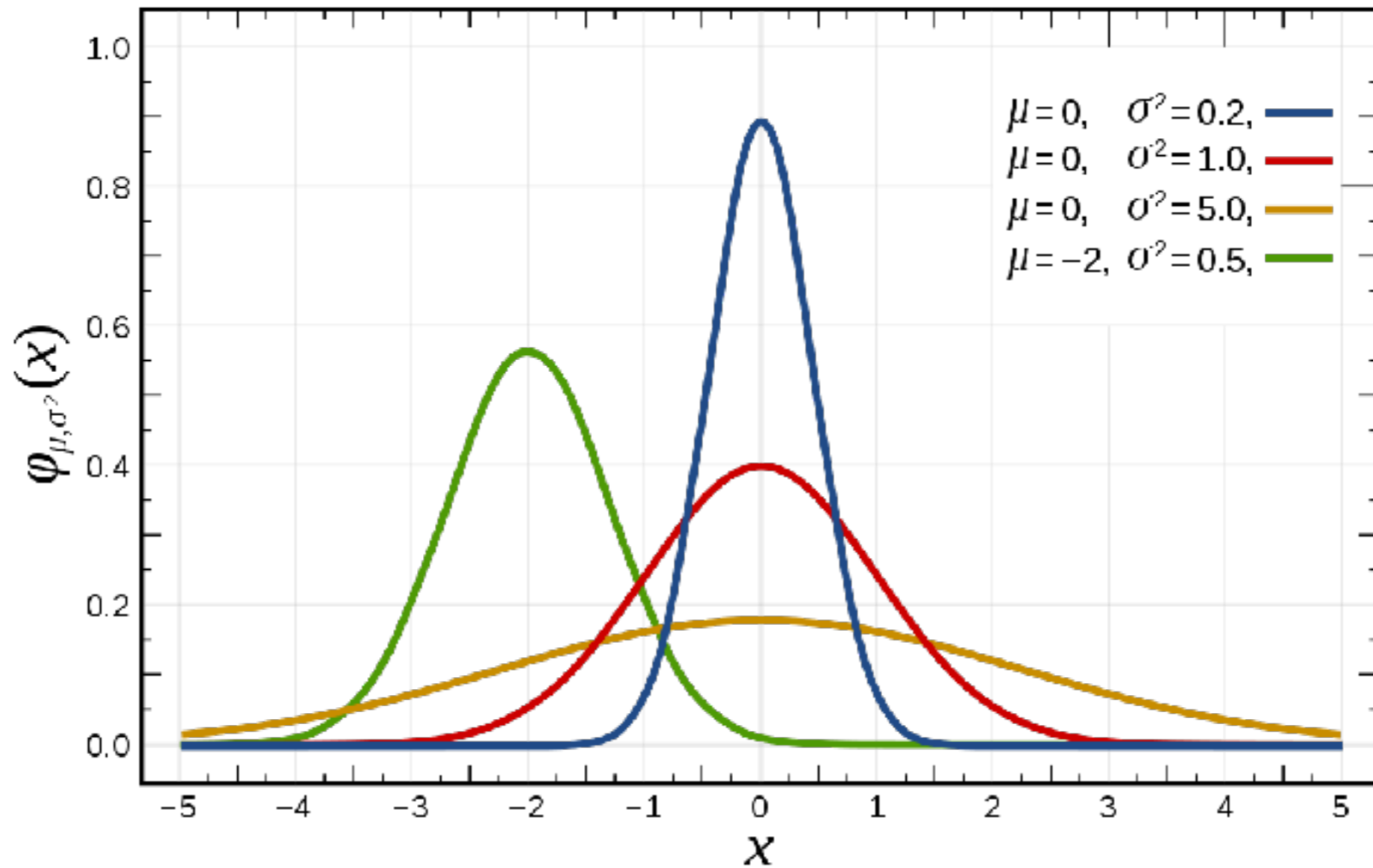
$$X \sim N(\mu, \sigma^2)$$



Figur 5.4 Tegning av en normalfordeling. Det er avmerket at $\mu \pm \sigma$ dekker 68 %, mens $\mu \pm 2\sigma$ dekker 95 %.

Normal distribution

Different locations and sd



Standard normal distribution

Definition

A **standard** normal distribution is a normal distribution with mean 0, variance 1

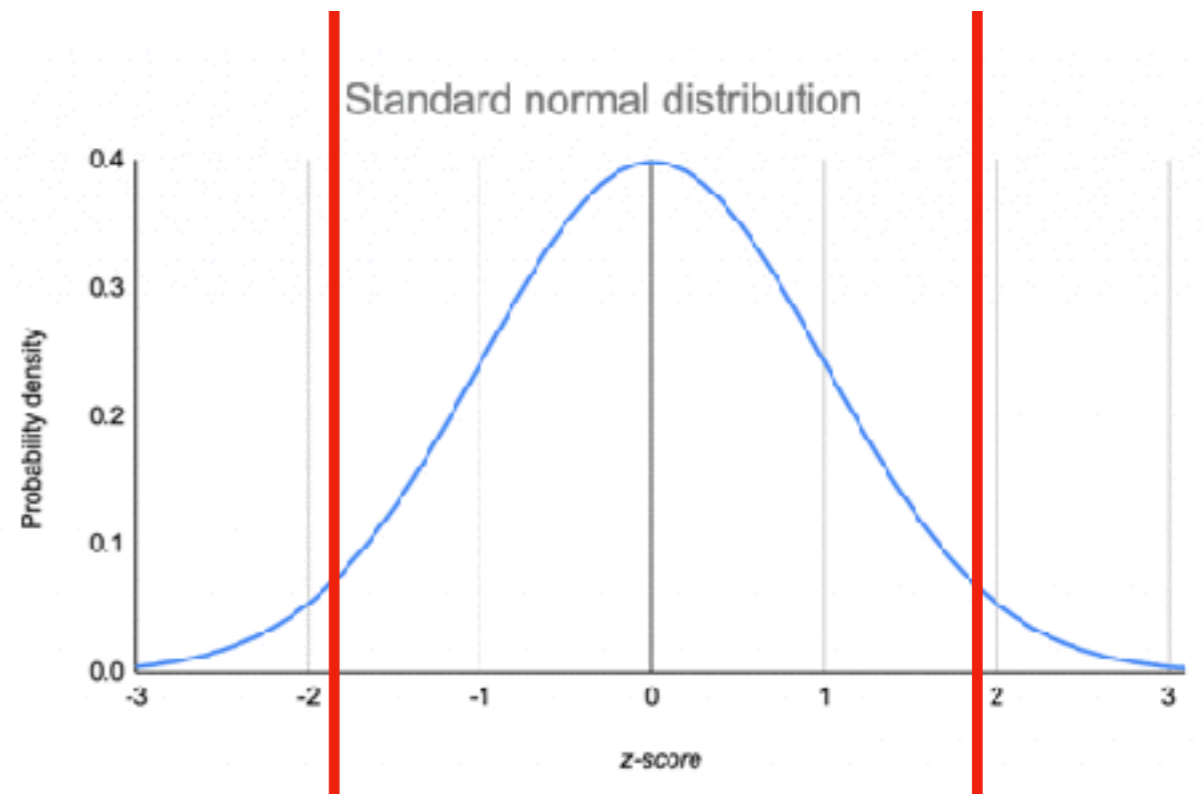
Noted as $N(0, 1)$

Any normal distribution can be transformed into a standard normal ...

By subtracting the mean, and dividing by the standard deviation

$$X \sim N(\mu, \sigma^2)$$

$$Y = \frac{X - \mu}{\sigma} \sim N(0, 1)$$



-1.96, 1.96 are 2.5% and 97.5% quantile for $N(0,1)$

Frequently used for computing 95% confidence intervals

Standard normal distribution

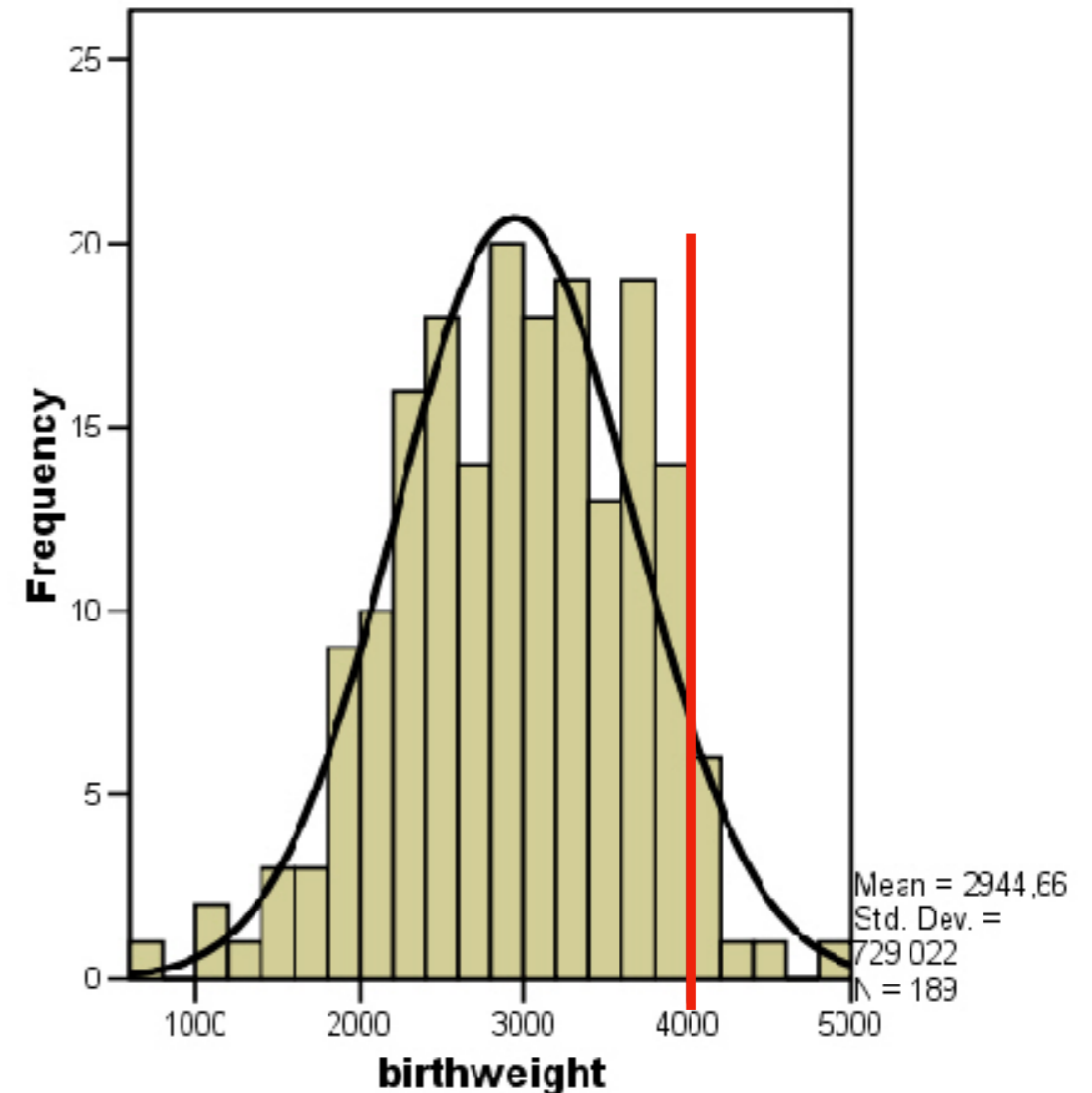
We calculated the 'empirical' probability of weights greater than 4000 by **counting**:

$P(X > 4000) = 9/189 = 4.7\%$, approximately

Try using the theoretical distribution (**curve**)

- mean (μ) = 2945, sd (σ) = 729
- transform X into Y, which is **N(0, 1)**
- then find the area for $P(Y > y)$

$$\begin{aligned} &P(X > 4000) \\ &= P[(X - 2945)/729 > (4000 - 2945)/729] \\ &= P[(X - 2945)/729 > 1.45] \\ &= P(Y > 1.45) \end{aligned}$$



Standard normal distribution

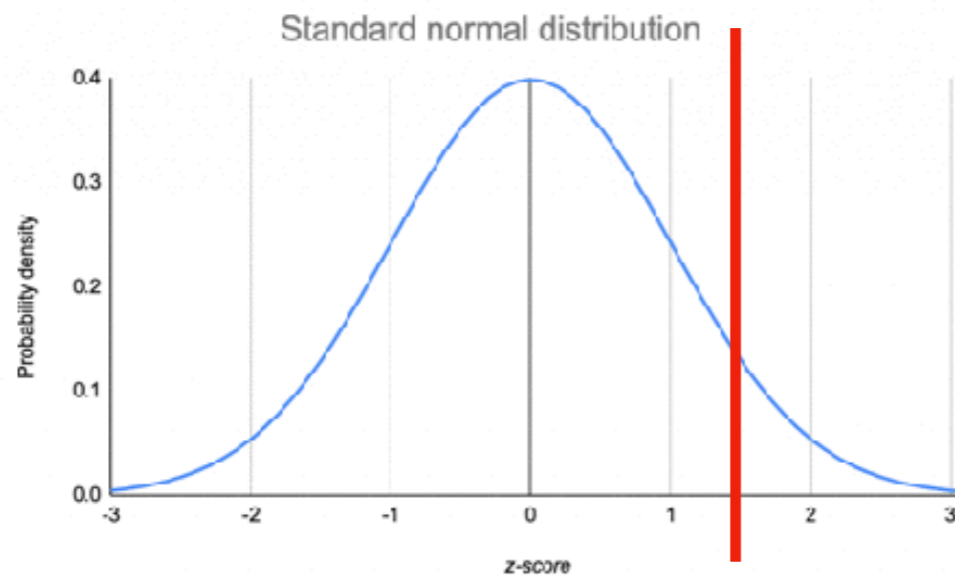
Find $P(Y > 1.45)$

Method II: in R

```
> pnorm(1.45, 0, 1)
[1] 0.9264707
> 1-pnorm(1.45, 0, 1)
[1] 0.07352926
```

Method III: in STATA

```
. display normal(1.45)
.92647074
. display 1-normal(1.45)
.07352926
```



Central limit theorem CLT

Important conclusion about the mean

Sample mean is normally distributed around the true, **population mean**, with 1 over n times the variance

Regardless of which distribution the population is.

- Binomial
- Poisson
- Uniform
- ...

This also applies for **sum** (n times mean; because mean is sum divided by n).

You will be able to make **inference** with CLT (relevant for confidence interval, t-test and more)

Central limit theorem CLT

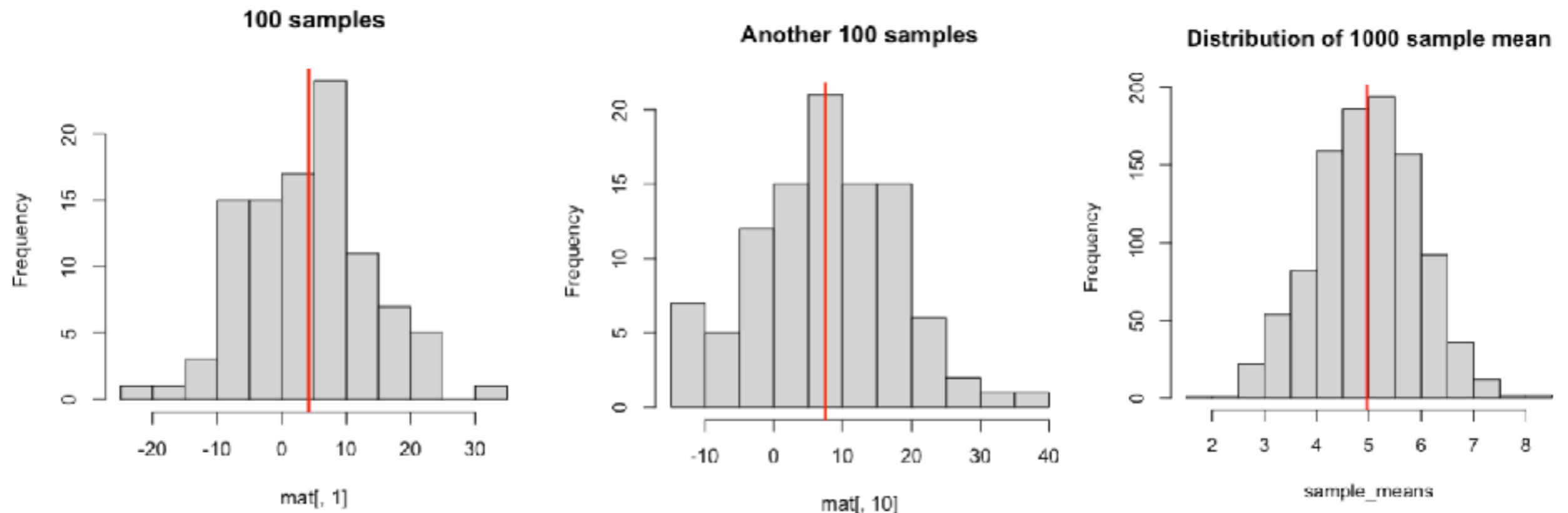
Draw 100 random samples with a known theoretical distribution: $N(5,100)$ - centered at 5, with variance 100 (standard deviation 10)

Mark their mean with red in figure 1, 2

Repeat the procedure 1000 times, compute 1000 means, plot the **sample means**

Centered at 5, with variance $100/100 = 1$ - a Normal Distribution!

First 100 is the sample size (100 random samples); second 100 is the variance



Central limit theorem CLT

Your data does not need to be normally distributed; it can be any shape

Their sample mean is normally distributed

This result is the foundation for day 3 and day 4, you will see that you can compute the range where your mean can vary!

