

# STA220H1: The Practice of Statistics I

Elena Tuzhilina

February 14, 2023

Please turn on your videos :)

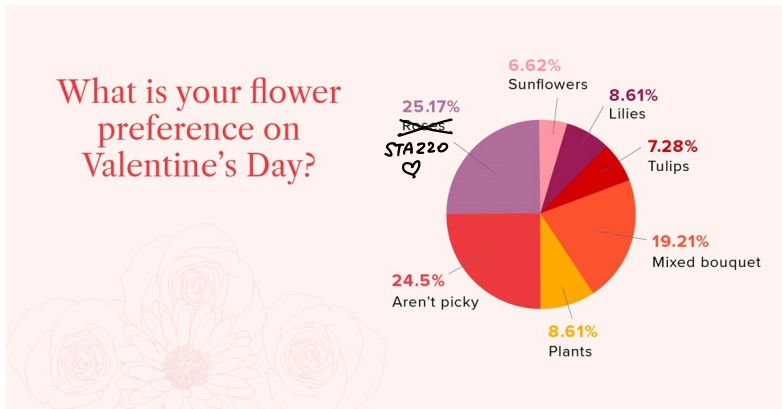


Figure 1: [picture source]

## Learning strategy

1. Attend lectures
2. Watch modules at <https://sta220.utstat.utoronto.ca>
3. Do practice sets, attend TAs office hours if something is not clear
4. Do Quiz, attend my office hours on Monday if something is still not clear
5. Post your questions on Piazza (not my personal email, pls!)

## Agenda for today

- ▶ Recap: expectation and variance
- ▶ Continuous random variables: uniform and normal distribution
- ▶ Sample mean distribution: normal sample and CLT

## Discrete: expected value

*Expected value measures the average value of random variable in long term.*

$$E(X) = \sum_x x \cdot P(X = x)$$

*Variance and standard deviation measure the spread of the values of a random variable in long term.*

$$\text{Var}(X) = \sum_x (x - E(X))^2 \cdot P(X = x)$$

$$\text{sd}(X) = \sqrt{\text{Var}(X)}$$

## Exercise

$$X = \begin{cases} 0 \\ 1 \end{cases} \quad P(X=1) = 0.1$$

What is the expectation and variance of Bernoulli random variable with  $p = 0.1$ ? What is the formula for general  $p$ ?

$X$	$1$	$0$
$P(X)$	$p$	$1-p$

$$\begin{aligned} E(X) &= 1 \cdot 0.1 + 0 \cdot 0.9 = \\ &= 0.1 \end{aligned}$$

$$\begin{aligned} \text{Var}(X) &= (1-0.1)^2 \cdot 0.1 + \\ &\quad (0-0.1)^2 \cdot 0.9 = \\ &= 0.09 \end{aligned}$$

$$E(X) = 1 \cdot p + 0 \cdot (1-p) = p$$

$$\begin{aligned} \text{Var}(X) &= (1-p)^2 \cdot p + (0-p)^2 \cdot (1-p) = \\ &= (1-p) \cdot p (1-p + p) = p \cdot (1-p) \end{aligned}$$

## Expected value vs. sample mean

$$s_x^2 \approx \text{var}(x)$$

$$\bar{x} \approx E(x)$$

$$E(x) = 0.1$$
$$E(x) = 0 \cdot 0.9 + 1 \cdot 0.1$$

- We have a random variable  $X = \begin{cases} 1 \\ 0 \end{cases}$   $P(X=1) = 0.1$

1 0 1 0 0 1 0 0 1 ... 0  
└────────────────────────────────┘  
100

- We generate a sample of size  $n$  using this random variable  
 $x_1, \dots, x_n$

$$\begin{array}{|c|c|} \hline 0 & 1 \\ \hline 91 & 9 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 91/100 & 9/100 \\ \hline \end{array}$$

What is the relationship between  $E(X)$  and  $\bar{x} = \frac{x_1 + \dots + x_n}{n}$ ?

$$\bar{x} = \frac{1+0+1+\dots+1+\dots+0}{100} = \frac{9 \cdot 1 + 91 \cdot 0}{100} =$$
$$= 1 \cdot \left(\frac{9}{100}\right) + 0 \cdot \left(\frac{91}{100}\right) = 1 \cdot (\text{prop of } 1) + 0 \cdot (\text{prop of } 0)$$

## Important rules

### Expectation

- ▶ If  $X$  is a random variable and  $a, b$  are some numbers then

$$E(a \cdot X + b) = a \cdot E(X) + b \quad E(\theta) = \theta$$

- ▶ If  $Y$  is also a random variable then

$$E(X + Y) = E(X) + E(Y)$$
$$E(X - Y) = E(X) - E(Y)$$

### Variance

- ▶ If  $X$  is a random variable and  $a, b$  are some numbers then

$$\text{Var}(-X) = \text{Var}(X) \quad \text{Var}(a \cdot X + b) = a^2 \cdot \text{Var}(X) \quad \text{Var}(b) = 0$$

- ▶ If  $Y$  is also a random variable and it is independent of  $X$  then

$$\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y)$$
$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$



## Exercise

$$y = X_1 + \dots + X_n \sim \text{Bin}(n, p) \quad X_i \sim \text{Bern}(0.1)$$

Ⓢ If  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  then  
 $y = X_1 + \dots + X_n \sim \text{Binomial}(n, p)$

What is the expectation and variance of Binomial random variable with  $n = 5$  and  $p = 0.1$ ? What is the formula for general  $n$  and  $p$ ?

$$\begin{aligned} E(y) &= E(X_1 + \dots + X_n) = \\ &= E(X_1) + \dots + E(X_n) = \\ &= \underbrace{p + \dots + p}_n = n \cdot p \end{aligned}$$

$$\begin{aligned} \text{Var}(y) &= \text{Var}(X_1 + \dots + X_n) = \\ &= \text{Var}(X_1) + \dots + \text{Var}(X_n) = \\ &= \underbrace{p \cdot (1-p) + \dots + p \cdot (1-p)}_n = n \cdot p(1-p) \end{aligned}$$

## Random variable: discrete

**Discrete** - takes one of a countable list of distinct values

- ▶ sex of a baby
- ▶ number of heads when tossing ten coins
- ▶ number of zeroes in your student ID

Special cases:

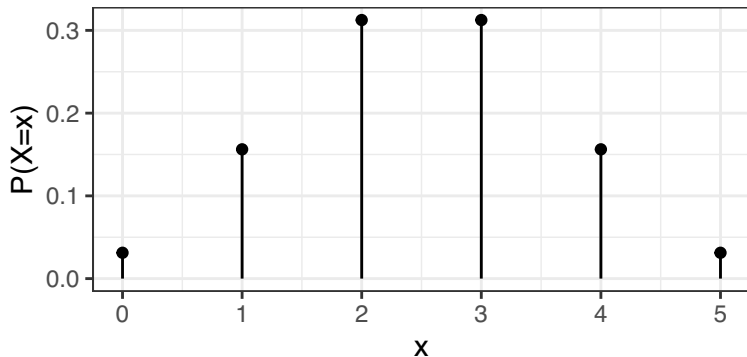
$$X \sim \text{Bernoulli}(p) = \begin{cases} 0 \\ 1 \end{cases} \quad P(X=1) = p$$

$$X \sim \text{Binomial}(n, p) = \begin{array}{l} \# \text{ successes} \\ \text{among } n \text{ exp} \\ 0, 1, \dots, n \end{array}$$

## Distribution: discrete

Two ways to present the distribution: a table and a plot

$X$	0	1	2	3	4	5
$P(X)$	0.03125	0.15625	0.3125	0.3125	0.15625	0.03125



## Random variable: continuous

**Continuous** - takes any value in an interval or collection of intervals

- ▶ the wait time for the next bus  $[0, 10]$
- ▶ birth weight of a baby  $[2000g, 5000g]$   $2000.1g$   
 $2000.01g$
- ▶ life expectancy in Canada  $[40\text{years}, 100\text{years}]$

Special cases:

$$X \sim \text{Uniform}(l, u)$$

$$X \sim \text{Normal}(\mu, \sigma^2)$$

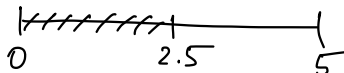
$[40.05\text{ years}]$

## Continuous: uniform

You live in a building that has an elevator. Once you push the button to call the elevator, it takes between 0 and 5 seconds (equally likely) for the elevator to arrive.

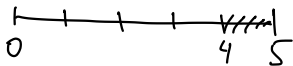
$X = \text{wait time (in seconds)} \sim \text{Uniform}(0, 5)$

lower  
↓  
upper ↙



$$P(0 \leq X \leq 2.5) = 0.5$$

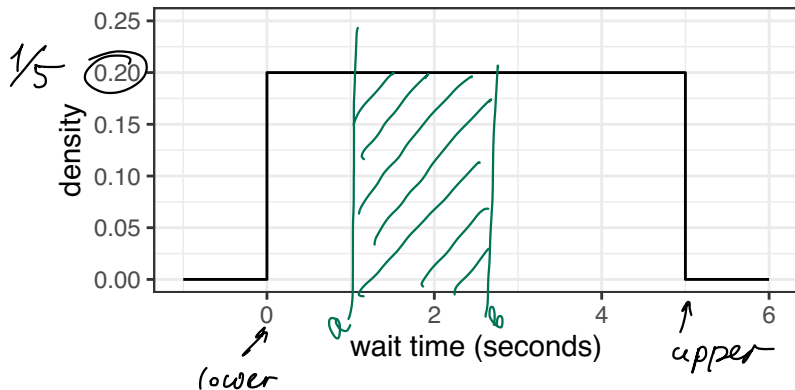
$$P(4 \leq X \leq 5) = 0.2$$



## Continuous: uniform

Since continuous variables can take *so many* possible values, we cannot use distribution tables. Instead we use **probability density functions**.

$P(a \leq X \leq b) =$  area under the curve bounded by  $a$  and  $b$  vertical lines



## Exercise

$$\approx 0.2 \cdot (4.5 - 2.5) = 0.4$$

Find  $P(2.5 \leq X \leq 4.5)$  and  $P(X \geq 4)$ .

$$\leq 0.2(5 - 4) = 0.2$$

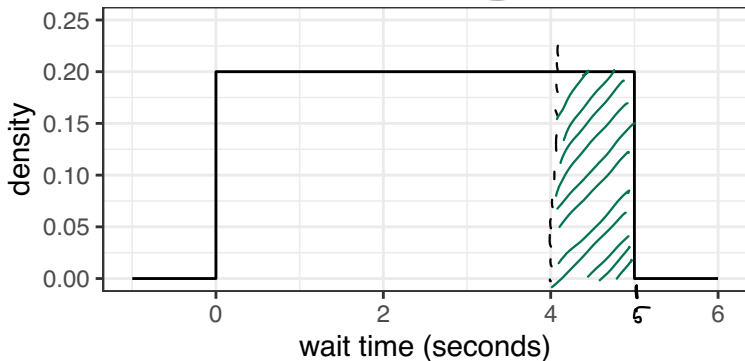
What is  $P(X = 4)$ ? What is  $P(X > 4)$ ?

$\approx 0$

$\approx 0.2$

$$P(X \geq 4) = P(X = 4) + P(X > 4)$$

$\approx 0$



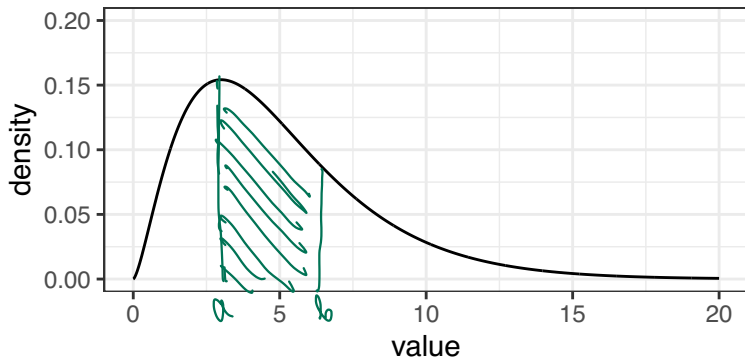
$$\text{total area} = 0.2 \cdot 5 = 1$$

## Continuous: density function

A curve is a valid **density function** if:

- ▶ It is greater or equal to 0
- ▶ The total area under the curve is 1

$P(a \leq X \leq b) =$  area under the curve bounded by  $a$  and  $b$  vertical lines

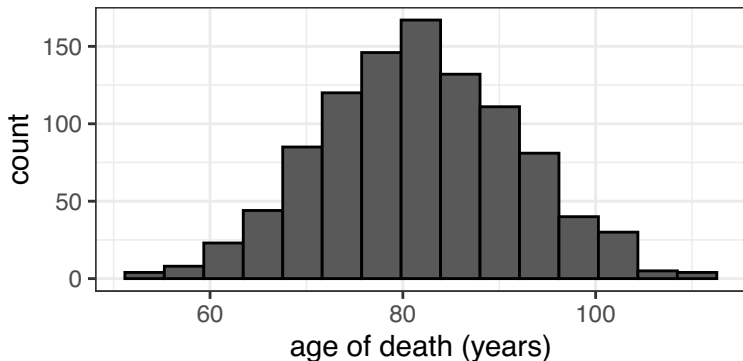




## Histogram vs density curve

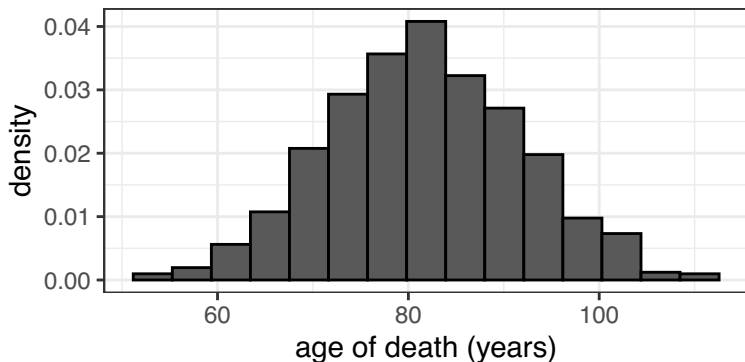
Suppose you are interested in the life expectancy in Canada. You record the age of death for 1000 Canadians.

```
## [1] 72.03085 82.84849 96.87845 69.69624 80.19748 82.32420
```



## Histogram vs density curve

You can convert a histogram to an approximate density by changing the scale of y-axis.



## Histogram vs density curve

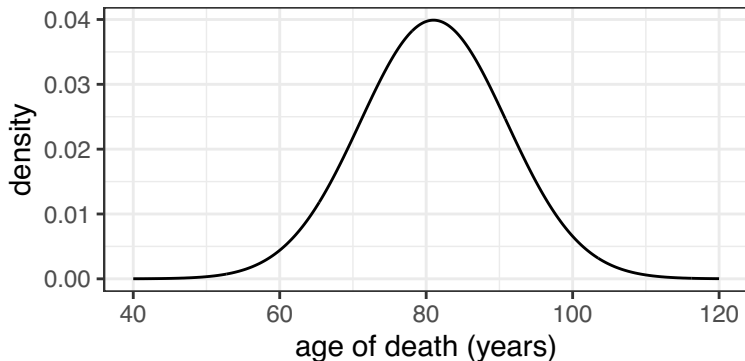
The “smoothed” version of your histogram is density.



## Continuous: normal

**Normal** random variable (or Gaussian) has symmetric, bell-shaped and unimodal distribution.

$$X = \text{age of death (in years)} \sim \text{Normal}(81, 100)$$



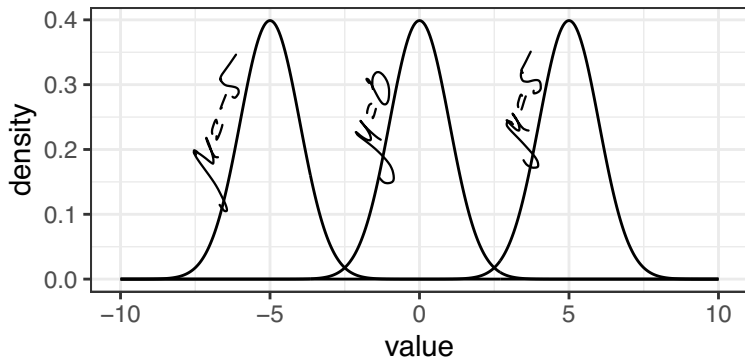
## Continuous: normal

Normal distribution  $X \sim \text{Normal}(\mu, \sigma^2)$  has **two parameters**

$$\mu = E(X) \text{ and } \sigma^2 = \text{Var}(X)$$

$$\left\{ \sigma = \text{Sd}(X) \right\}$$

- ▶  $\mu$  controls the “center” of the distribution

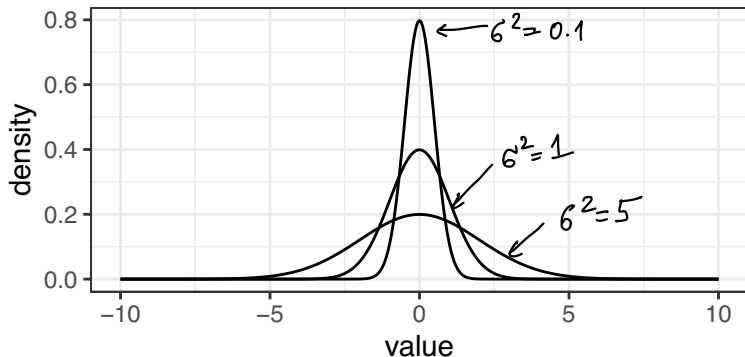


## Continuous: normal

Normal distribution  $X \sim \text{Normal}(\mu, \sigma^2)$  has **two parameters**

$$\mu = E(X) \text{ and } \sigma^2 = \text{Var}(X)$$

- ▶  $\sigma^2$  controls the “spread” of the distribution

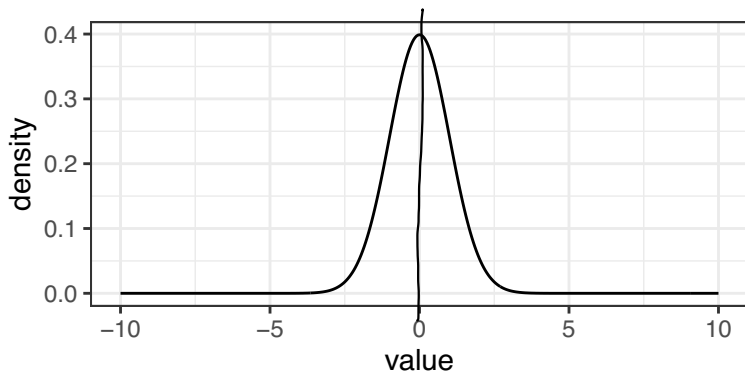


## Standard normal

$$X \sim \text{Normal}(0, 1)$$

$$\left\{ E(x) = 0 \quad \text{Sd}(x) = 1 \right\}$$

**Standard normal** distribution has  $\mu = 0$  and  $\sigma^2 = 1$ .

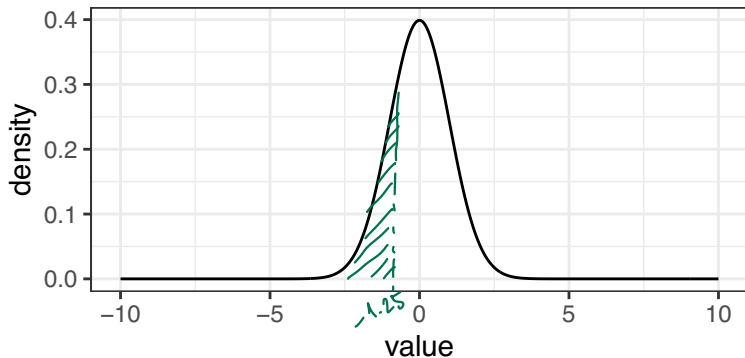


## Standard normal

To find the probabilities for standard normal we use the distribution table.

$$1 - P(X \leq -1.25)$$

If  $X \sim \text{Normal}(0, 1)$  what is the probability  $P(X \leq -1.25)$ ?



```
pnorm(-1.25)
```

```
## [1] 0.1056498
```

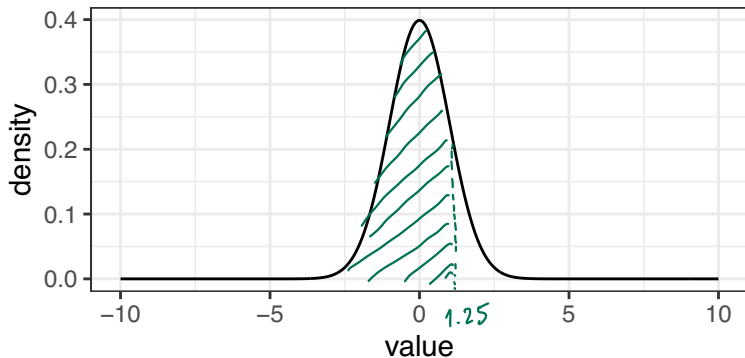


## Standard normal

To find the probabilities for standard normal we use the distribution table.

//  $1 - P(X \leq -1.25)$

If  $X \sim \text{Normal}(0, 1)$  what is the probability  $P(X \leq 1.25)$ ?



```
pnorm(1.25)
```

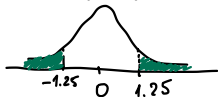
```
## [1] 0.8943502
```

## Exercise

$$P(X \leq -1.25) = 0.1056$$

$$P(X \leq 1.25) = 0.8944$$

If  $X \sim \text{Normal}(0, 1)$  what is the probability  $P(X > 1.25)$ ?



If  $X \sim \text{Normal}(0, 1)$  what is the probability  $P(X = 1.25)$ ?

If  $X \sim \text{Normal}(0, 1)$  what is the probability  $P(X \geq 1.25)$ ?

$$P(X \leq -1.25)$$

$$P(X > 1.25)$$

$$P(X \geq 1.25)$$

$$0$$

$$P(X > 1.25)$$

## Normal: properties

$$X \sim \text{Normal}(\mu, \sigma^2)$$

$$P(\mu - \sigma \leq X \leq \mu + \sigma) = 0.68$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 0.95$$

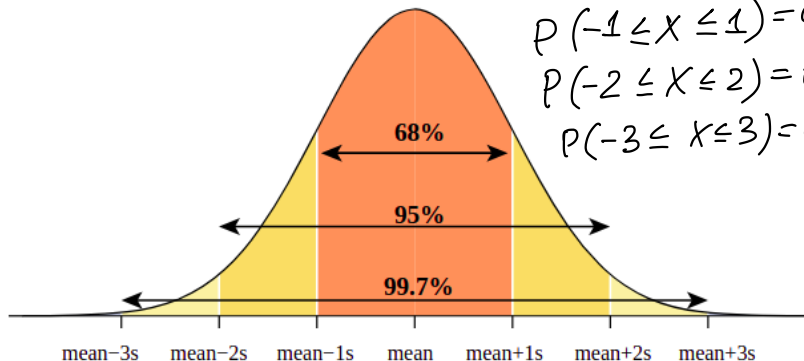
$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) = 0.997$$

$$X \sim \text{Normal}(0, 1)$$

$$P(-1 \leq X \leq 1) = 0.68$$

$$P(-2 \leq X \leq 2) = 0.95$$

$$P(-3 \leq X \leq 3) = 0.997$$



## Normal: properties

$$\begin{array}{l} E(x+y) = E(x) + E(y) \\ \text{Var}(x+y) = \text{Var}(x) + \text{Var}(y) \end{array} \left| \begin{array}{l} E(ax+b) = a \cdot E(x) + b \\ \text{Var}(ax+b) = a^2 \cdot \text{Var}(x) \end{array} \right.$$

- ▶ if  $X$  is normal then  $a \cdot X + b$  is also normal

$$X \sim \text{Normal}(\mu, \sigma^2) \Rightarrow E(X) = \mu \quad \text{Var}(X) = \sigma^2$$

$$Y = a \cdot X + b \sim \text{Normal}(a\mu + b, a^2 \cdot \sigma^2)$$

- ▶  $X$  and  $Y$  are independent and normal then  $X + Y$  is normal

$$X \sim \text{Normal}(\mu_x, \sigma_x^2)$$

$$Y \sim \text{Normal}(\mu_y, \sigma_y^2)$$

$$X + Y \sim \text{Normal}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

## Standard normal

$\sim \text{Normal}(0, 1)$

**Standardization:** if  $X \sim \text{Normal}(\mu, \sigma^2)$  then  $\frac{X-\mu}{\sigma}$  is standard normal.

For example, if life expectancy  $\overset{\text{Normal}(0,1)}{\uparrow}$   $y = \frac{X-\mu}{\sigma} \sim \text{Normal}(0, 1)$

$$X \sim \text{Normal}(81, 100) \quad E(y) = E\left(\frac{X-\mu}{\sigma}\right) = E\left(\underbrace{\frac{1}{\sigma}}_a \cdot X - \underbrace{\frac{\mu}{\sigma}}_b\right) =$$

then

$$P(X \leq 75) = a \cdot E(X) + b = \frac{1}{\sigma} \cdot E(X) - \frac{\mu}{\sigma} =$$

```
pnorm(75, mean = 81, sd = 10)
```

$$= P\left(\frac{X-81}{10} \leq \frac{75-81}{10}\right) = \frac{\mu}{\sigma} - \frac{\mu}{\sigma} = 0$$
$$= P(Y \leq -0.6)$$

```
## [1] 0.2742531
```

$$\text{Var}(y) = \text{var}\left(\underbrace{\frac{1}{\sigma}}_a X - \underbrace{\frac{\mu}{\sigma}}_b\right) = a^2 \cdot \text{Var}(X) = \frac{1}{\sigma^2} \cdot \sigma^2 = 1$$

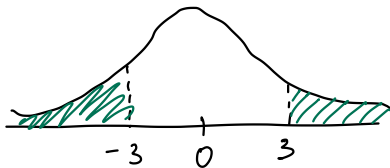
## Exercise

The life expectancy in Canada follows normal distribution with <sup>expectation</sup> mean 81 and standard deviation 3. What is the probability to live longer than 90 years?

$$X \sim \text{Normal}(81, 9)$$

$$P(X \geq 90) = P\left(\frac{X-81}{3} \geq \frac{90-81}{3}\right) = P(Y \geq 3) \quad Y \sim \text{Normal}(0, 1)$$

What is the life expectancy corresponding to the first quartile?



$$P(Y \geq 3) = P(Y \leq -3) = 0.0044$$

## Sample mean distribution

The theoretical model of the life expectancy is

$X \sim \text{Normal}(81, 100)$ . We record the ages of death for ~~the~~ 5 Canadians and average them.

$$X \rightarrow x_1, x_2, x_3, x_4, x_5 \rightarrow \bar{x} = \frac{x_1 + x_2 + x_3 + x_4 + x_5}{5}$$

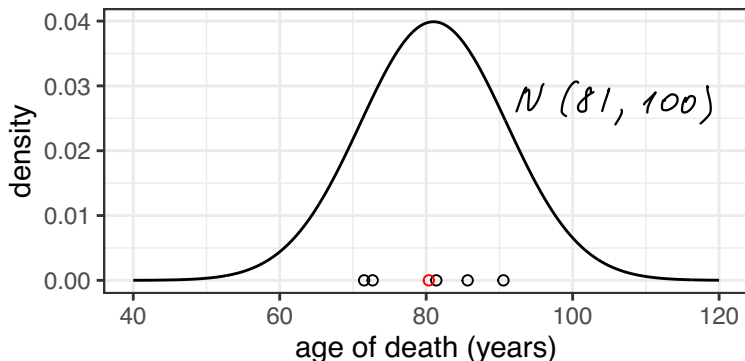
## Sample mean distribution

```
ages = rnorm(5, mean = 81, sd = 10)
ages
```

```
## [1] 85.66200 90.54666 71.52794 81.38563 72.69118
```

```
mean(ages)
```

```
## [1] 80.36268
```





## Sample mean distribution

```
ages = rnorm(5, mean = 81, sd = 10)  
ages
```

```
## [1] 86.57172 79.87897 71.42362 76.31927 91.44757
```

```
mean(ages)
```

```
## [1] 81.12823
```



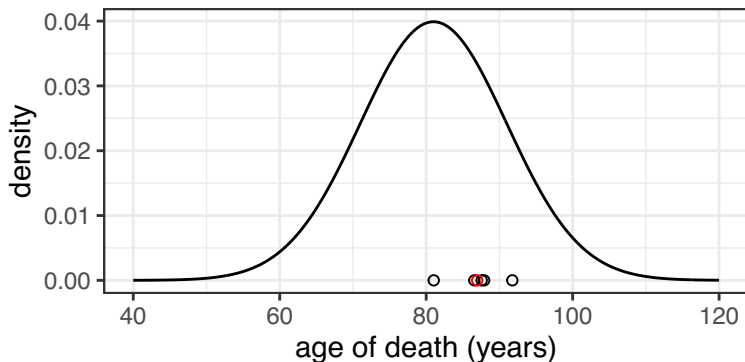
## Sample mean distribution

```
ages = rnorm(5, mean = 81, sd = 10)
ages
```

```
## [1] 81.03306 91.76645 86.58143 87.91968 87.57774
```

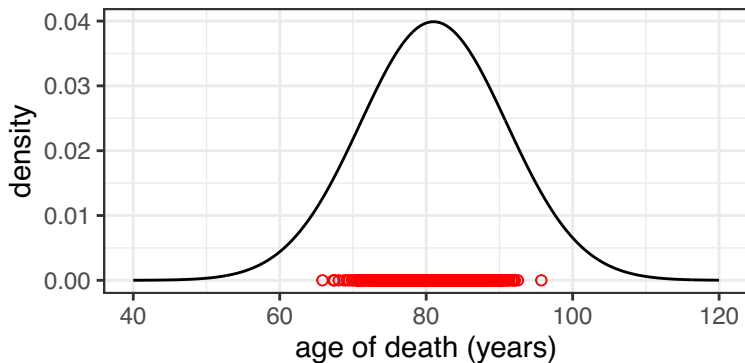
```
mean(ages)
```

```
## [1] 86.97567
```



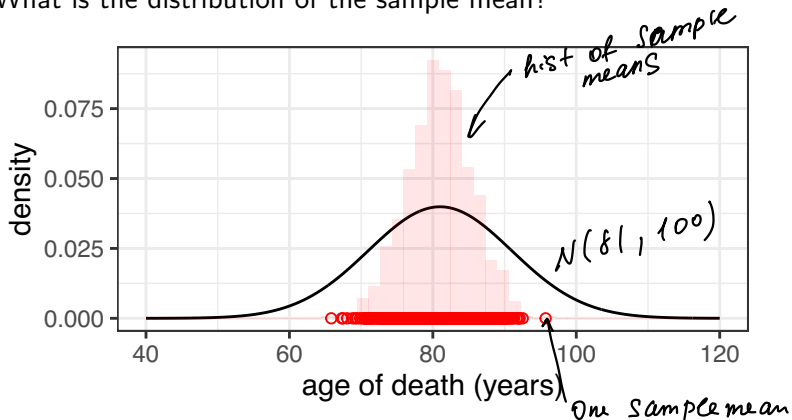
## Sample mean distribution

Let's repeat this experiment many times.



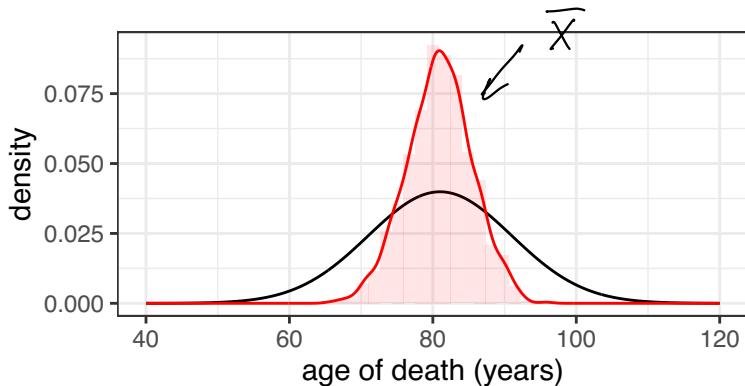
## Sample mean distribution

What is the distribution of the sample mean?



## Sample mean distribution

What is the distribution of the sample mean?



## Exercise

$$\begin{array}{ccccccccc} X_1 & , & X_2 & , & X_3 & , & X_4 & , & X_5 & \sim N(81, 100) \\ \downarrow & & \downarrow & & & & & & \downarrow & \\ x_1 & & x_2 & & \dots & & & & x_5 & \rightarrow \bar{x} \end{array}$$

We have 5 observations, each of them were generated from  $Normal(81, 100)$ .

If  $X_1, \dots, X_5 \sim Normal(81, 100)$  and independent, then

$$\bar{X} = \frac{X_1 + \dots + X_5}{5} \sim Normal(81, 20).$$

$$X \sim N(81, 100) \rightarrow x_1, x_2, x_3, x_4, x_5 \rightarrow \bar{x}$$

$$x_1, x_2, x_3, x_4, x_5 \rightarrow \bar{x}$$

$$x_1, x_2, x_3, x_4, x_5 \rightarrow \bar{x}$$

$$x_1, x_2, x_3, x_4, x_5 \rightarrow \bar{x}$$

$$\begin{array}{l} \bar{X} = \frac{X_1 + \dots + X_5}{5} \rightarrow \bar{x} \\ \rightarrow \bar{x} \\ \rightarrow \bar{x} \end{array}$$

## Exercise

We have 5 observations, each of them were generated from  $Normal(81, 100)$ .

If  $X_1, \dots, X_5 \sim Normal(81, 100)$  and independent, then

$$\bar{X} = \frac{X_1 + \dots + X_5}{5} \sim Normal(81, 20)$$

$$\begin{aligned} E(\bar{X}) &= E\left(\frac{X_1 + \dots + X_5}{5}\right) = \frac{1}{5} E(X_1 + \dots + X_5) = \\ &= \frac{1}{5} \cdot (E(X_1) + \dots + E(X_5)) = \frac{1}{5} (81 + \dots + 81) = 81 \end{aligned}$$

$$\begin{aligned} Var(\bar{X}) &= Var\left(\frac{X_1 + \dots + X_5}{5}\right) = \frac{1}{5^2} Var(X_1 + \dots + X_5) = \\ &= \frac{1}{25} \cdot [Var(X_1) + \dots + Var(X_5)] = \frac{1}{25} (100 + \dots + 100) = \frac{100}{5} \end{aligned}$$

## Sample mean distribution

We have  $n$  observations, each of them were generated from  $Normal(\mu, \sigma^2)$ .

If  $X_1, \dots, X_n \sim Normal(\mu, \sigma^2)$  and independent, then

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} \sim Normal\left(\mu, \frac{\sigma^2}{n}\right).$$



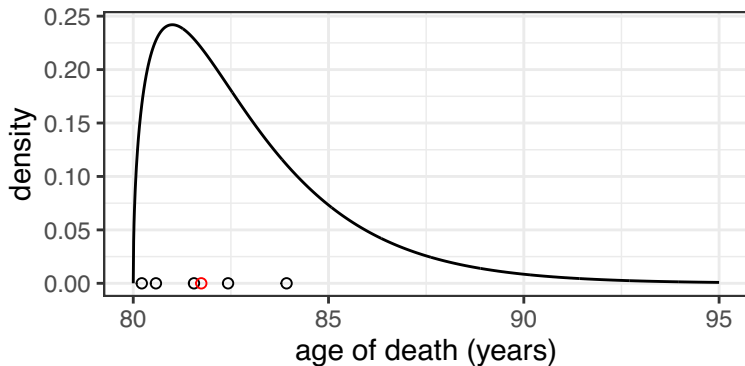
## Sample mean distribution

What if the distribution of age is not normal?



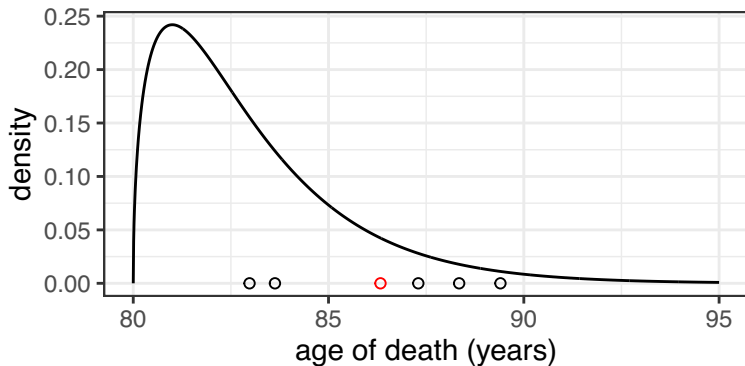
## Sample mean distribution

What if the distribution of age is not normal?



## Sample mean distribution

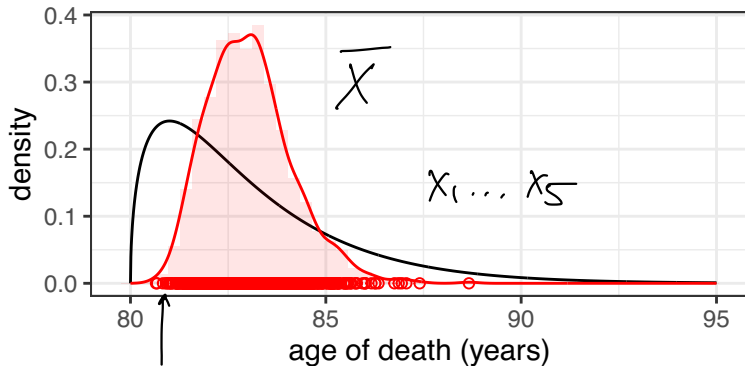
What if the distribution of age is not normal?



## Sample mean distribution

What if the distribution of age is not normal?

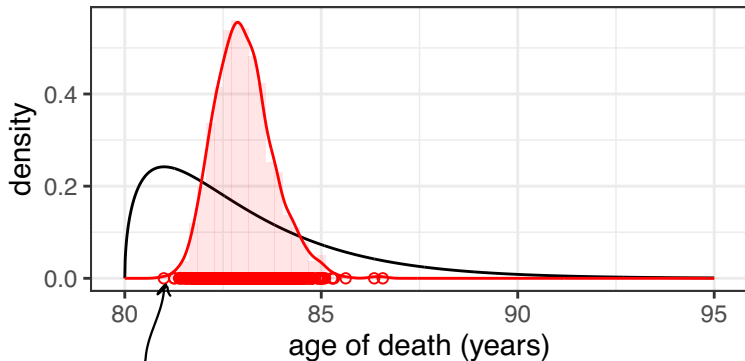
$$n = 5$$



$$\frac{x_1 + \dots + x_5}{5}$$

## Sample mean distribution

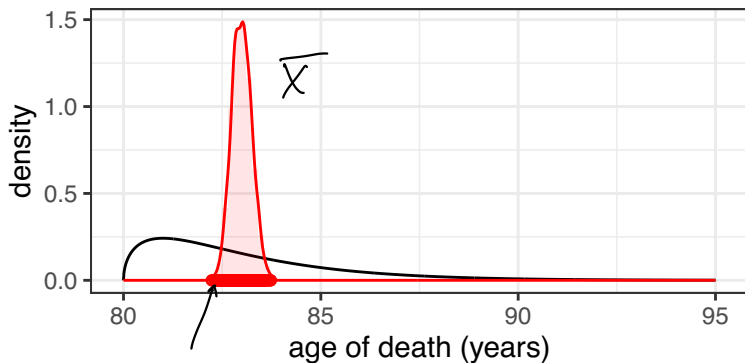
Now, let's increase the sample size from 5 to 10.  $n = 10$



$$\frac{x_1 + \dots + x_{10}}{10}$$

## Sample mean distribution

Now, let's increase the sample size from 10 to 100.  $n = 100$



$$\frac{x_1 + \dots + x_{100}}{100}$$

## Central limit theorem

If  $X_1, \dots, X_n$  have the same distribution (not necessary normal!) with  $E(X_i) = \mu$  and  $Var(X_i) = \sigma^2$ , then for  $n$  large enough

- ▶ the distribution of  $\bar{X} = \frac{X_1 + \dots + X_n}{n}$  is well approximated by normal
- ▶ the expectation of normal distribution is  $\mu$
- ▶ the variance of normal distribution is  $\frac{\sigma^2}{n}$

$$\bar{X} \text{ approximately } \sim \text{Normal} \left( \mu, \frac{\sigma^2}{n} \right) \quad \begin{array}{l} n \uparrow \\ \frac{\sigma^2}{n} \downarrow \end{array}$$

## Central limit theorem

CLT works even if  $X_1, \dots, X_n$  have discrete distribution.

For example, if  $X_i \sim \text{Bernoulli}(p)$ , then

$Y = X_1 + \dots + X_n \sim \text{Binomial}(n, p)$  total number of successes.

The proportion of successes  $\bar{X} = \frac{Y}{n}$  has approximately normal distribution for large  $n$ .



## Exercise

What are the parameters of this normal distribution?

# TO DO

1. Module 2. Probability: Random Variables and Module 3. Sampling Distributions
2. Quiz 6 due Monday (February 20) @ 11:59 PM (EST)
3. Practice Problem Set 6