



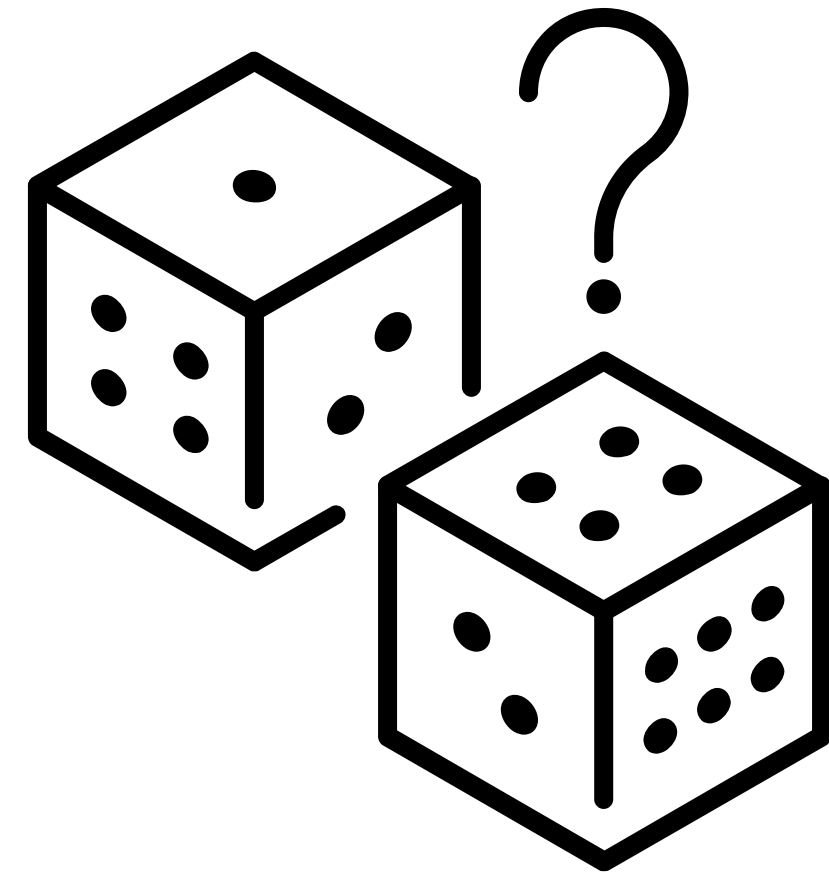
WELCOME



# STATISTIC WITH PYTHON

EMBARKING ON A JOURNEY  
INTO DATA SCIENCE

YA MANON



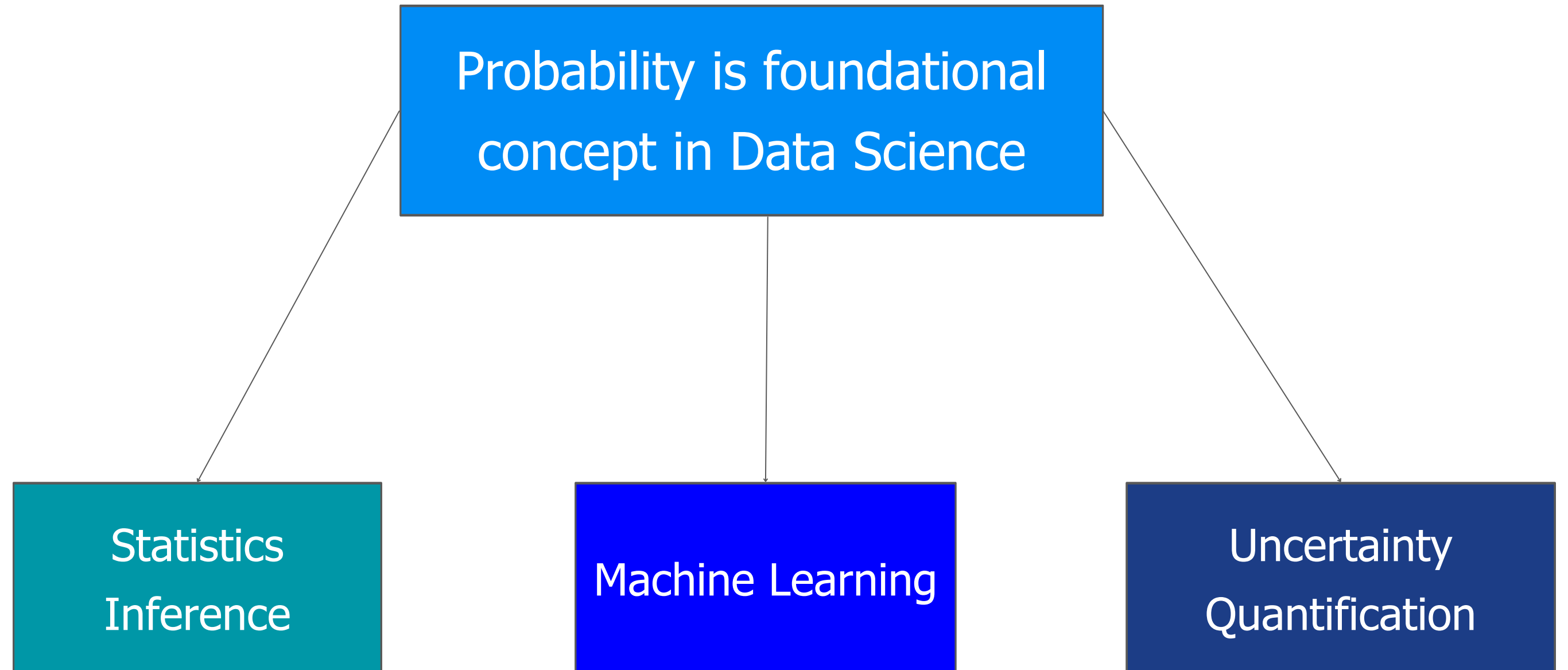
# PROBABILITY REVIEW

# WHY PROBABILITY?

Why Probability?

Calculate Probab

Normal Distribution

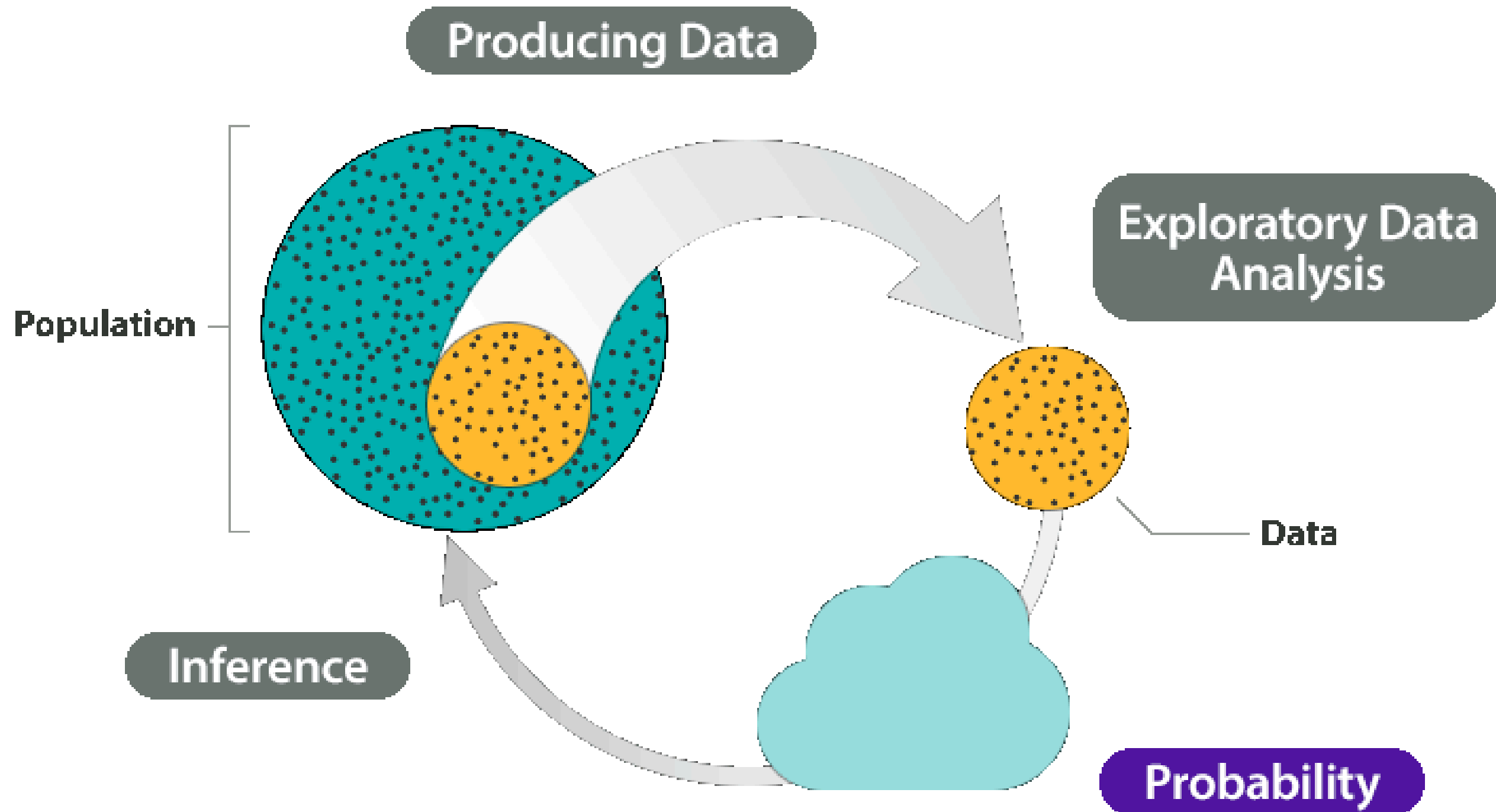


# WHY PROBABILITY?

Why Probability?

Calculate Probab

Normal Distribution



# WHY PROBABILITY?

**Statement** “There is a 20% probability of rain today.”

**Interpretation 1**      **It will rain for 20% of the day 4.8 hours.**

**Interpretation 2**      **There is a 1 in 5 chance that it will rain.**

**Interpretation 3**      **We can be 20% confident that it will rain today.**



Why Probability?

Calculate Probab

Normal Distribution

# PROBABILITY CALCULATION

$$P(E) = n(E) / n(S)$$

**Sample space** of an experiment, denoted by **S**, is the set of all possible outcomes of that experiment.

**Example:** - Sample space for single coin toss = {Heads, Tails}  
- Sample space for single die roll = {1, 2, 3, 4, 5, 6}



**Event A** is any collection (subset) of outcomes contained in the sample space **S**.  
That is, if **A** is an event then  $A = \{ \omega : \omega \in S \}$ .

Why Probability?

Calculate Probab

Normal Distribution

# RANDOM

Why Probability?

Calculate Probab

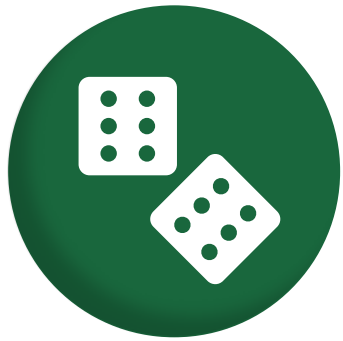
Normal Distribution



# PROBABILITY DISTRIBUTIONS



# PROBABILITY DISTRIBUTIONS



In this section we'll cover modeling data with **probability distributions**, and use the normal distribution to calculate probabilities and make estimates about normal populations

## TOPICS WE'LL COVER:

Distribution Basics

Distribution Tyes

Normal Distribution

Z-Scores

Probabilities

Values Estimates

## GOALS FOR THIS SECTION:

- *Understand the concept of a probability distribution, and its relationship with frequency distributions*
- *Learn about the different types of probability distributions, and their main differences*
- *Identify the properties of the normal distribution*
- *Calculate probabilities, values, and z-scores from normal distributions using Excel functions*

# PROBABILITY DISTRIBUTIONS

## Distribution Basics

## Distribution Types

## Normal Distribution

## Z-Scores

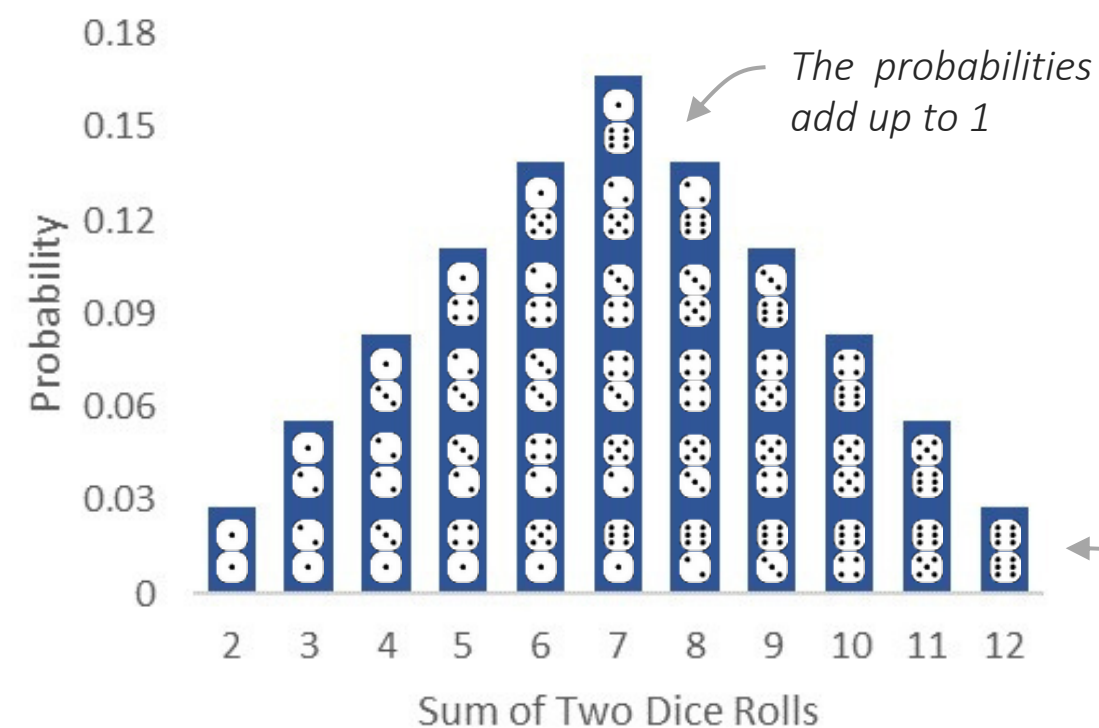
## Probabilities

A **probability distribution** represents a variable's idealized frequency distribution. It shows all the possible values a variable can take, and their chances of occurring.

- Frequencies in a sample are based on the underlying probabilities of those values occurring.

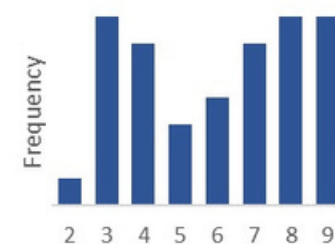
**EXAMPLE** *Results of rolling two dice*

**PROBABILITY DISTRIBUTION** (Population):

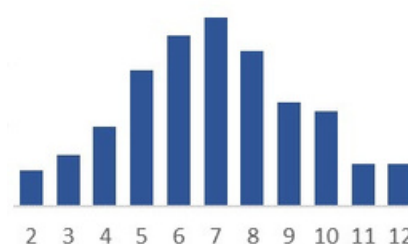


**FREQUENCY DISTRIBUTION** (Sample):

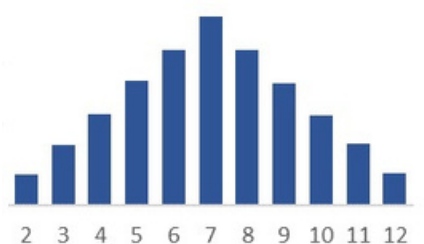
Sample size = 50



$n = 500$



$n = 50,000$



In an infinite sample, a variable's relative frequency distribution is equal to its probability distribution!

This is known as a **binomial distribution**, and it can be used to calculate probabilities on the outcome of rolling two dice (without rolling them fifty thousand times!)

# TYPES OF PROBABILITY DISTRIBUTIONS

There are two **types of probability distributions**: Discrete & Continuous

Distribution Basics

Distribution Tyes

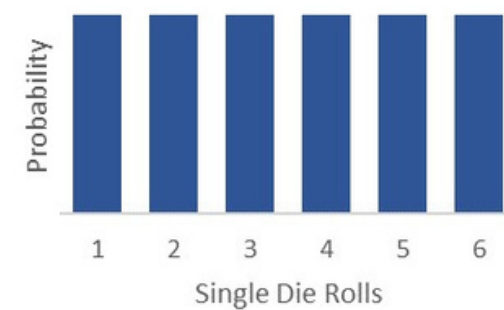
Normal Distribution

Z-Scores

Probabilities

## 1) Discrete probability distributions

*Uniform*



*Binomial*



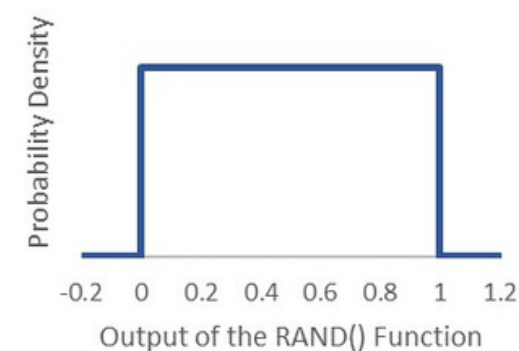
*Poisson*



*The height of each bar is its probability. There are "gaps" between the numbers.*

## 2) Continuous probability distributions

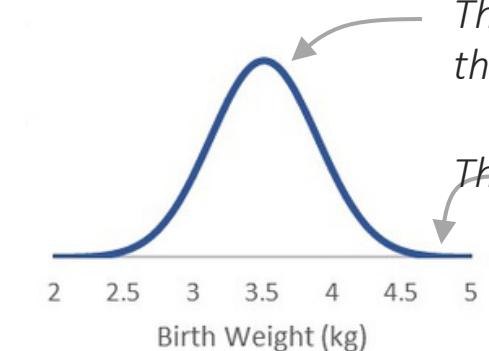
*Uniform*



*Exponential*



*Normal*



*The height of the curve is NOT its probability, the **area under the curve** is (more on this later!)*

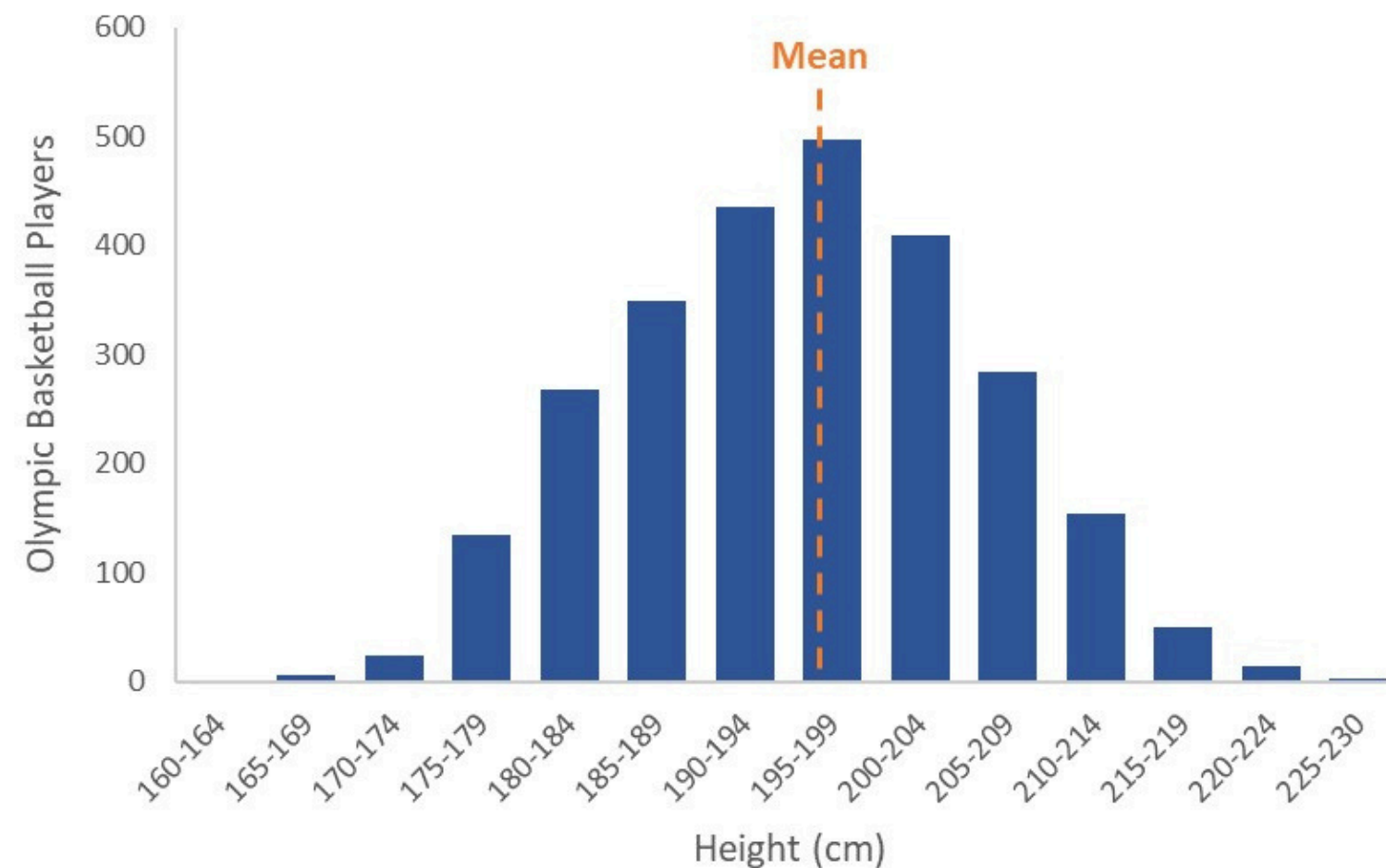
*The numbers can take any value*

# THE NORMAL DISTRIBUTION

Many numerical variables naturally follow a **normal distribution**, or “bell curve”

- Normal distributions are symmetrical around the mean and have no skew (*mean = median*), with most data concentrated around its center and flaring out in “tails” on both ends

**EXAMPLE** | *Olympic Basketball Player Heights*



## HEY THIS IS IMPORTANT!

Since they are so common, many statistical tests are designed for normally distributed populations, which is why we'll mostly focus on the normal distribution in the course

# THE NORMAL DISTRIBUTION

The normal distribution is described by two values: the **mean & standard deviation**

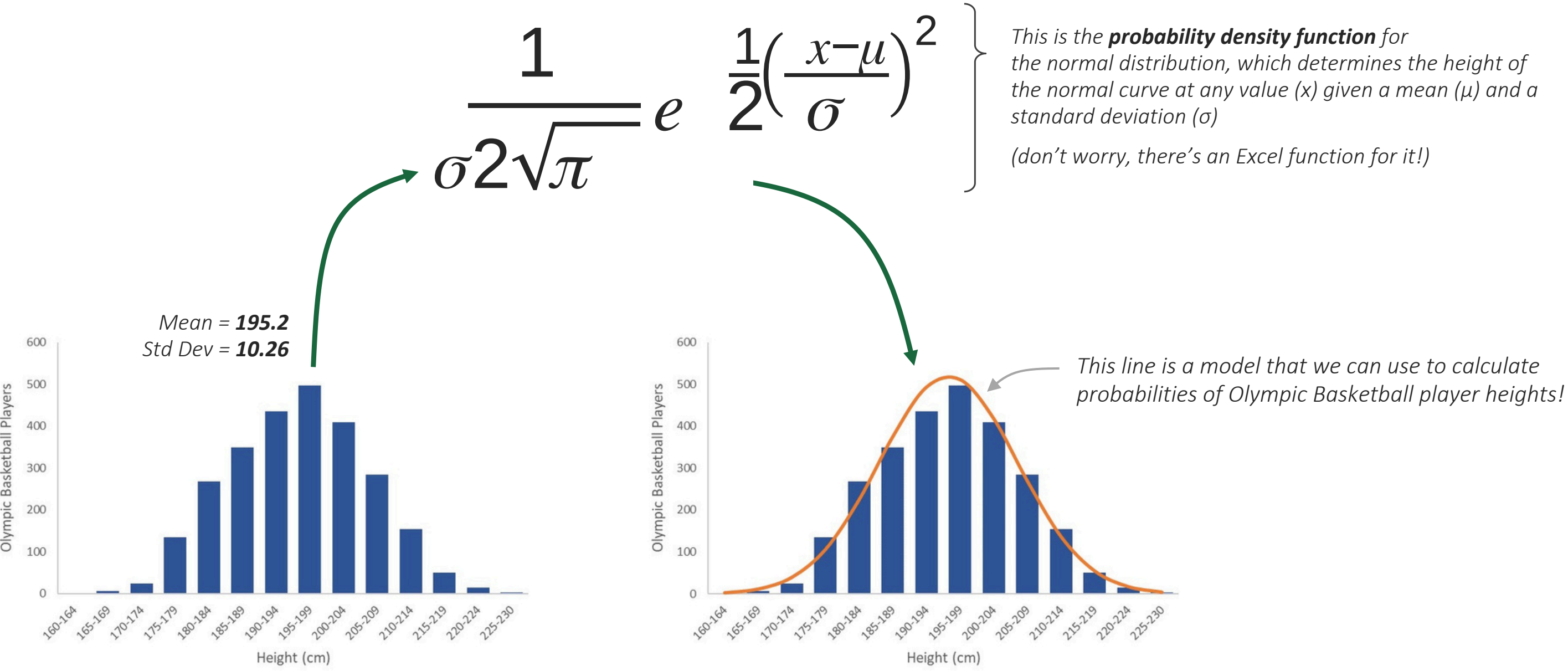
Distribution Basics

Distribution Tyes

Normal Distribution

Z-Scores

Probabilities

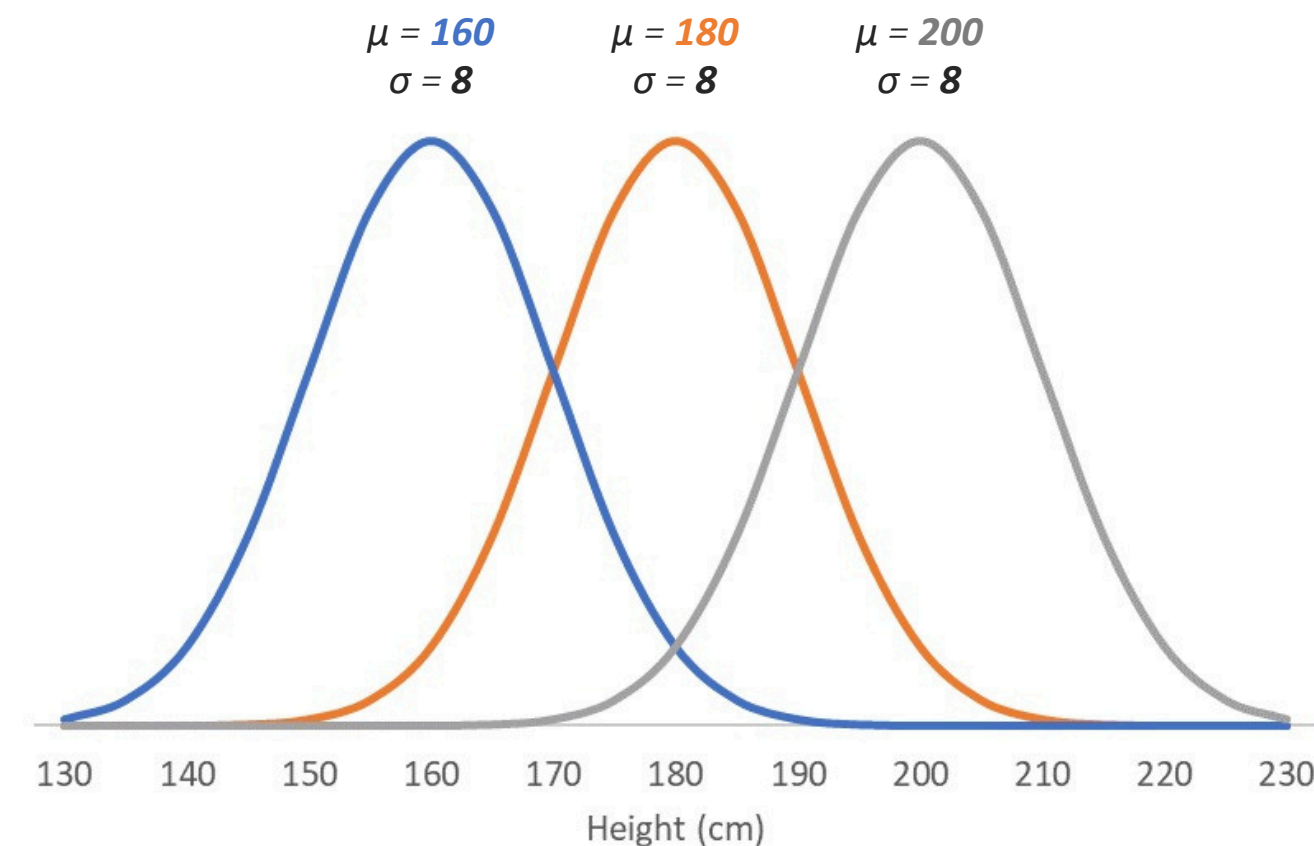




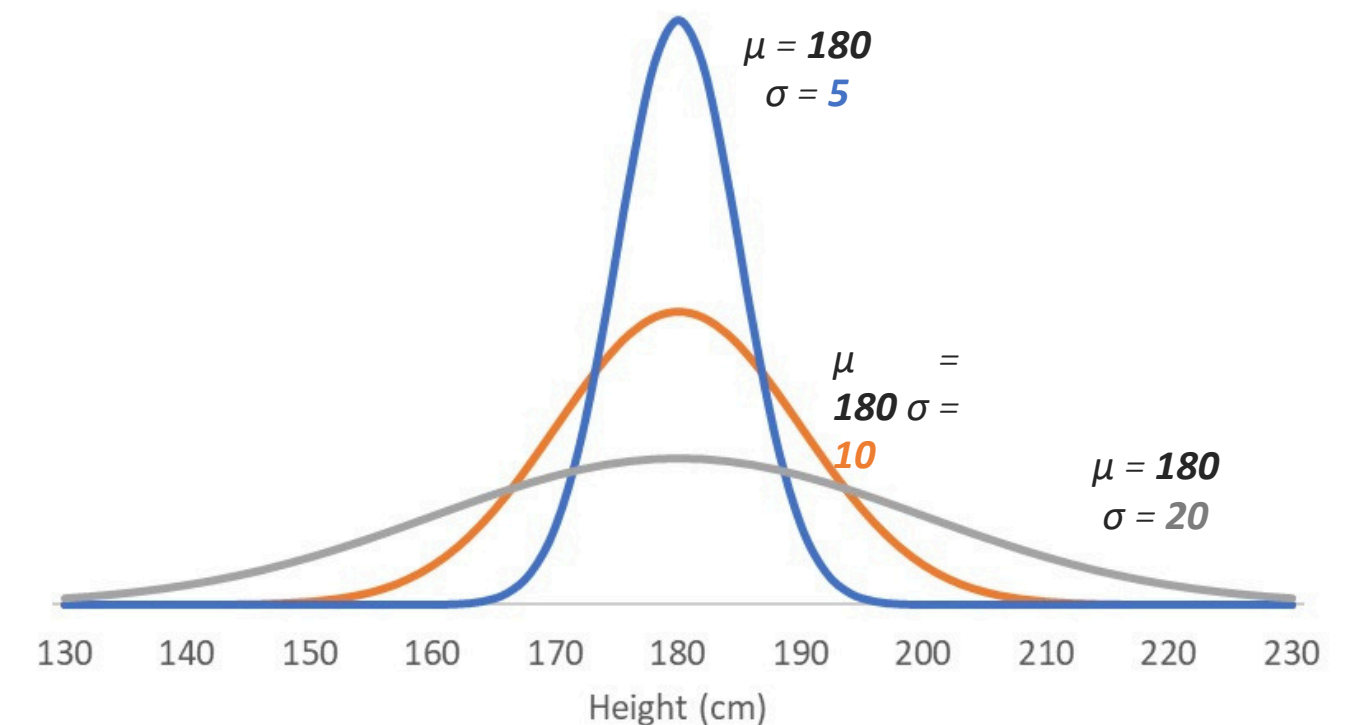
# THE NORMAL DISTRIBUTION

The normal distribution is described by two values: the **mean & standard deviation**

- The mean determines the *center* of the distribution, and the standard deviation its *width*



Changing the mean **shifts** the curve along the x axis



Changing the standard deviation **squeezes** or **stretches** the curve

Distribution Basics

Distribution Types

Normal Distribution

Z-Scores

Probabilities

# Z-SCORES

A **z-score** indicates how many standard deviations away from the mean a value lies

Distribution Basics

Distribution Tyes

Normal Distribution

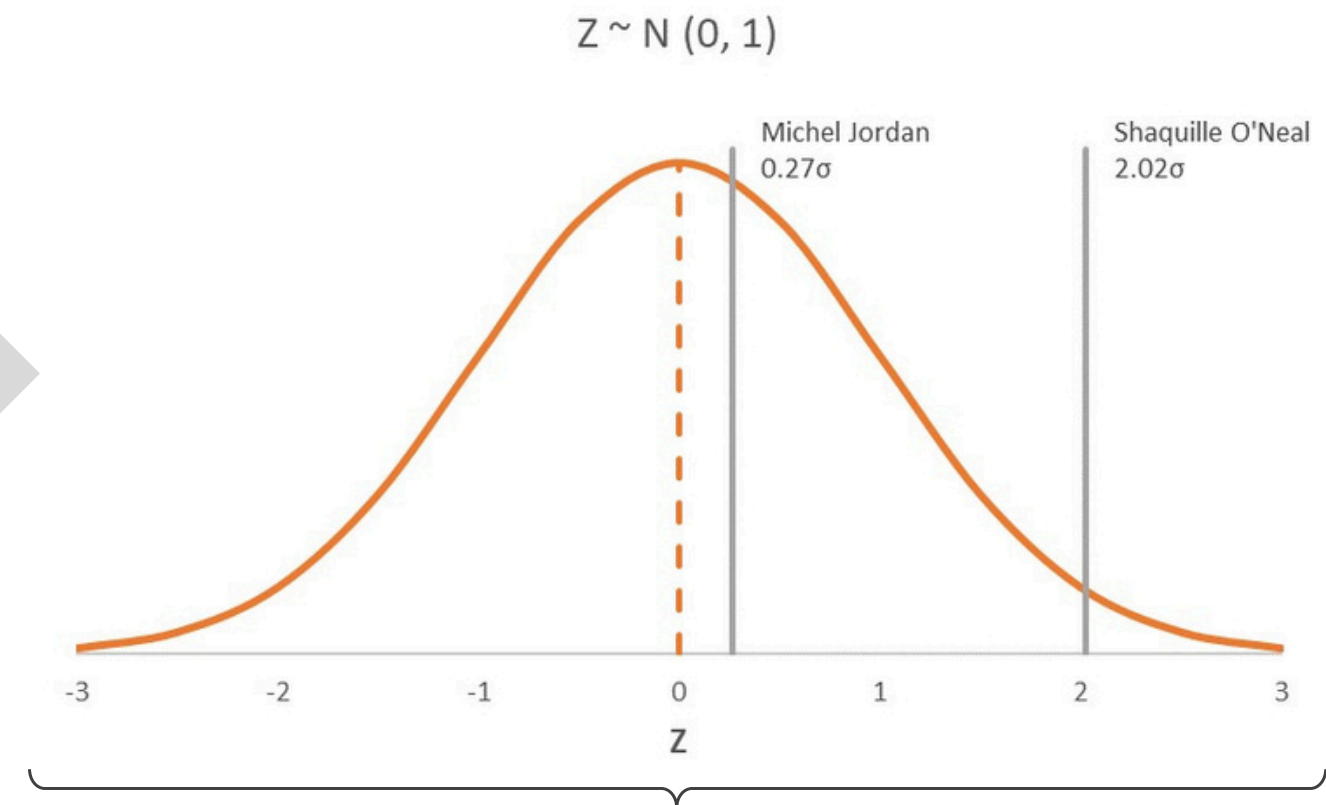
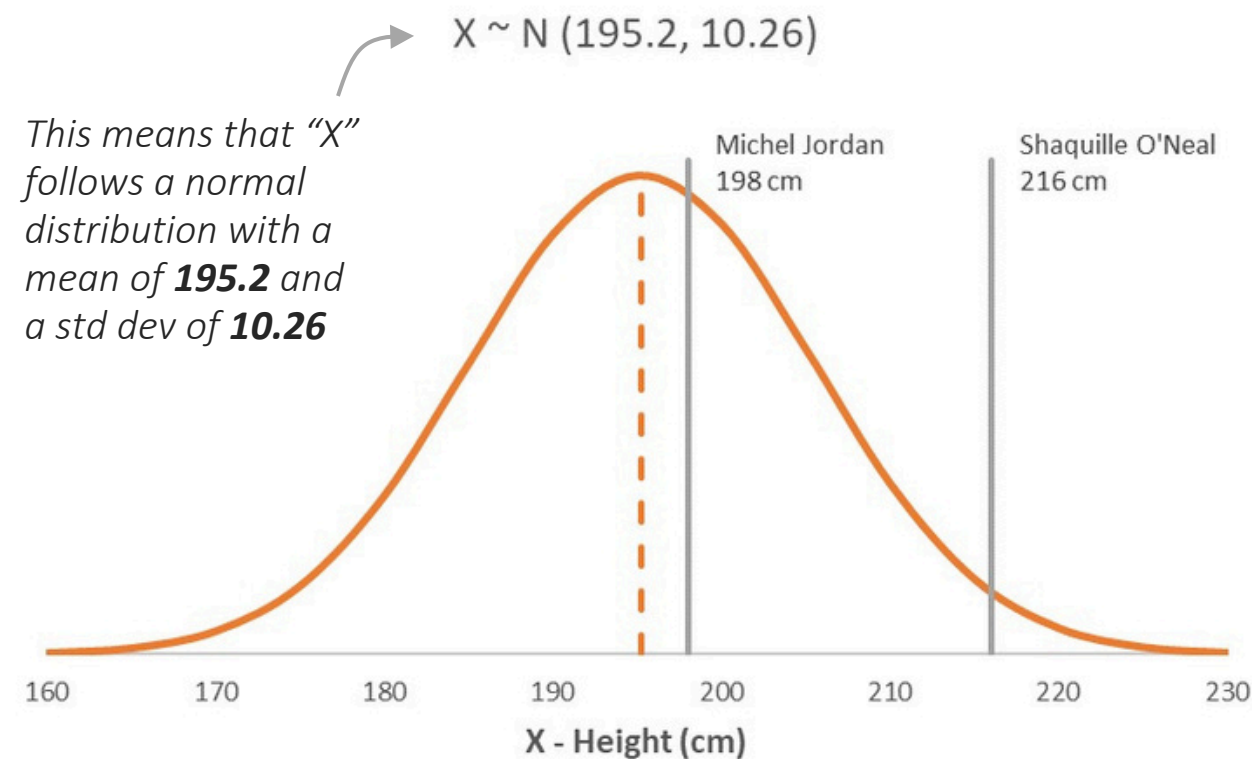
Z-Scores

Probabilities

$$z = \frac{x - \mu}{\sigma}$$

To calculate a z-score for a value, simply subtract the mean and divide by the standard deviation (or use the STANDARDIZE function)

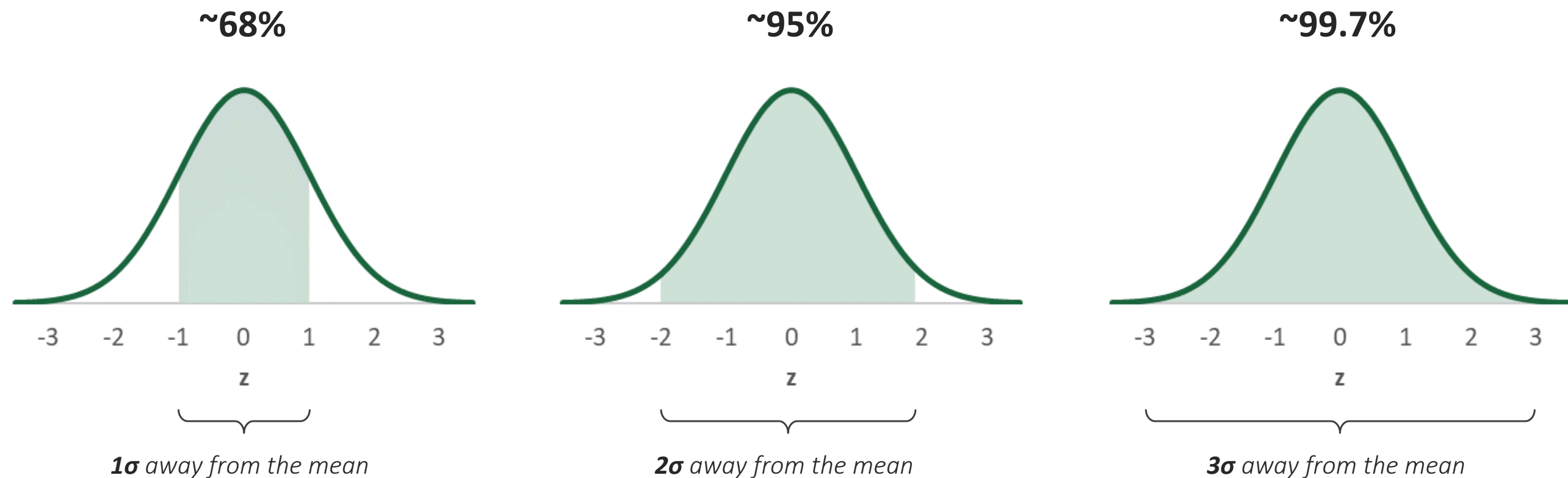
$$z = \frac{198 - 195.2}{10.26} = 0.27$$



This is known as the **standard normal distribution**, or z-distribution, and has a mean of 0 and a standard deviation of 1

# THE EMPIRICAL RULE

The **empirical rule** outlines where most values fall in a normal distribution



**PRO TIP:** Beyond using a histogram to determine whether your data is distributed normally, check if it follows the empirical rule



# EXCEL NORMAL DISTRIBUTION FUNCTIONS

These **Excel functions** help make calculations related to the normal distribution:

Distribution Basics

Distribution Tyes

Normal Distribution

Z-Scores

Probabilities

**NORM.DIST()**

Returns the cumulative probability or the probability density at an  $x$  value from a given normal distribution

=**NORM.DIST**( $x$ ,  $\mu$ ,  $\sigma$ , cumulative)

**NORM.INV()**

Returns the  $x$  value in a given normal distribution at a specified cumulative probability

=**NORM.INV**(probability,  $\mu$ ,  $\sigma$ )

**STANDARDIZE()**

Returns the  $z$ -score for a specified  $x$  value in a given normal distribution

=**STANDARDIZE**( $x$ ,  $\mu$ ,  $\sigma$ )

**NORM.S.DIST()**

Returns the cumulative probability or the probability density at a  $z$ -score from the standard normal distribution

=**NORM.S.DIST**( $z$ , cumulative)

**NORM.S.INV()**

Returns the  $z$ -score in the standard normal distribution at a specified cumulative probability

=**NORM.S.INV**(probability)

# CALCULATING PROBABILITIES

Distribution Basics

Distribution Types

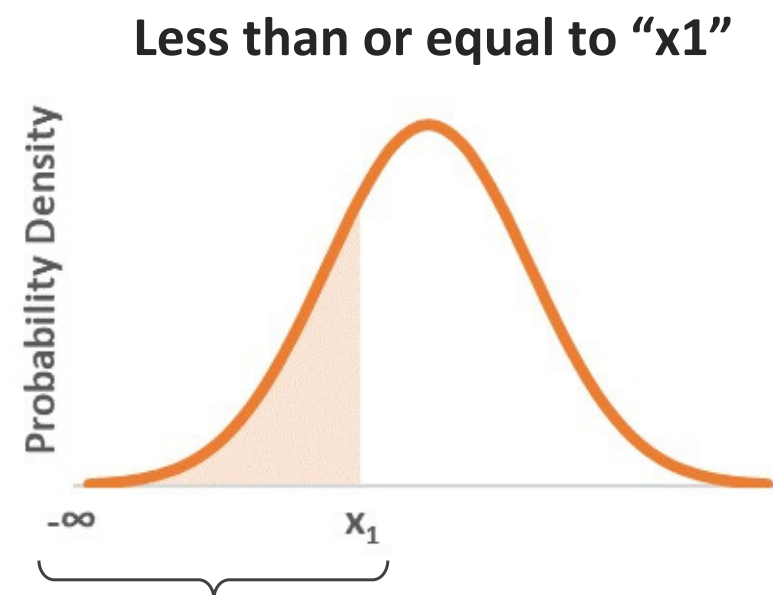
Normal Distribution

Z-Scores

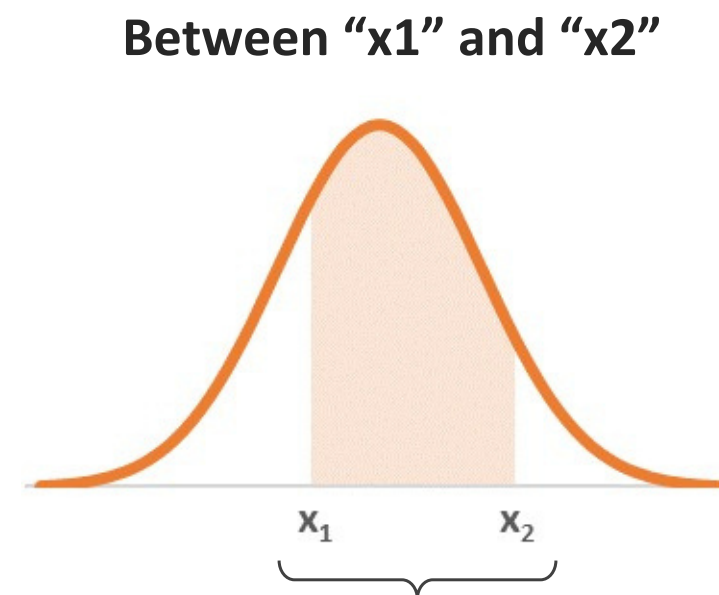
Probabilities

If a variable follows a normal distribution, you can **calculate the probability** of randomly obtaining a value within a specified range

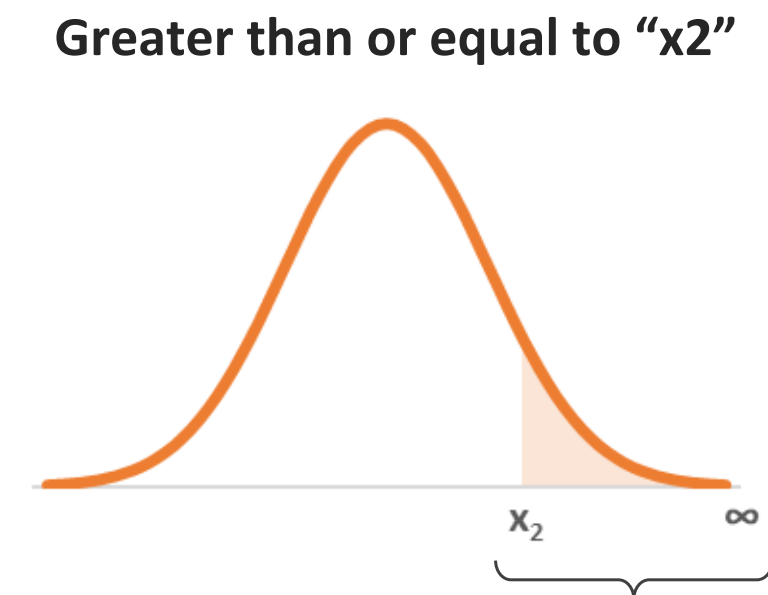
- This is determined by the area under the curve in that range



The area from negative infinity to “ $x_1$ ” is the **cumulative probability**



This is the cumulative probability of “ $x_2$ ” minus the cumulative probability of “ $x_1$ ”



This is 1 (the entire area under the curve) minus the cumulative probability of “ $x_2$ ”



**HEY THIS IS IMPORTANT!**

You CANNOT calculate the probability of obtaining an  $x$  value *exactly* – there’s no area under a single point!

# THE NORM.DIST FUNCTION

## NORM.DIST()

Returns the cumulative probability or the probability density at “x” from a normal distribution

=**NORM.DIST**(x, mean, standard\_dev, cumulative)

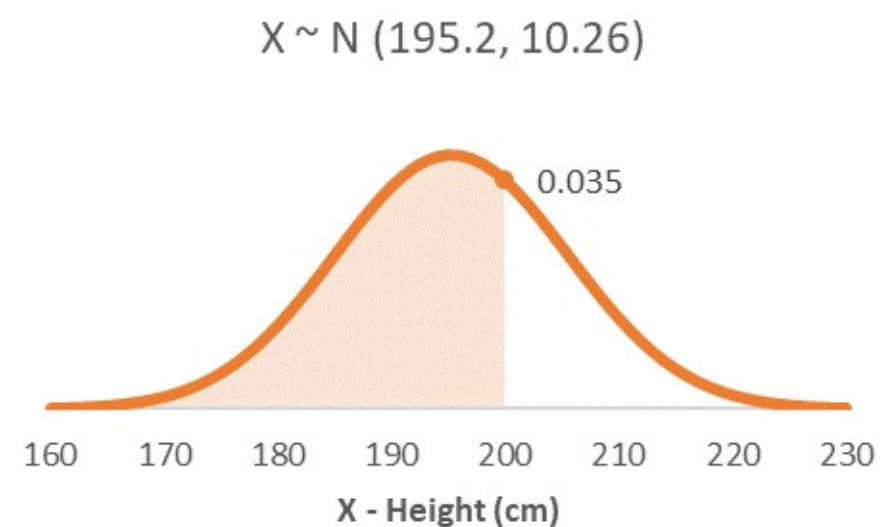
The **value** to calculate the probability for

The **mean & standard deviation** for the normal distribution of the population

**TRUE:** The area under the curve  
**FALSE:** The height of the curve

Possible question:

“What’s the probability of an Olympic Basketball Player being **2 meters tall or shorter**?”



=**NORM.DIST**(200, 195.2, 10.26, TRUE) = **0.68**

This is the probability!

=**NORM.DIST**(200, 195.2, 10.26, FALSE) = **0.035**

This is just the height of the curve

# THE NORM.DIST FUNCTION

## NORM.DIST()

Returns the cumulative probability or the probability density at “x” from a normal distribution

=**NORM.DIST**(x, mean, standard\_dev, cumulative)

The **value** to calculate the probability for

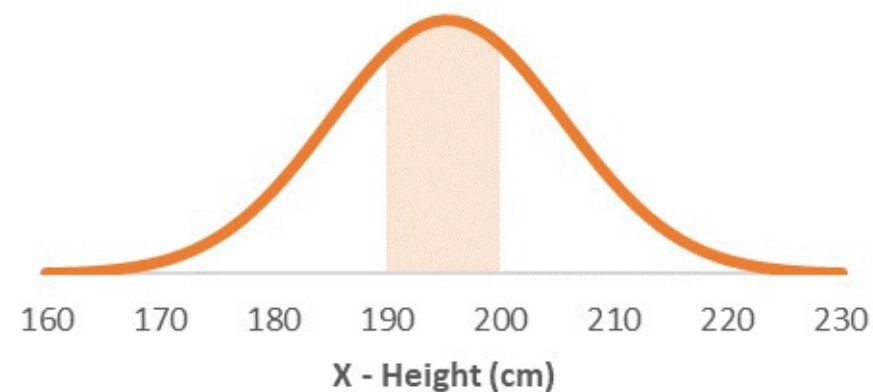
The **mean & standard deviation** for the normal distribution of the population

**TRUE:** The area under the curve  
**FALSE:** The height of the curve

Possible question:

“What’s the probability of an Olympic Basketball Player being **between 1.9 and 2 meters tall?**”

$X \sim N(195.2, 10.26)$



=**NORM.DIST**(200, 195.2, 10.26, TRUE) = **0.68**

=**NORM.DIST**(190, 195.2, 10.26, TRUE) = **0.3061**

=0.68-0.306 = **0.3739**

← This is the probability!

# THE NORM.DIST FUNCTION

**NORM.DIST()**

Returns the cumulative probability or the probability density at “x” from a normal distribution

**=NORM.DIST**(x, mean, standard\_dev, cumulative)

The **value** to calculate the probability for

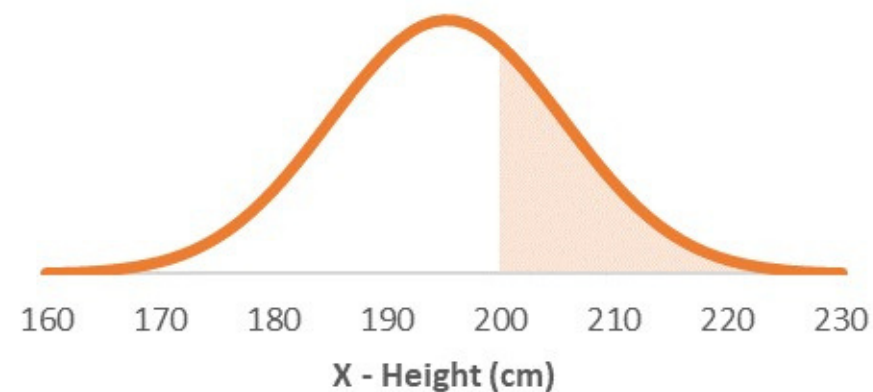
The **mean & standard deviation** for the normal distribution of the population

**TRUE:** The area under the curve  
**FALSE:** The height of the curve

Possible question:

“What’s the probability of an Olympic Basketball Player being **at least 2 meters tall**?”

$X \sim N(195.2, 10.26)$



**=NORM.DIST(200, 195.2, 10.26, TRUE) = 0.68**

**=1-NORM.DIST(190, 195.2, 10.26, TRUE) = 0.32**

The cumulative probability under the entire curve is equal to 1  
(it’s every value possible!)

This is the probability!



# THE NORM.S.DIST FUNCTION

## NORM.S.DIST()

Returns the cumulative probability or the probability density at “z” from the z-distribution

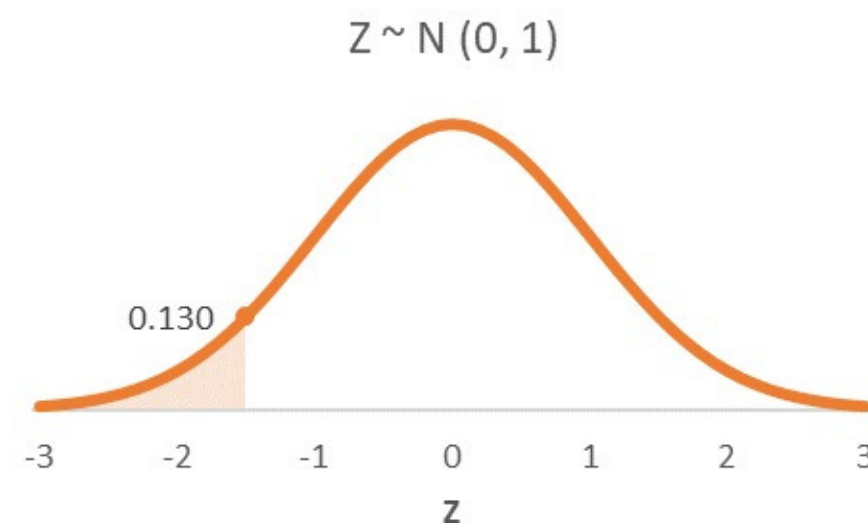
=**NORM.S.DIST**(z, cumulative)

The **z-score** to calculate the probability for

**TRUE:** The area under the curve  
**FALSE:** The height of the curve

Possible question:

“What’s the probability of an Olympic Basketball Player being **at least 1.5 standard deviations shorter than the mean?**”



=**NORM.S.DIST**(-1.5, TRUE) = **0.066** *This is the probability!*

=**NORM.S.DIST**(-1.5, FALSE) =  
**0.130** *This is just the height of the curve*

# KEY TAKEAWAYS: PROBABILITY DISTRIBUTIONS

---

★ A probability distribution is an **idealized frequency distribution**

- *It shows all the possible values the variable can take, and the probability of each value occurring*

★ Many variables naturally follow a **normal distribution**

- *The data is symmetrical around its mean, and flares out in “tails” (the width depends on the standard deviation)*

★ The probability in a normal distribution is the **area under its curve**

- *It can only be calculated in intervals, not for exact values!*

★ There are **Excel functions** to solve normal probability problems

- *NORM.DIST and NORM.S.DIST let you calculate the probability of randomly obtaining values in specified ranges*
- *NORM.INV and NORM.S.INV let you estimate values or z-scores based on their cumulative probabilities*