



---

# **Module 6 Lecture - Multiple Comparisons for Kruskal-Wallis**

## Analysis of Variance

---

Quinton Quagliano, M.S., C.S.P

Department of Educational Psychology

## Table of Contents

<b>1</b>	<b>Overview and Introduction</b>	<b>2</b>
1.1	Textbook Learning Objectives . . . . .	2
1.2	Instructor Learning Objectives . . . . .	2
1.3	Introduction . . . . .	2
<b>2</b>	<b>Introducing the Normal Curve</b>	<b>3</b>
2.1	Notation . . . . .	3
2.2	Probability Density Function . . . . .	3
2.3	Basic Characteristics . . . . .	4
<b>3</b>	<b>Calculating AUC for the Normal Distribution</b>	<b>5</b>
3.1	Introduction . . . . .	5
3.2	Calculating . . . . .	6
<b>4</b>	<b>The Standard Normal Distribution</b>	<b>6</b>
4.1	Introduction and Z-scores . . . . .	6
4.2	Notation . . . . .	8
<b>5</b>	<b>Brief Example of Data ‘Coercion’</b>	<b>8</b>
5.1	Introduction . . . . .	8
<b>6</b>	<b>Conclusion</b>	<b>10</b>
6.1	Recap . . . . .	10

# 1 Overview and Introduction

## 1.1 Textbook Learning Objectives

- Recognize the normal probability distribution and apply it appropriately.
- Recognize the standard normal probability distribution and apply it appropriately.
- Compare normal probabilities by converting to the standard normal distribution.

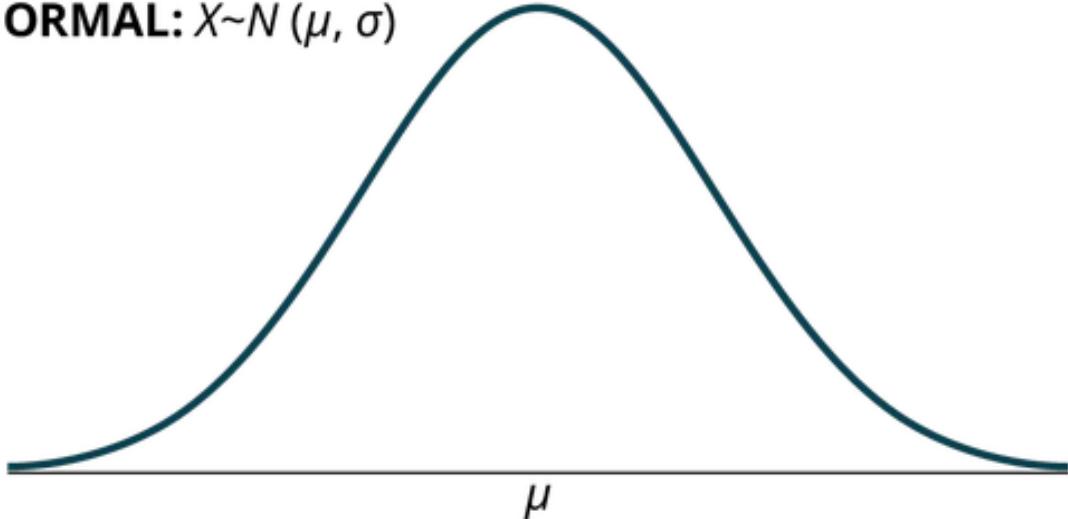
## 1.2 Instructor Learning Objectives

- Understand the normal distribution within the broader context of “ideal” distributions for continuous variables
- Be able to calculate z-scores, and understand the role of the mean and standard deviation in the case of the normal distribution

## 1.3 Introduction

- The \_\_\_\_\_ distribution is arguably the single most prevalent way to describe continuous variables.
  - It is often called the \_\_\_\_\_ curve due to its curved, symmetrical shape when plotted

**NORMAL:  $X \sim N(\mu, \sigma)$**



- Like the \_\_\_\_\_ or the \_\_\_\_\_ distributions, the normal distribution is an “ideal”
  - “Real” data will \_\_\_\_\_ ever be perfectly normal
  - However, much like with the other ideal distributions the normal curve serves a purpose for \_\_\_\_\_

"All models are wrong, but some are useful" - George Box

- Unlike the prior distributions that have been discussed, the normal distribution is used often in common \_\_\_\_\_ statistics

## 2 Introducing the Normal Curve

### 2.1 Notation

- The normal curve's notation is similar to others:  $X \sim N(\mu, \sigma)$  where:
  - $X$  is the \_\_\_\_\_ variable
  - $N$  is the designation of the \_\_\_\_\_ curve
  - $\mu$  is the population \_\_\_\_\_ parameter
  - $\sigma$  is the population standard \_\_\_\_\_ parameter
- Thus, if we assume that a normal distribution has a mean of 20 and a standard deviation of 2, this would be written as:  $X \sim N(20, 2)$

 Discuss: Now you try it: write the notation for a normal distribution with a mean of 12 and standard deviation of 3.

 Discuss: Review: Try writing notation for an exponential distribution with a decay parameter of 0.01 and a separate notation for uniform distribution with minimum value 2 and maximum value 20.

### 2.2 Probability Density Function

- The probability \_\_\_\_\_ function (pdf) for the normal curve is:

---

*"I don't mind not knowing. It doesn't scare me." — Richard P. Feynman*

$$f(x) = \frac{1}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot e^{-0.50 \cdot (\frac{x-\mu}{\sigma})^2}$$

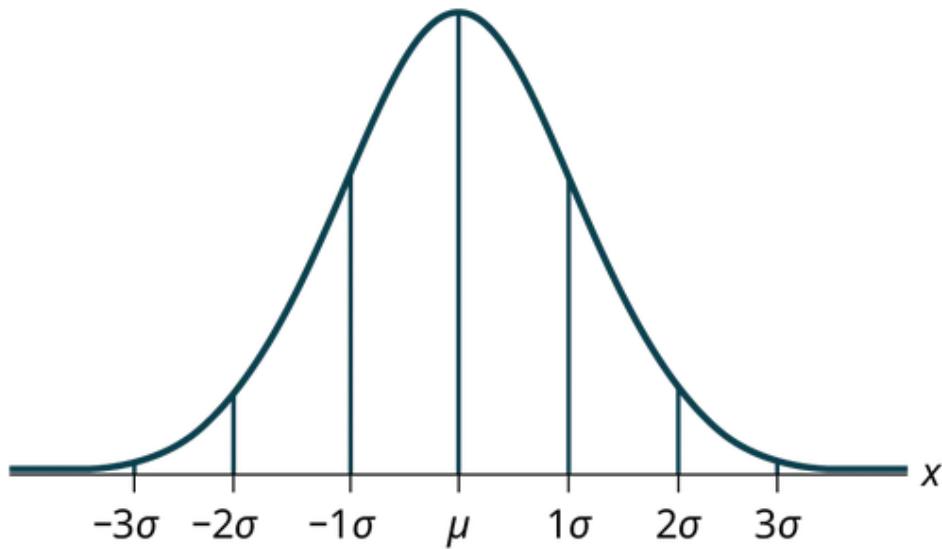
- Compare this to the relatively \_\_\_\_\_ pdf for a uniform distribution:

$$f(x) = \frac{1}{b - a}$$

- Due to the complexity of this equation, we often don't calculate \_\_\_\_\_ under the curve (AUC) for normal distributions by hand
  - Historically, one can use \_\_\_\_\_ that approximate the AUC, but in modern practice, this is handled fully by computers
  - We'll introduce the concept of calculating probability by hand in [Calculating AUC for the Normal Distribution](#), but later will use SPSS to handle this for us

## 2.3 Basic Characteristics

- The normal curve can \_\_\_\_\_ in appearance quite a bit, as:
  - Changes in \_\_\_\_\_ move the curve to the left and right
  - Changes in standard deviation can make it more \_\_\_\_\_ or \_\_\_\_\_
- However, it is:
  - Always \_\_\_\_\_
  - The mean, median, and mode of the distribution are all the \_\_\_\_\_
  - Following the \_\_\_\_\_ rule
- The **empirical rule**, also known as 68-95-99.7 rule, states that:
  - 68% of  $x$  values lie between  $-1\sigma$  and  $1\sigma$  or  $z = -1$  to  $1$
  - 95% of  $x$  values lie between  $-2\sigma$  and  $2\sigma$  or  $z = -2$  to  $2$
  - 99.7% of  $x$  values lie between  $-3\sigma$  and  $3\sigma$  or  $z = -3$  to  $3$



### 3 Calculating AUC for the Normal Distribution

#### 3.1 Introduction

- As mentioned prior, the \_\_\_\_\_ of the normal distribution is often calculated with computers due to complexity
  - This is often the case with the \_\_\_\_\_ and exponential distributions as well - simply too time consuming and complex to calculate for large datasets or complex scenarios

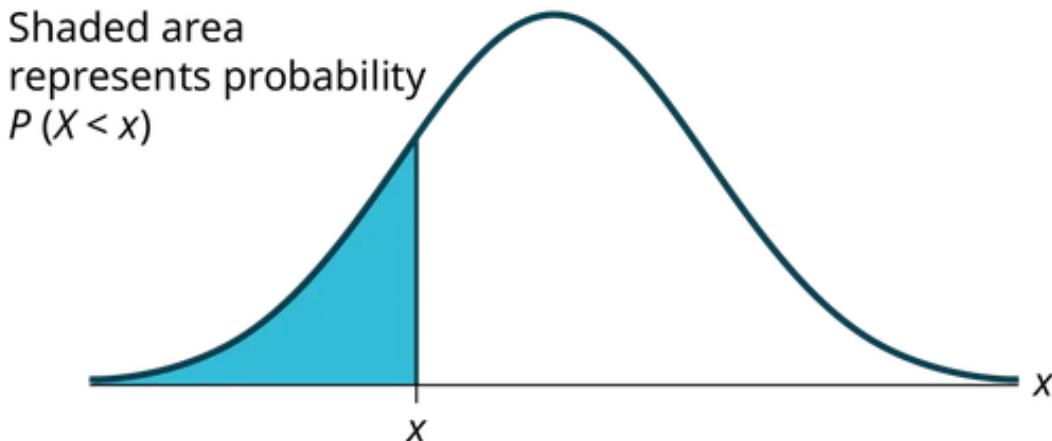
##### ! Important

Its worth mentioning that the field of statistics has rapidly grown in possible methods with advancements in technology and computing

Discuss: Review: Under what circumstances can we use a binomial distribution?

### 3.2 Calculating

- For the following applications  $X$  and  $x$  mean the same thing they have in previous modules



- To find area to the \_\_\_\_\_ of a specified point  $x$ , we find  $P(X < x)$ 
  - Thus, to find the complement or inverse or area to the \_\_\_\_\_ of specified point  $x$ , we will do  $P(X > x) = 1 - \underline{P(X < x)}$
  - Remember that individual points of  $x$  don't have \_\_\_\_\_ in a continuous distribution, so this is functionally the same as saying  $P(X \leq x)$  for area to the left
- Your book shows how to perform these calculations using a calculator function, but we will use SPSS in the next practical assignments

## 4 The Standard Normal Distribution

### 4.1 Introduction and Z-scores

- There is a special case of the normal distribution with more clearly specified characteristics: the \_\_\_\_\_ **normal distribution**
- The standard/standardized normal distribution is made up of \_\_\_\_\_, instead of whatever "raw" continuous values would be used

► Discuss: Review from descriptive statistics before the next part, what is a z-score? Explain in your own words

- Review: Z-scores can be practically \_\_\_\_\_ as, “number of standard deviations a point is from the mean of the data”
  - Thus, a data point with a  $z = +2.00$  is 2 standard deviations \_\_\_\_\_ the mean and a data point with  $z = -1.20$  is 1.2 standard deviation \_\_\_\_\_ the mean
  - A data point directly at the mean of the data will have  $z = 0.00$ .
  - Z-scores are best understood as \_\_\_\_\_ indicators of position within the dataset

! Important

If you ever hear that data was ‘standardized’ or ‘mean-centered’, that means that each data point was transformed into its respective z-score.

- Review: the \_\_\_\_\_ for z-scores in a sample is  $z = \frac{x-\bar{x}}{s}$ 
  - Application example in data of 1, 2, 3, 4, 5
  - $\bar{x} = 3$
  - $s = 1.581$
  - To get z-score of 4:  $z = \frac{4-3}{1.581} = 0.633$
- To determine what value a particular z-score is (in a sample) you can use:  $(z * s) + \bar{x}$ 
  - Taking from example above to find what point that  $z = 1.5$ :
  - $(1.5 * 1.581) + 3 = 5.37$

! Important

Remember there is a subtle difference in notation and formulas for sample statistics vs population parameters!

► Discuss: Rewrite this above problem and recompute using the population parameters instead of sample statistics

- One notable benefit of z-scores is that they allow us to compare variables on the same \_\_\_\_\_
  - However, a drawback along this same line is that interpreting z-scores in practical analysis is more \_\_\_\_\_ from a realistic interpretation.

## 4.2 Notation

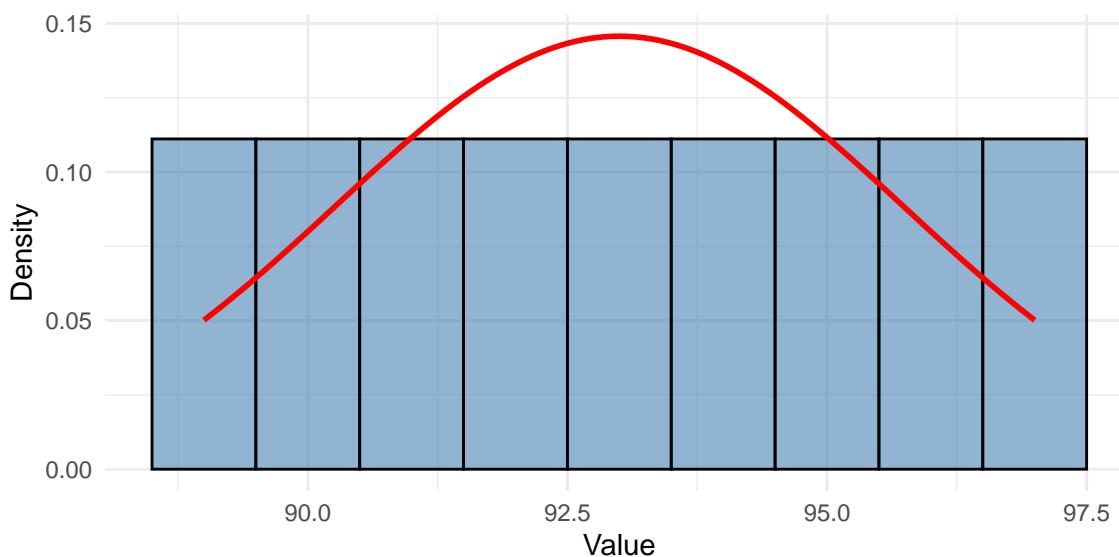
- Based on the \_\_\_\_\_ of z-scores, a standard normal curve always takes the notation  $Z \sim N(0, 1)$ 
  - This is because a set of z-scores from a normally-distributed variable will always have  $\mu = 0$  and a  $\sigma = 1$

## 5 Brief Example of Data ‘Coercion’

### 5.1 Introduction

- One \_\_\_\_\_ is that one can readily treat any continuous data as normal - this is not wise
- Take for example class grade percentage data: {93, 92, 91, 90, 89, 94, 95, 96, 97}
  - Treated “as usual”, this would mean lots of As and a B+
  - But, what if we treat this data as if it comes from a normal distribution?

Figure 1: Class Grade Percentage Histogram



🔊 Discuss: Instead of the normal distribution, what other continuous distribution does this plot remind you of?

- If I were to “grade on a curve”, it means the person who got an 89 would be graded as if they failed, because they were lower \_\_\_\_\_ to the other data points
  - This is \_\_\_\_\_ only if the assumption that the scores of all students is perfectly normal, but feels unfair if that assumption isn’t met
  - On the other hand, this system \_\_\_\_\_ help if all students did poorly, but some did just marginally better
- We’ll revisit the idea of improper application of \_\_\_\_\_ when we cover specific inferential tests

## 6 Conclusion

### 6.1 Recap

- The normal distribution is a useful and commonly used continuous variable distribution that meets several specific conditions. These characteristics make it readily predictable and applicable
- Much like the other distributions and density functions, we can use the characteristics of the normal distribution to calculate AUC and understand the relative spread and placement of the data. This is aided by use of z-score and the standard normal curve as a special case
- *However, it is often mis-used and misunderstood, and we must take caution as we continue in the semester*

*The instructor-provided glossary may not include all terms worth memorizing, make sure you consider using the vocabulary list in your book and your own judgment to make sure you have all relevant terms*