# Applying the Normal Model

Ryan Miller

▶ The last presentation introduced the **Normal distribution** as a probability model for a continuous random variable

**X**

## Introduction

- ▶ The last presentation introduced the **Normal distribution** as a probability model for a continuous random variable
- ▶ This model is defined by two components:
  - ▶ The parameter $\mu$, a constant that defines the *center* of the bell-curve
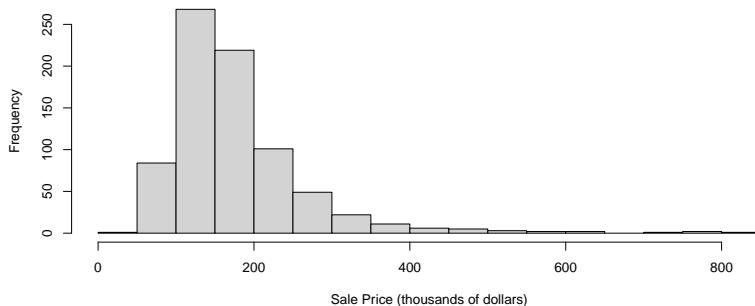  - ▶ The parameter $\sigma$, a constant that defines the *spread* of the bell-curve

# Introduction

- The last presentation introduced the **Normal distribution** as a probability model for a continuous random variable
- This model is defined by two components:
  - The parameter $\mu$, a constant that defines the *center* of the bell-curve
  - The parameter $\sigma$, a constant that defines the *spread* of the bell-curve
- This presentation will go through an example that illustrates where this model is and is not appropriate

# Example

- In this example, we'll look at the sale prices of all homes in Iowa City, IA between 2005-2008
  - The mean sale price was \$180.1k, and the standard deviation was \$90.65k

**Home Sales in Iowa City (2005–2008)**



Sale Price (thousands of dollars)

- Let $X$ be a random variable denoting the sale price of a randomly selected home
  - What might you consider using as a probability model for $X$?

- Let $X$ be a random variable denoting the sale price of a randomly selected home
  - What might you consider using as a probability model for $X$?
- Because $X$ is a continuous random variable, it seems reasonable to take the mean and standard deviation in our dataset and use $N(180.1, 90.65)$ as a probability model for $X$
  - How would you use this model to estimate $P(X \geq \$400k)$?

▶ Using StatKey, we could directly input our mean and standard deviation then calculate this right-tail probability to be 0.0076

## Example

▶ Using StatKey, we could directly input our mean and standard deviation then calculate this right-tail probability to be 0.0076
  ▶ We also could standardize $400k into a Z-score of $z = 400 - 180.190.65 = 2.426$ and use the Standard Normal distribution to arrive at the same estimated probability
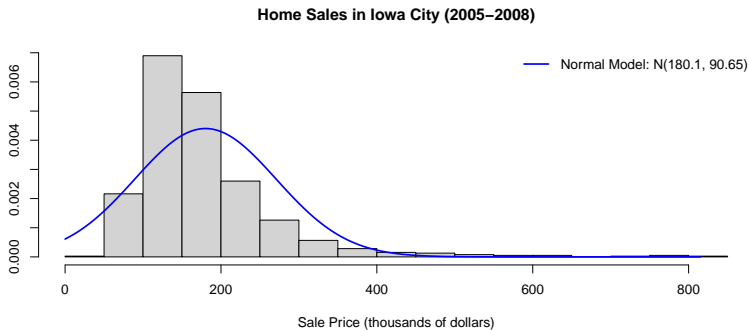
**X**

## Example

▶ Using StatKey, we could directly input our mean and standard deviation then calculate this right-tail probability to be 0.0076
  ▶ We also could standardize $400k into a Z-score of $z = 400 - 180.190.65 = 2.426$ and use the Standard Normal distribution to arrive at the same estimated probability
▶ However, both of these calculations assume the Normal model is a perfect representation of these data (or the population represented by them)
  ▶ Is that an appropriate assumption?

## Example

- The *empirical probability* of a randomly selected home selling for more than $400k$ is 0.0283 (22 of 777 homes)
  - This discrepancy might not seem like much, but this is 3.7 times larger than what the Normal model suggested! (0.0076)

**X**

# Example

- The *empirical probability* of a randomly selected home selling for more than $400*k* is 0.0283 (22 of 777 homes)
  - This discrepancy might not seem like much, but this is 3.7 times larger than what the Normal model suggested! (0.0076)



**Home Sales in Iowa City (2005–2008)**

Normal Model: N(180.1, 90.65)

Sale Price (thousands of dollars)

▶ Primary issue with this application is that distribution of the data doesn't match the *shape* of the normal curve

▶ Primary issue with this application is that distribution of the data doesn't match the *shape* of the normal curve
  ▶ That is, even if we center and scale our normal model appropriately (ie: good choices of $\mu$ and $\sigma$), the model is incapable of representing the underlying distribution of these data

- Primary issue with this application is that distribution of the data doesn't match the *shape* of the normal curve
  - That is, even if we center and scale our normal model appropriately (ie: good choices of $\mu$ and $\sigma$), the model is incapable of representing the underlying distribution of these data
- As an aside, notice these data contain $n = 777$ cases
  - A common misconception is that larger amounts of data tend to be normally distributed

# Appropriateness of the Normal Model

- ▶ Primary issue with this application is that distribution of the data doesn't match the *shape* of the normal curve
  - ▶ That is, even if we center and scale our normal model appropriately (ie: good choices of $\mu$ and $\sigma$), the model is incapable of representing the underlying distribution of these data
- ▶ As an aside, notice these data contain $n = 777$ cases
  - ▶ A common misconception is that larger amounts of data tend to be normally distributed
- ▶ That said, more data (larger sample sizes) *does* impact the distributional shape of certain random variables
  - ▶ Next week, we will extend the normal model to the distribution of *sample averages*

▶ The Normal distribution provides a useful probability model for many, but not all, continuous random variables
  ▶ Proper application of the Normal model requires the specification the bell-curve's center, $\mu$, and it's spread, $\sigma$

**X**

- The Normal distribution provides a useful probability model for many, but not all, continuous random variables
  - Proper application of the Normal model requires the specification the bell-curve's center, $\mu$, and it's spread, $\sigma$
  - However, variables with skewed distributions cannot be appropriately modeled by the normal curve, even when using reasonable values of $\mu$ and $\sigma$

# Conclusion

- ▶ The Normal distribution provides a useful probability model for many, but not all, continuous random variables
  - ▶ Proper application of the Normal model requires the specification the bell-curve's center, $\mu$, and it's spread, $\sigma$
  - ▶ However, variables with skewed distributions cannot be appropriately modeled by the normal curve, even when using reasonable values of $\mu$ and $\sigma$
- ▶ In general, having more data does not make a random variable more normally distributed
  - ▶ However, if the random variable represents the *sample average* (rather than the data-points themselves), having more data *does* have an important impact
  - ▶ We will explore the distribution of sample averages next week