# Biostat 100A Lab #3

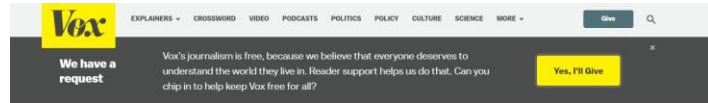## Section 1D: Wednesday 9:00 - 9:50 pm

## Power Transformation, Assessing Normality

```
LAB 2, C3: Apply Descriptive Techniques Commonly used to Summarize
Public Health Data
```

- Why do we care about having a GOOD write-up?
  - Not everyone knows how to INTERPRET statistics and especially in an age where media runs rampant and data is everywhere, it is EXTREMELY important that you can CORRECTLY INTERPRET DATA for others in a way that remains consistent with SCIENTIFIC TRUTH
  - Numbers mean NOTHING unless you assign some sort of meaning to them:
    - "Age is just a number"
      - In the context of **public health**, this is (usually) not true. Why?
        Example: During the COVID-19 pandemic, it was important to vaccinate high-risk populations, namely, **the elderly** because they are the most likely to become the most ill, have higher symptom severity, and more likely to be immunocompromised, compared to younger age groups

LAB 2, C3: Apply Descriptive Techniques Commonly used to Summarize Public Health Data

- When would you need to do a write-up?
  - Policy Briefs → People in Health Policy and Management are likely very familiar with this
  - Scientific Papers/Writing
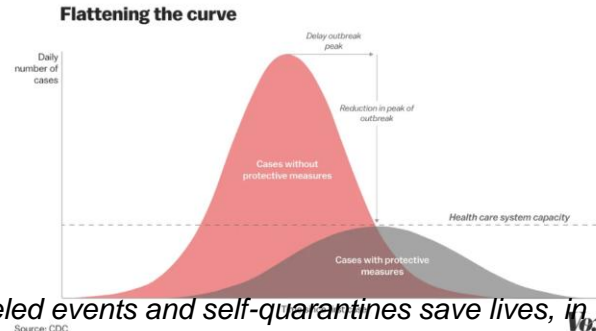  - If you are presenting in front of people who have no public health knowledge



Barclay, E., & Scott, D. (2020, March 10). *How canceled events and self-quarantines save lives, in one chart*. Vox. https://www.vox.com/2020/3/10/21171481/coronavirus-us-cases-quarantine-cancellation

LAB 2, C3: Apply Descriptive Techniques Commonly used to Summarize
Public Health Data

- ● What is YOUR ROLE?
  - ○ You act as a bridge between the scientific and non-scientific community
  - ○ You are a data **STORYTELLER**
    - ■ State what you **found/evidence**
    - ■ State the **Significance of your study**
    - ■ If you **can, explain why the data is the way it is (or, data generating process/the WHY)**
    - ■ Any **novel or surprising/not surprising finds**

# Example of an Excellent Write-Up

In acyanotic congenital heart disease, blood has normal oxygen content but is pumped throughout the body abnormally (1). Cyanotic congenital heart disease leads to less oxygen-rich blood to be delivered to the body because oxygen-rich and -poor blood mix in the heart (1, 2). Thus, children with cyanotic congenital disease tend to have a rise in hemoglobin levels to compensate for the decreased oxygen content (3), while children with acyanotic disease would not have this. We see this in the data provided: acyanotic children had a lower median hemoglobin level (13.1) than cyanotic children (14.8) [see C1 #3-4 for explanation on why medians are compared]. This is bolstered by the shapes of the frequency distributions for acyanotic children, which is centered around 11-15 g/cc, compared to 12.4-19.4 g/cc for cyanotic children. Both distributions are skewed by outliers; for acyanotic children, there is one abnormally large outlier (20.5) leading to right-skew, and for cyanotic children, there is one small outlier (5.4) leading to left-skew. The spread, as indicated by interquartile range, is also smaller for acyanotic children (1.55) than cyanotic children (3.3); this increased variability can potentially explained by differences in how low blood oxygen content is in each cyanotic child due to different disease etiologies—the lower the blood level, the greater the compensation needed. In conclusion. cyanotic children tend to have larger and more variable hemoglobin levels than acyanotic children, reflecting poor blood oxygen content.
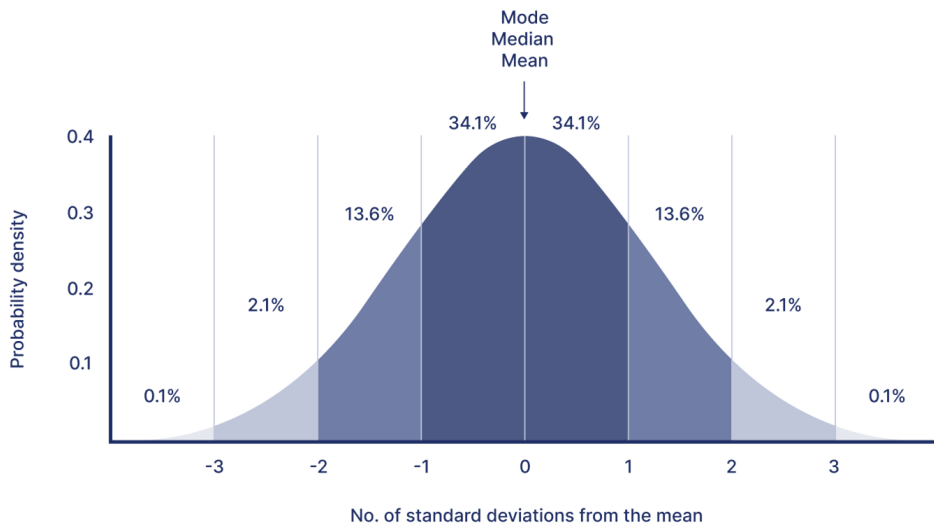
Green – Context, Claim/Hypothesis

Black – Supporting Evidence, how does the data support or undermine the claim/hypothesis?

Blue – Explanation behind the data, the WHY

Purple – Conclusion/Takeaway

# Normality



## Standard normal distribution

Usually, normally distributed data is the most ideal, It is well-understood and has many benefits:

- Easy to understand: mean as a central tendency and sd as a variance parameter
- Can be directly compared with other data
- Meet the assumptions of many statistical analysis.
- ….

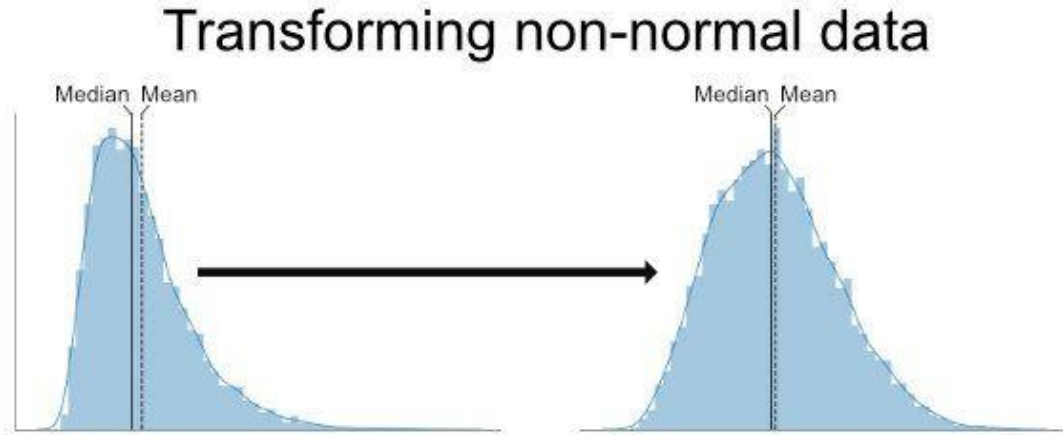**But in many cases, our data is far from a normal distribution…**

*Slides Courtesy of Qingyuan Liu*

# Power Transformation

Data transformation can transform non-normal data to approximately normal distribution. Common transformations include the **logarithm, square root, and square**

## Log-transformation Example:

- Log Transformation is usually used when the data is **positively skewed**, where data in the long tail is far away from the center
- Applying a logarithmic transformation compresses the range by **pulling in the extreme values more than the lower values,** which makes it more symmetric.

Transforming non-normal data

Median Mean          Median Mean

*Slides Courtesy of Qingyuan Liu*

# Caution: Zero Values in Power Transformation



- Many power transformation, such as log, is not applicable if the data contains 0's!
- It is important to check your original and transformed data to make sure there is no error produced

**Logarithm: Log(0) = UNDEFINED**
**Reciprocal: 1/0 = UNDEFINED**

*Slides Courtesy of Qingyuan Liu*

# Solutions for error messages in transformed data

- One way to omit the error: **IF()** combined with **ISNUMBER()** inside summary functions

  =SUM(IF(ISNUMBER(A1:A10), A1:A10))

  =PERCENTILE(IF(ISNUMBER(A1:A10), A1:A10), 0.25)

-In this case, IF() **evaluates whether the value is a valid number** for each value one by one, return True or False

-If True, this data is included, if false, it is not. The error message "#NUM!" is not a valid number, thus it will not be included in the summary functions
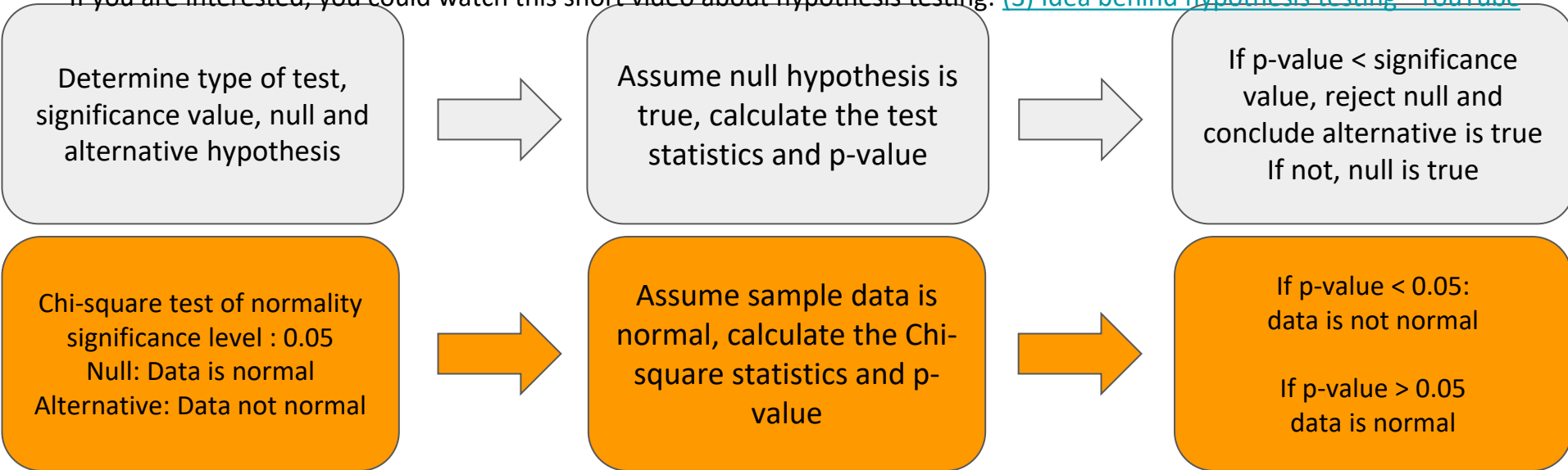
- For average(), you can also use **AVERAGEIFS()** instead:

  =AVERAGEIFS(A1:A10, "<>#DIV/0!")

*Slides Courtesy of Qingyuan Liu*

# Normality Test - Hypothesis Testing

The mechanism of normality test is **beyond the scope**, this diagram is a simple illustration of general hypothesis testings, you will learn more in the lectures. For this lab, you only need to know **how to interpret p-values.**

If you are interested, you could watch this short video about hypothesis testing: (3) Idea behind hypothesis testing - YouTube

| Determine type of test, significance value, null and alternative hypothesis | → | Assume null hypothesis is true, calculate the test statistics and p-value | → | If p-value < significance value, reject null and conclude alternative is true If not, null is true |

| Chi-square test of normality significance level : 0.05 Null: Data is normal Alternative: Data not normal | → | Assume sample data is normal, calculate the Chi-square statistics and p-value | → | If p-value < 0.05: data is not normal If p-value > 0.05 data is normal |

Note: p-value is not a golden standard, the result could be misleading

*Slides Courtesy of Qingyuan Liu*

# Venn Diagram

Each circle in a diagram indicates an event, and the area of their overlap means the intersection of events.

**Notations:**
A ∪ B: Union of A and B
A ∩ B: Intersection of A and B
$A^c$ or A' : Complement of A