

# Statistical Methodology for Software Engineers

Hadas Lapid, PhD

# Course Syllabus

- Statistical distributions
- Hypothesis testing
- Chisquare test for goodness of fit
- Extreme value distribution and its role in sequence analysis
- Multiple hypothesis testing
- Analysis of contingency tables
- Two-sample t-test
- Analysis of variance
- Non-parametric tests
- Correlation tests
- Linear regression

# Course Requirements

- 4 HW Assignments, 20%
- Exam, 80%

## Success guaranteed if

- Review class materials before each lecture
- Practice and freely explore R software

## Resources

- All course materials on Moodle
- “An Introduction to Statistical Methods and Data Analysis” by R. Lyman Ott and Michael Longnecker (on moodle)
- Introduction to Statistics and Data Analysis \_ With Exercises, Solutions and Applications in R - Heumann · Schomaker
- R Software: <https://rstudio.com/products/rstudio/download/#download>
- Web-textbook: <http://www.statsoft.com/Textbook>

# Contents Today

- **Variables**

- types
- Basic concepts

- **Distributions**

- Binomial Distribution
- Poisson Distribution
- Normal Distribution
- Standardized Normal Density function

- **Data example**

# Why do we need statistics?

## Typical Problems:

- Data variability
- Repeated measurements may yield different results

## Statistical Solutions:

- Provide predictions of outcome
- Draw **conclusions** about data variables (inverse problem)

## Related Subjects:

- [Probability](#) (direct problem)
- [Business analytics](#)
- [Machine learning](#) (predicting)
- [Operations research](#) (optimization)

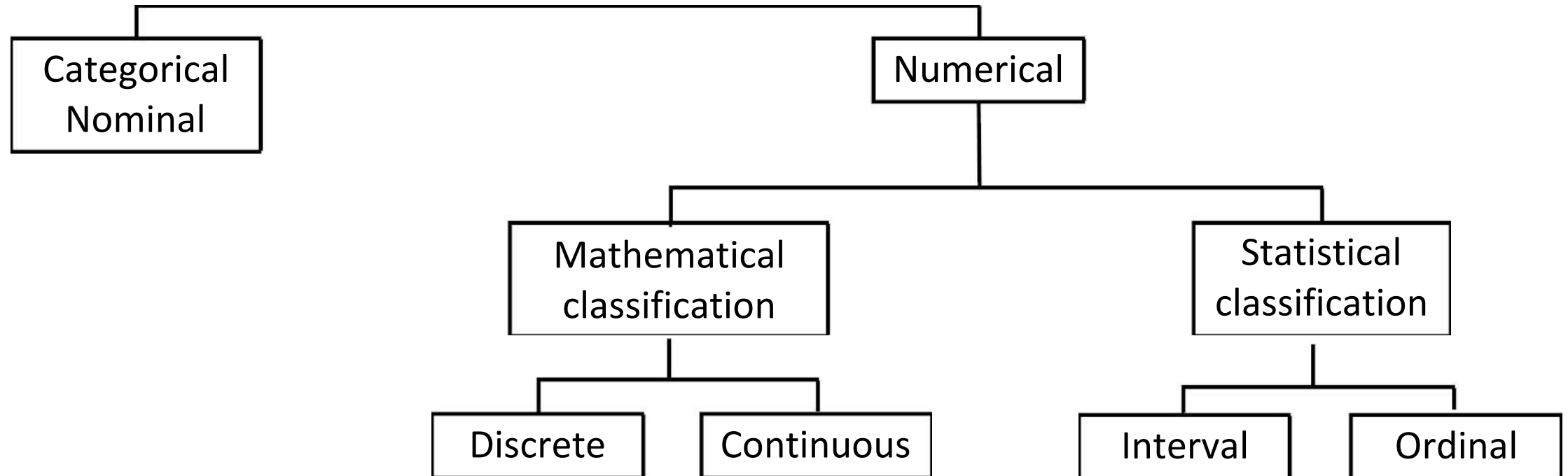
# What are variables?

- Objective/subjective measures that can be monitored, controlled or manipulated.
- Variables differ in the quantity of measurable information their scale provides

## Dependent vs. Independent variables

- Independent variables – inputs, explanatory, manipulated in the research (e.g. gender)
- Dependent variables – output, response only measured/monitored (e.g. white cell blood count)

# Variables Classification

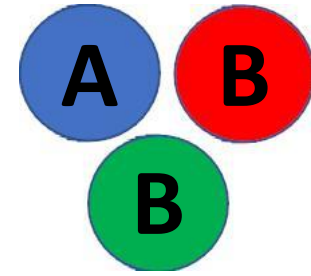


# Nominal (Categorical) Variables

Enable qualitative classification

i.e., whether an item belongs to a distinct category

e.g., city, gender, color

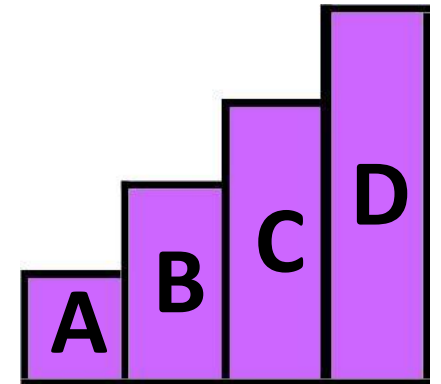


# Discrete Variables

Enable order ranking

i.e., whether an item has more or less of the measured quality

e.g., Olympic rank, grades rank, disease severity level

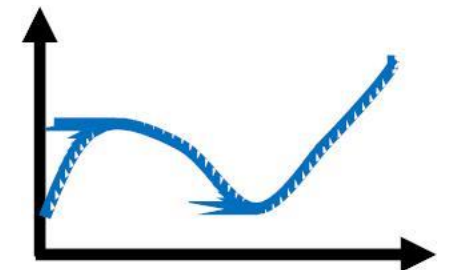


# Continuous Variables

Quantitative measurement

i.e., how much of the measured characteristic the item displays

e.g., height, weight, temperature





# Basic Concepts

Let  $(x_1, x_2, x_3 \dots x_n)$ , be  $n$  samples of  $X$  measurements in some population

$\mu \equiv$  Population mean (Expectation Value)

$\bar{x} = \langle x \rangle \equiv$  Sample mean

$$\langle x \rangle = \frac{\sum_{i=1}^n x_i}{n}$$

# Basic Concepts

**$\sigma^2$**   $\equiv$  population variance

**$s^2$**   $\equiv$  sample variance,  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

We can assume, on average:  **$S^2 \approx \sigma^2$**

**$\sigma$**   $\equiv$  population standard deviation

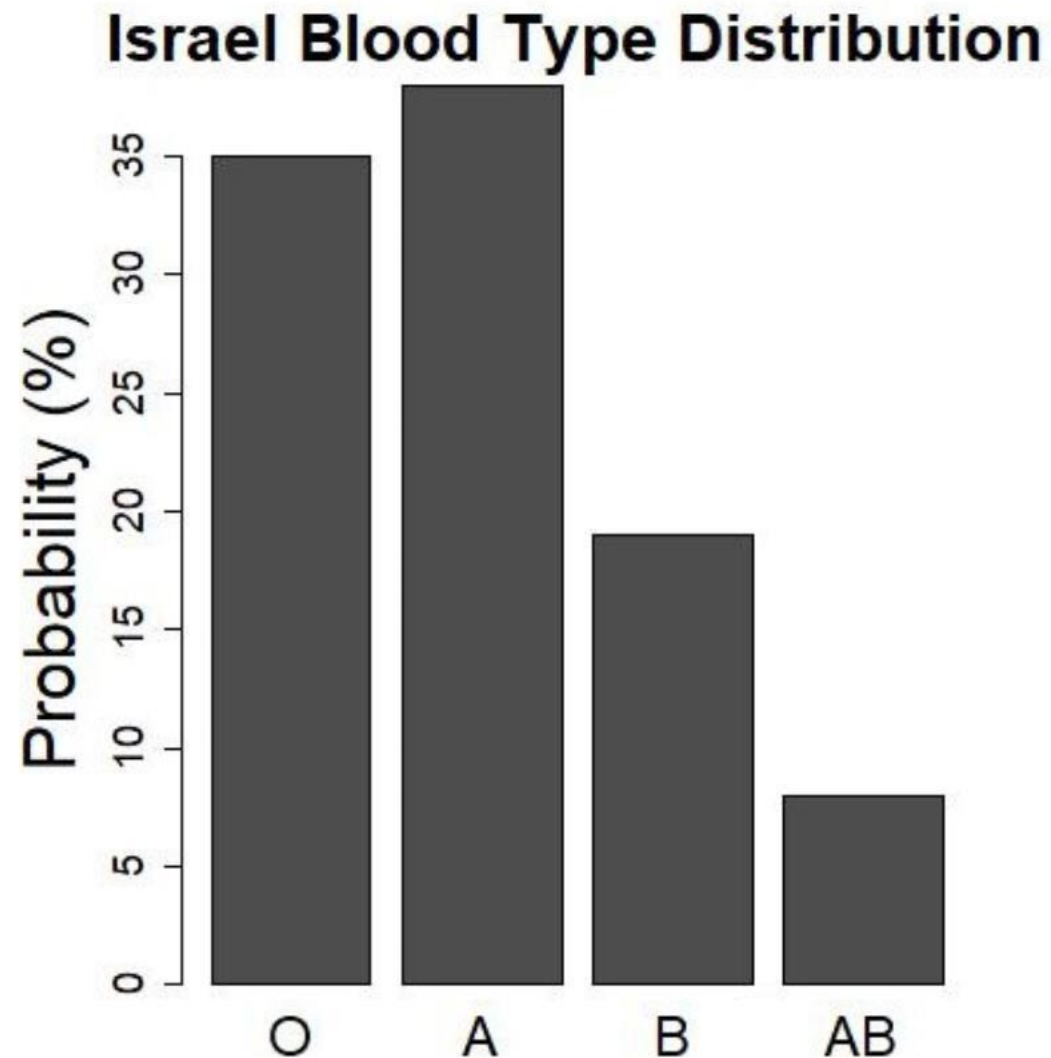
**$S$**   $\equiv$  sample standard deviation,  $S = \sqrt{S^2}$

**Stderr**  $\equiv$  sample standard error of the mean,

$$s = \frac{S}{\sqrt{n}} = \sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2}$$

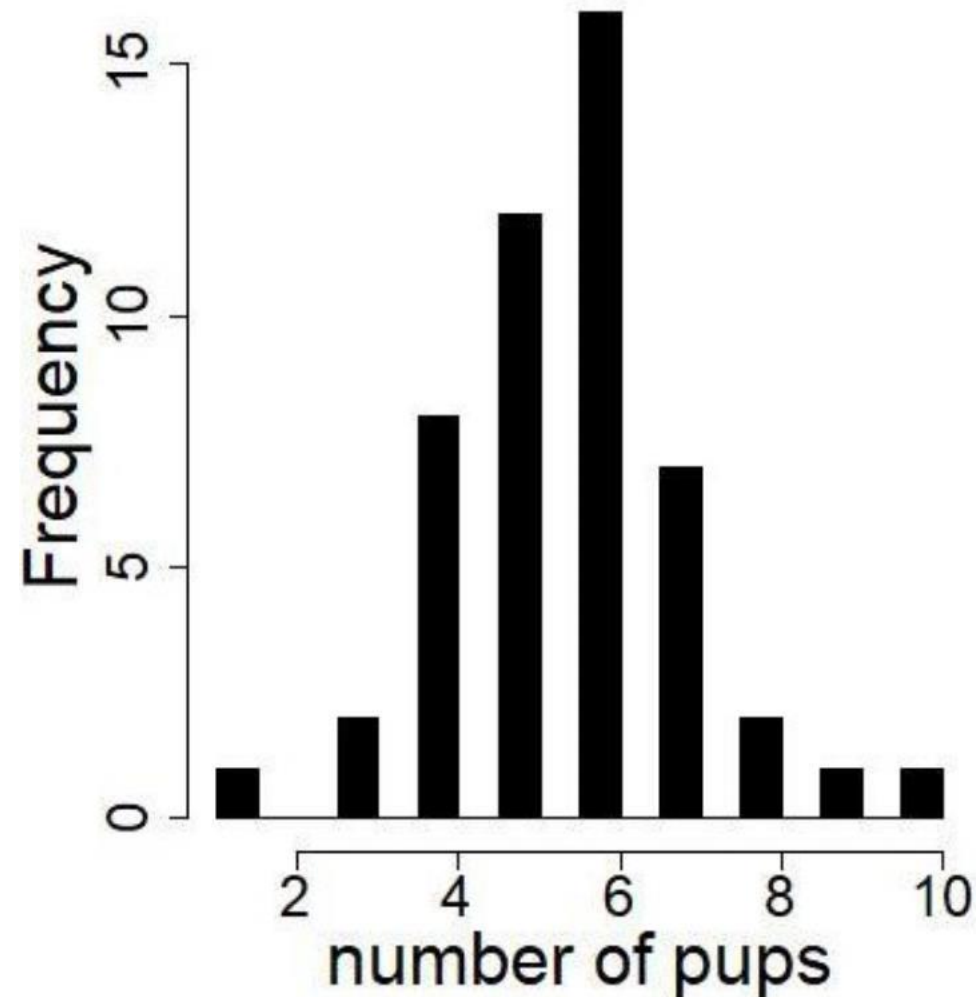
# Distributions

## Categorical Variable Distribution



# Discrete Variable Distribution

pups per litter (sample size = 50)



# Binomial Variables

- Binomial variable is a measured outcome of an experiment, which can be either **Success (True)** or **Failure (False)**.

$p$  = Probability of Success (roh)

$q = 1-p$  = Probability of Failure

# Binomial Distribution

- Binomial distribution results from  **$n$**  independent trials of a Bernoulli experiment.
- **$x$**   $\equiv$  The number of successes measured in an experiment.
- **$n$**   $\equiv$  The number of trials
- Notation:  **$x \sim B(n, \rho)$**

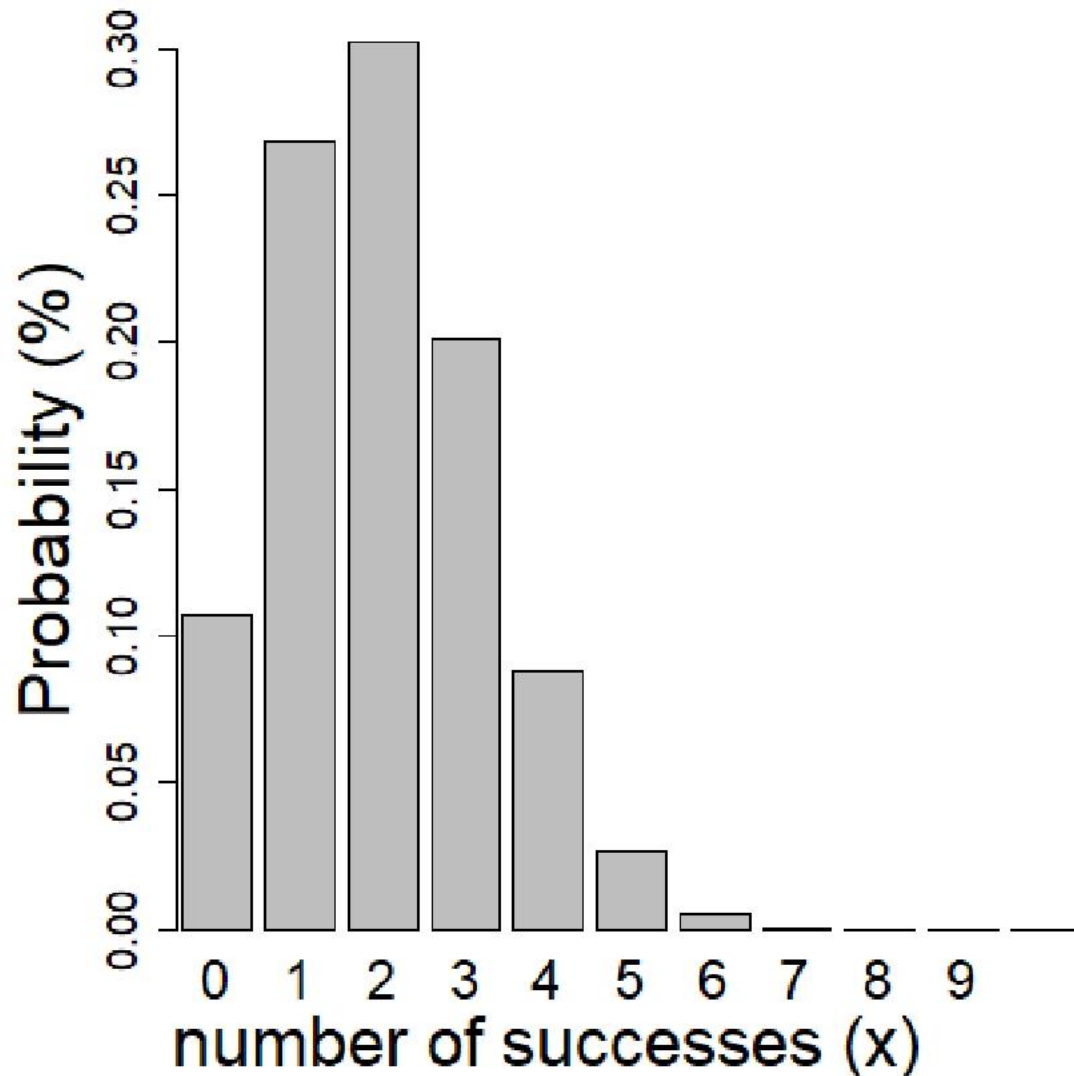
The **probability of  $x$  successes** out of  $n$  trials provided success rate  $\rho$ :

- **$P(X = x) = \binom{n}{x} \rho^x q^{n-x}$**

# Binomial Distribution - Summary

- $x \sim B(n, \rho)$
- $x \equiv$  The number of successes measured in an experiment.
- $n \equiv$  The number of trials
- $\rho \equiv$  probability of success
- $P(X = x) = \binom{n}{x} \rho^x q^{n-x} = \frac{n!}{x! \cdot (n-x)!} \rho^x q^{n-x}$
- $\mu = np$  Expectation Value (population mean)
- $\sigma^2 = npq$  Variance
- $\sigma = \sqrt{npq}$  Standard deviation

# Binomial Distribution - Example



$$X = \{x | 0 < x < 10\}$$

$$n = 10$$

$$p = 0.2$$

$$q = 0.8$$

- Calculate the probability for 3 successes out of 10  $P(X = 3)$
- Calculate the expectation value,  $\mu$
- Calculate the standard deviation,  $\sigma$



# Poisson Distribution

- **Discrete probability distribution** that describes the probability of a given number of events occurring in a fixed interval of space or time
- **Examples:**
  - # of calls to an emergency service per hour
  - # of traffic accidents in a given street per week
  - # of incoming text messages per day
  - # of laser photons hitting a detector in a particular time interval

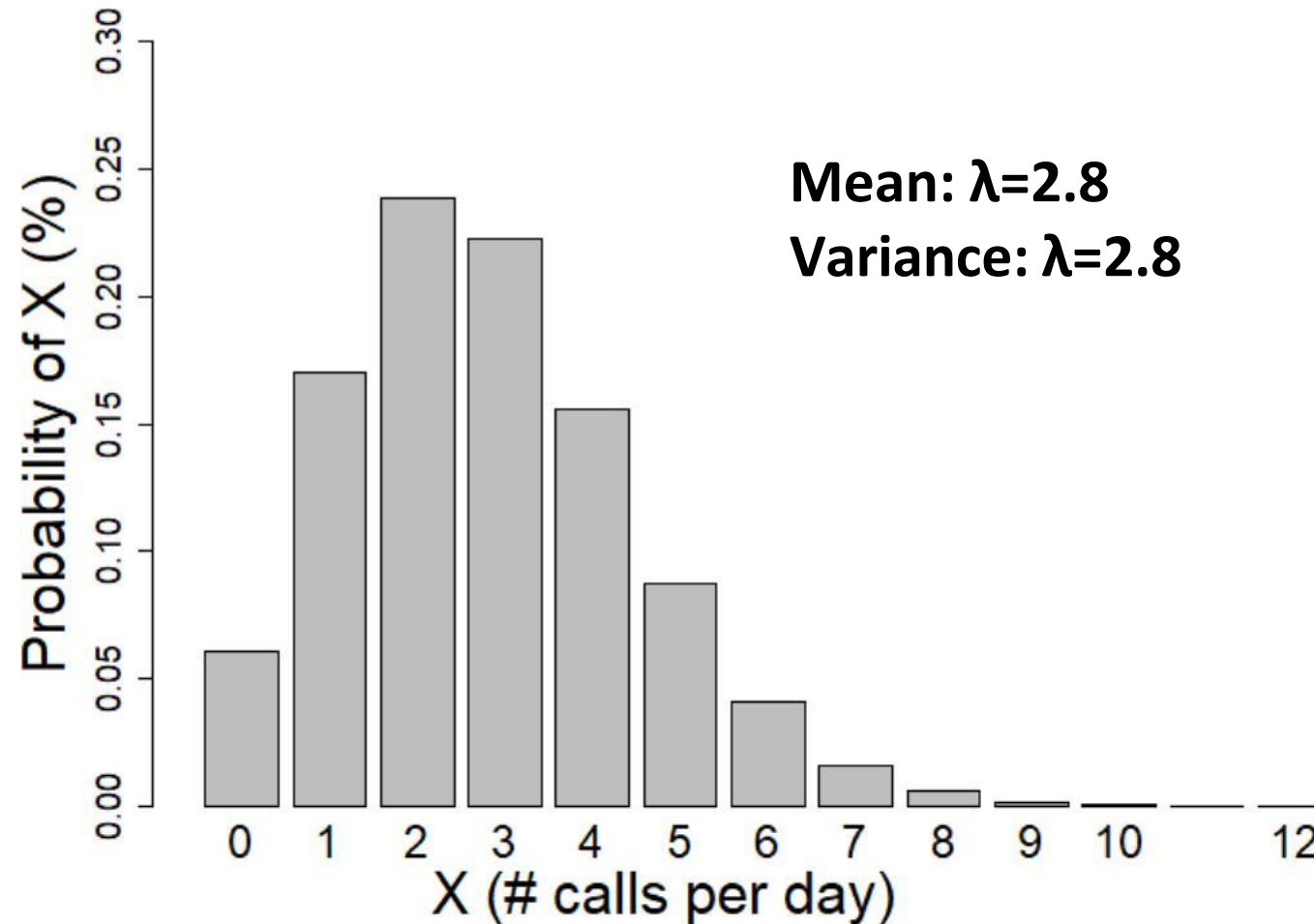
# Poisson Distribution

## mathematical definition

- $x \sim P(\lambda, \lambda)$
- $\lambda \equiv$  **Poisson Variable**
- $\mu = \lambda$  , population mean
- $\sigma^2 = \lambda$  , variance
- $P_{i \geq 0} = \frac{e^{-\lambda} \lambda^i}{i!}$  probability of i occurrences per unit
- When the number of trials (n) goes to infinity, and the expected number of successes ( $p_x$ ) is small then

$$F_{binomial}(x; n, p_x) = F_{Poisson}(x; \lambda = np_x)$$

# Poisson Distribution with Poisson variable $\lambda=2.8$



# Normal Distribution

- Most common distribution for the probability of a **continuous, real, random variable**

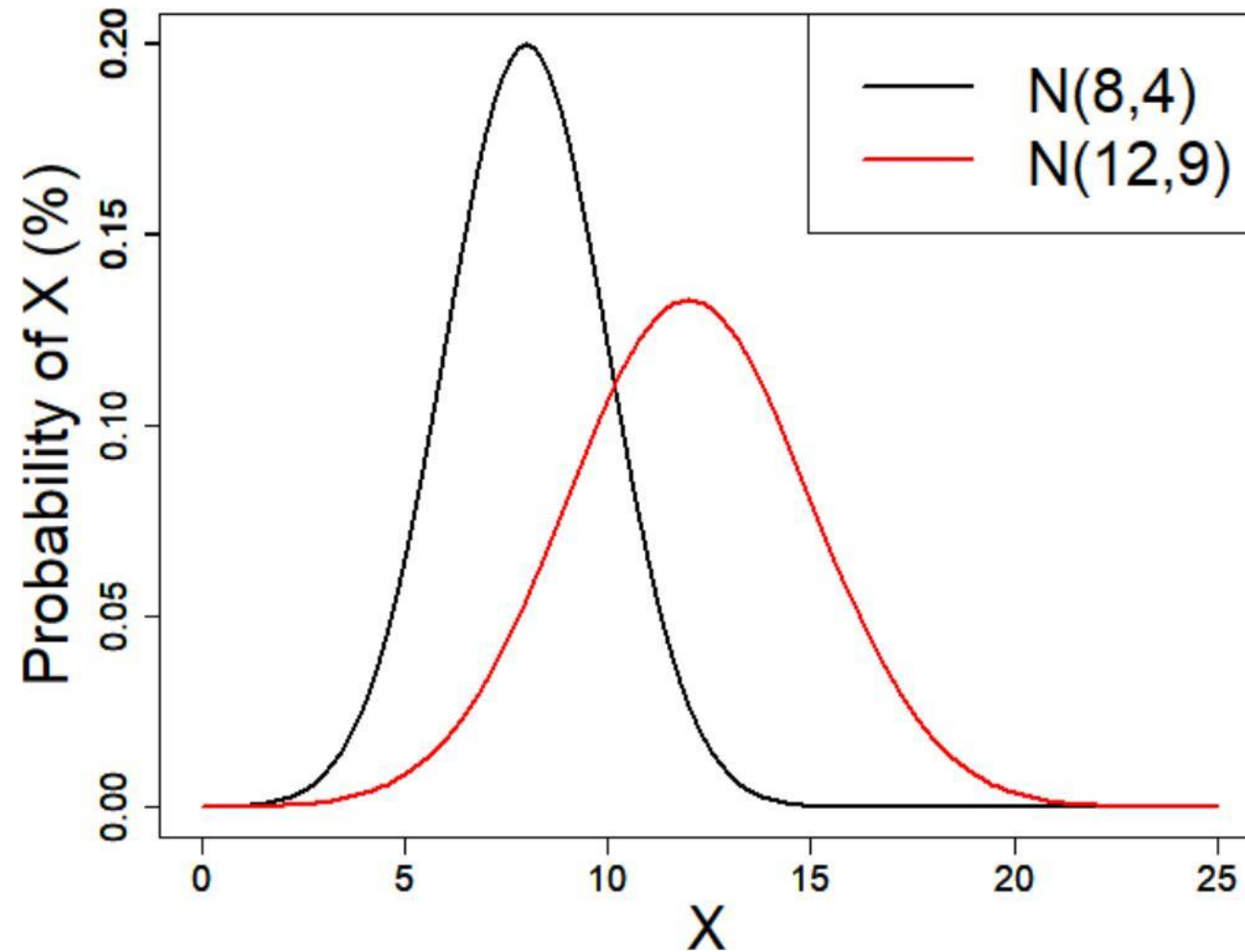
$$X \sim N(\mu, \sigma^2)$$

$\mu$   $\equiv$  Mean (and median)

$\sigma$   $\equiv$  Standard deviation

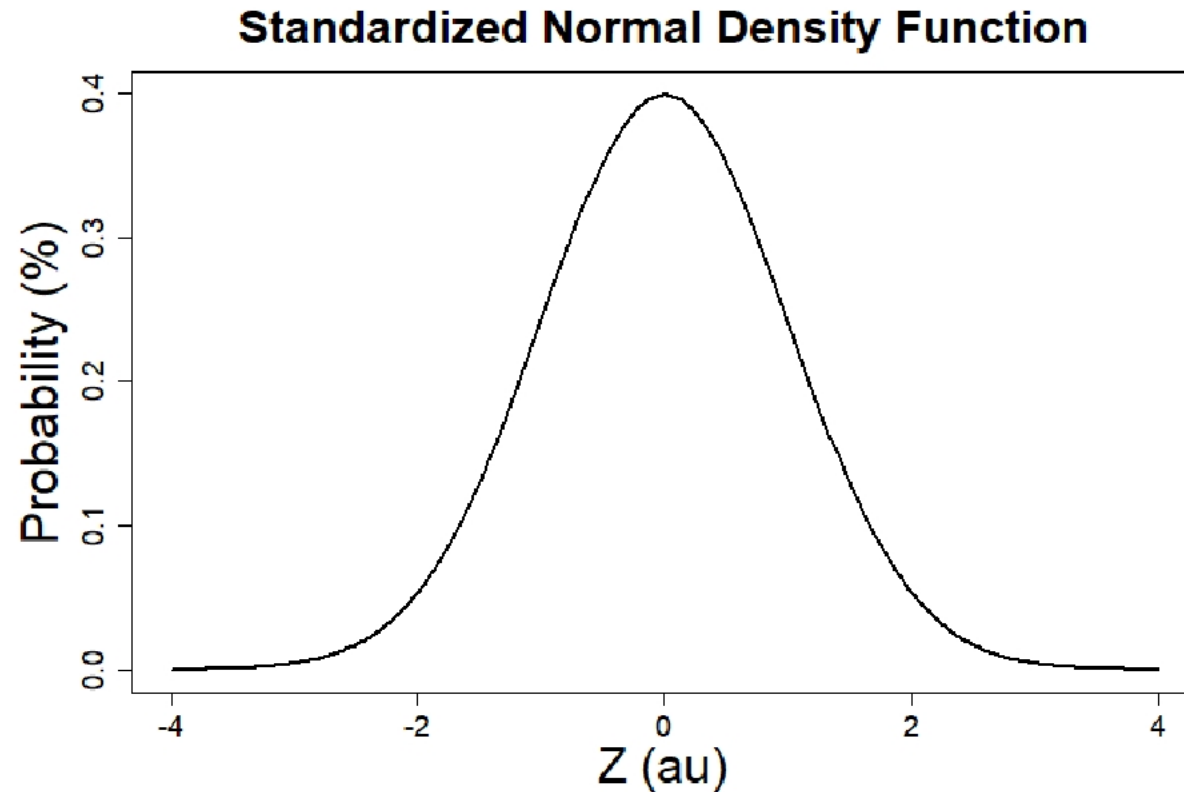
$$f(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \equiv \text{Distribution function}$$

# NormalDistribution

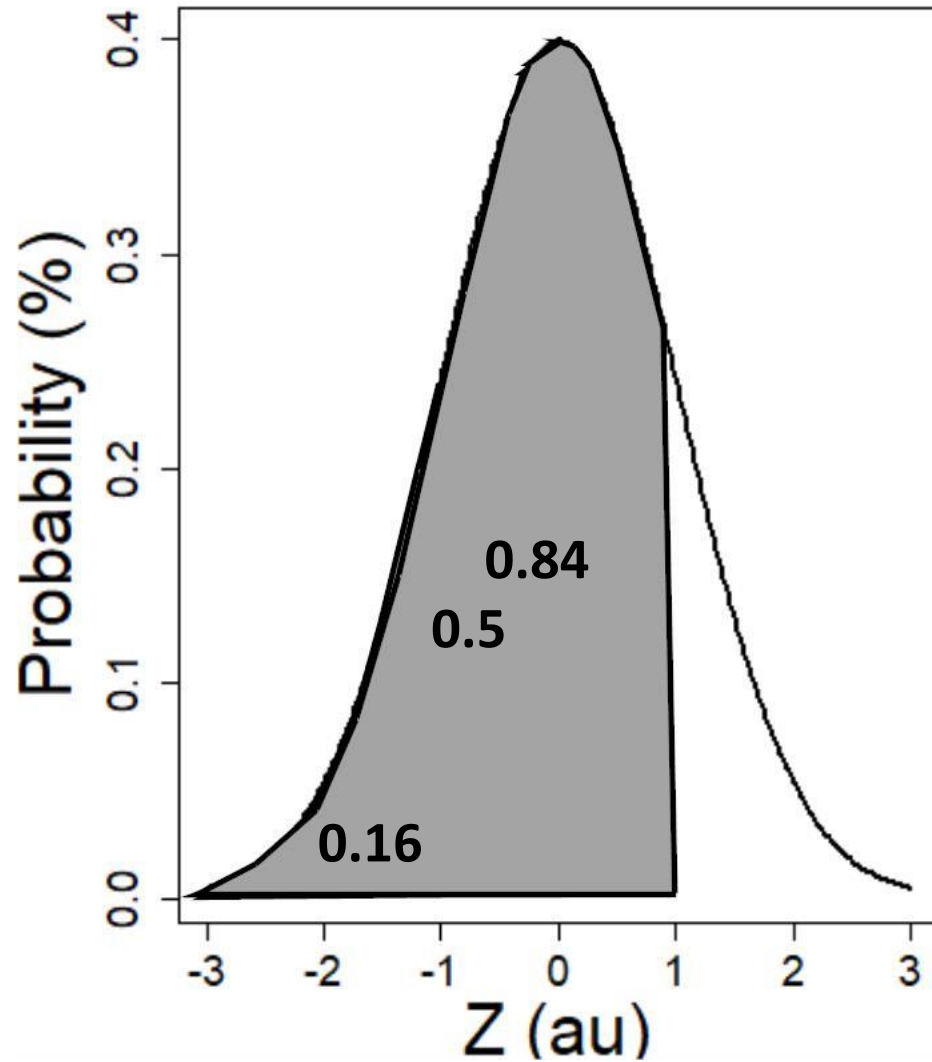


# The Standardized Normal Density Function

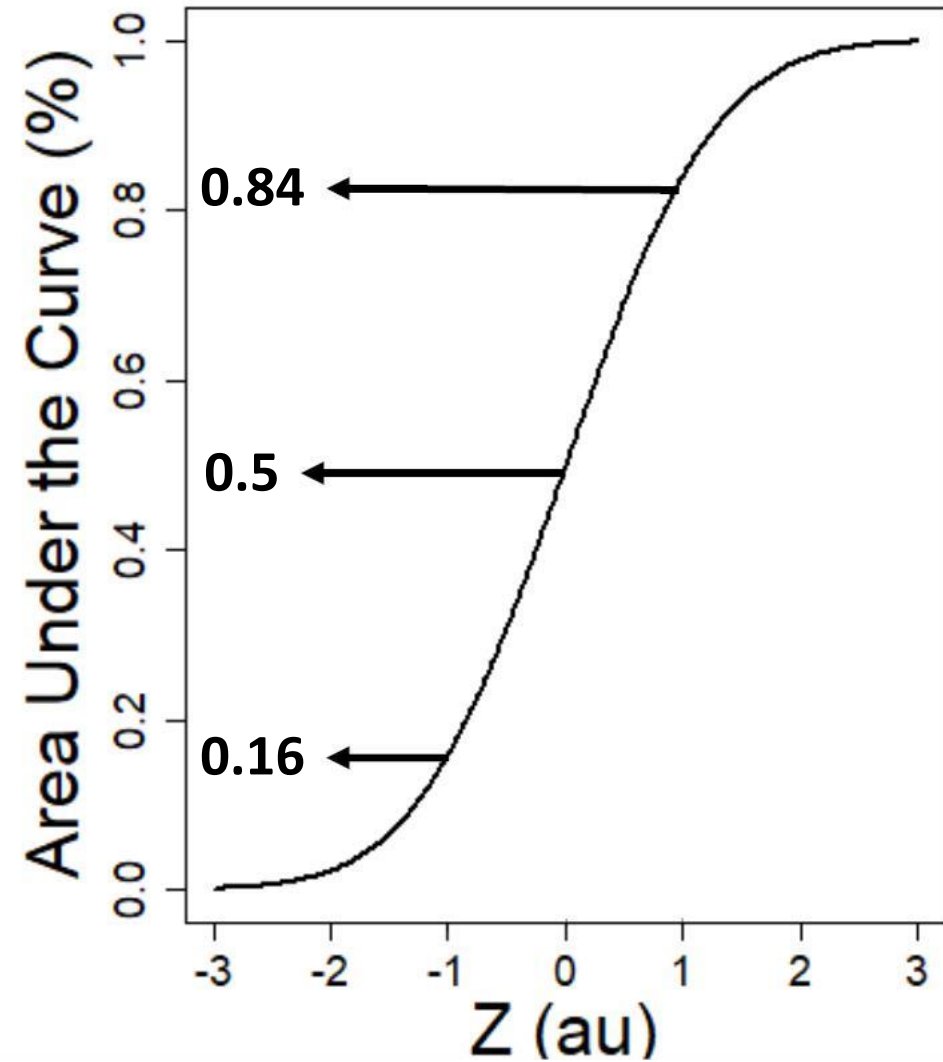
$$Z(x|\mu = 0, \sigma = 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad \text{Standard Distribution function}$$



### Normal Density Function

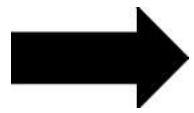


### Cumulative Distribution Function



# Why is the "Normal Distribution" important?

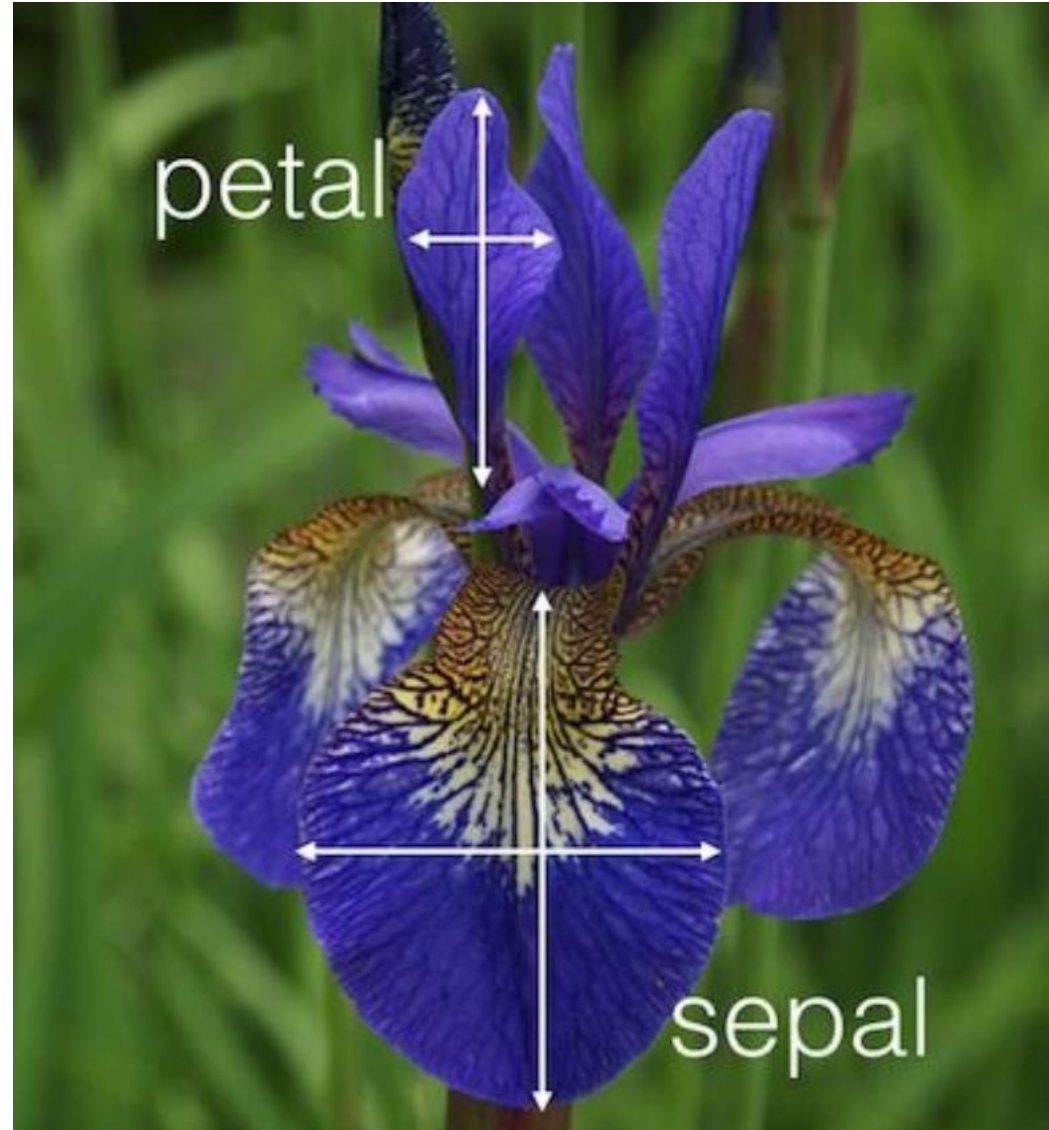
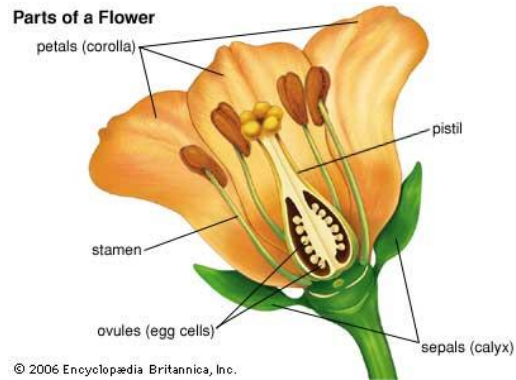
- The distribution of many continuous observables in nature is approximately normal
- We can describe the normal distribution in an equation



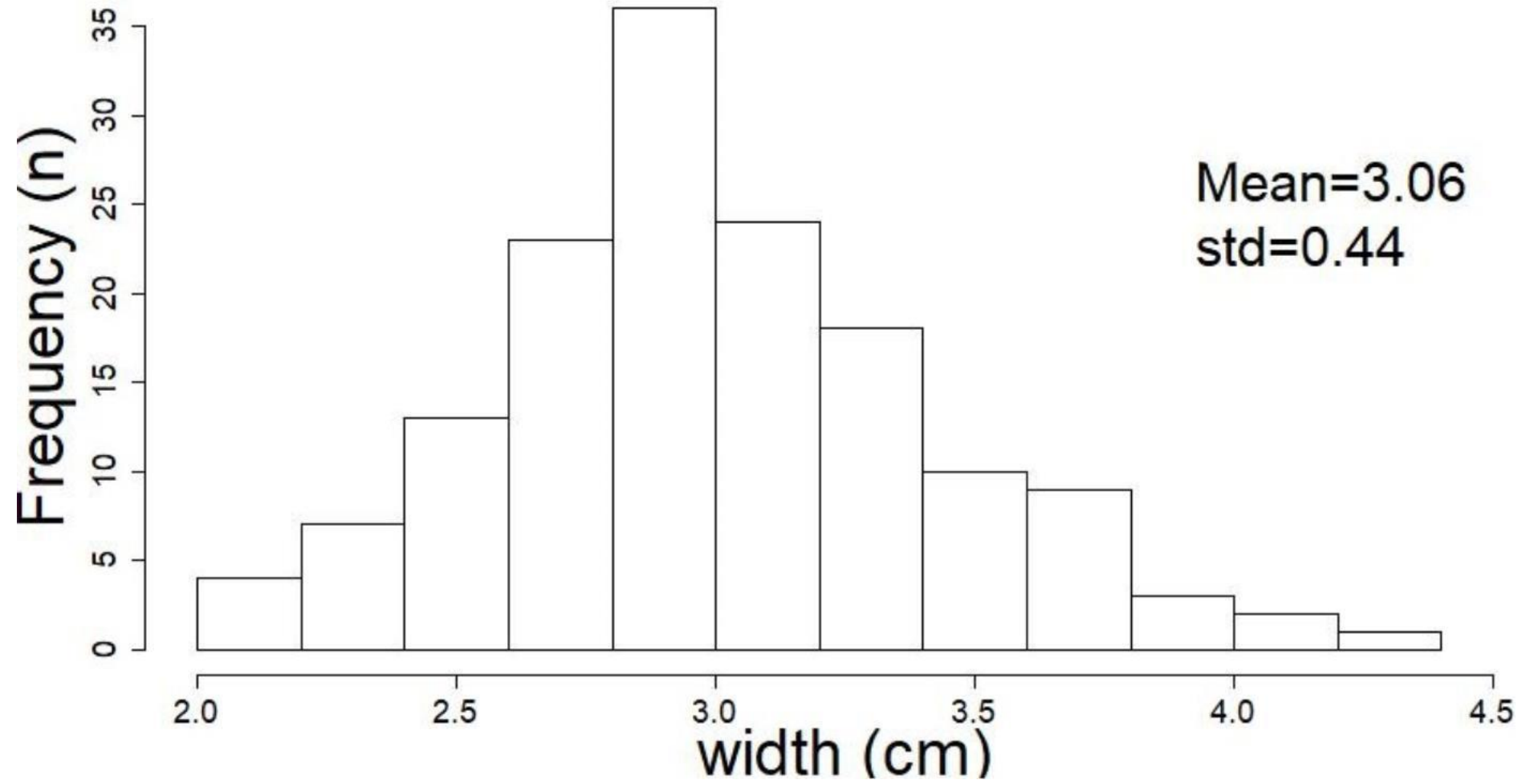
We can draw conclusions about the relations between variables with **Magnitude** and **Significance** estimates



# The Irisdataset



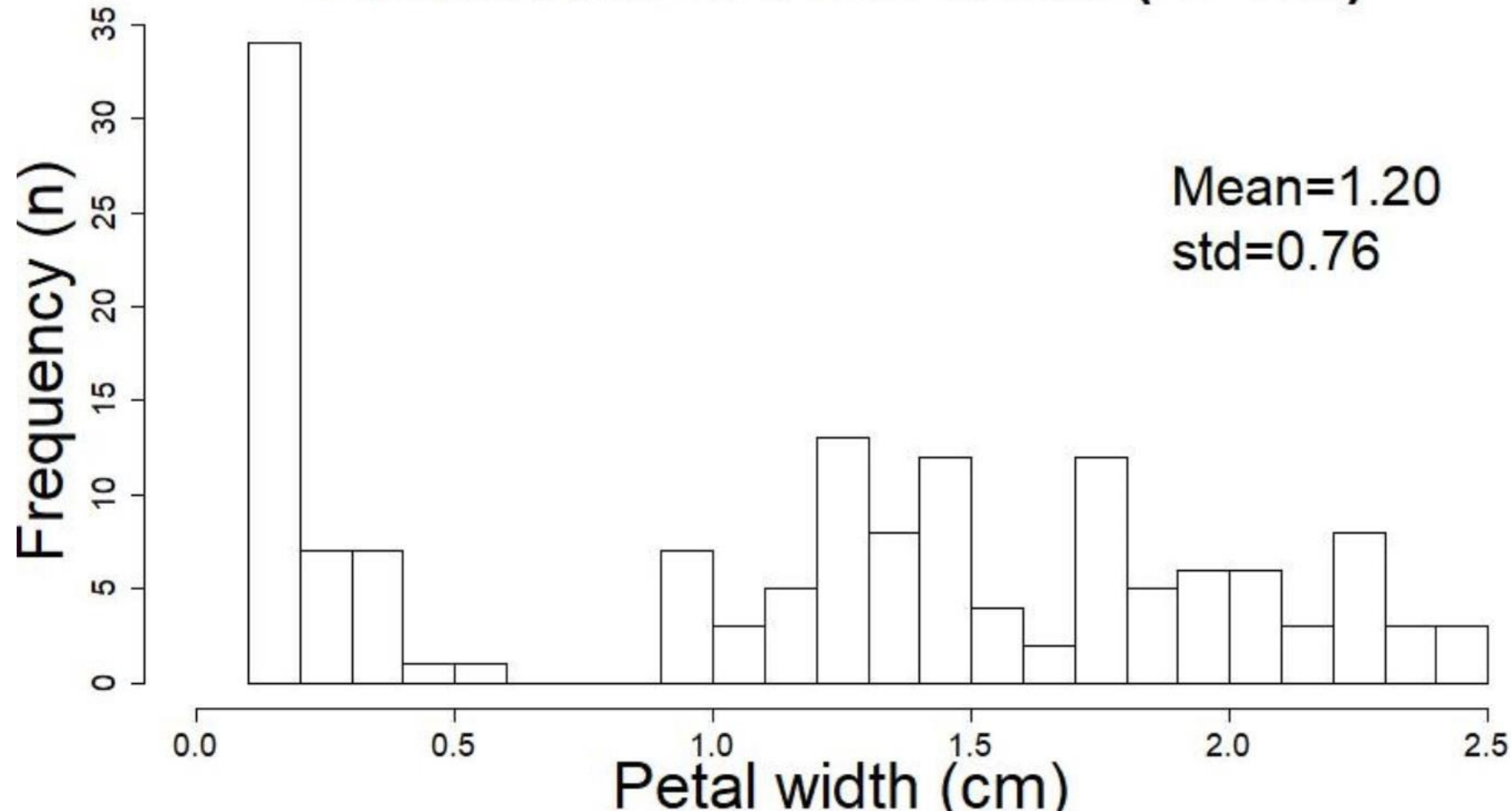
## Distribution of Speal Width (N=150)



# Conclusion:

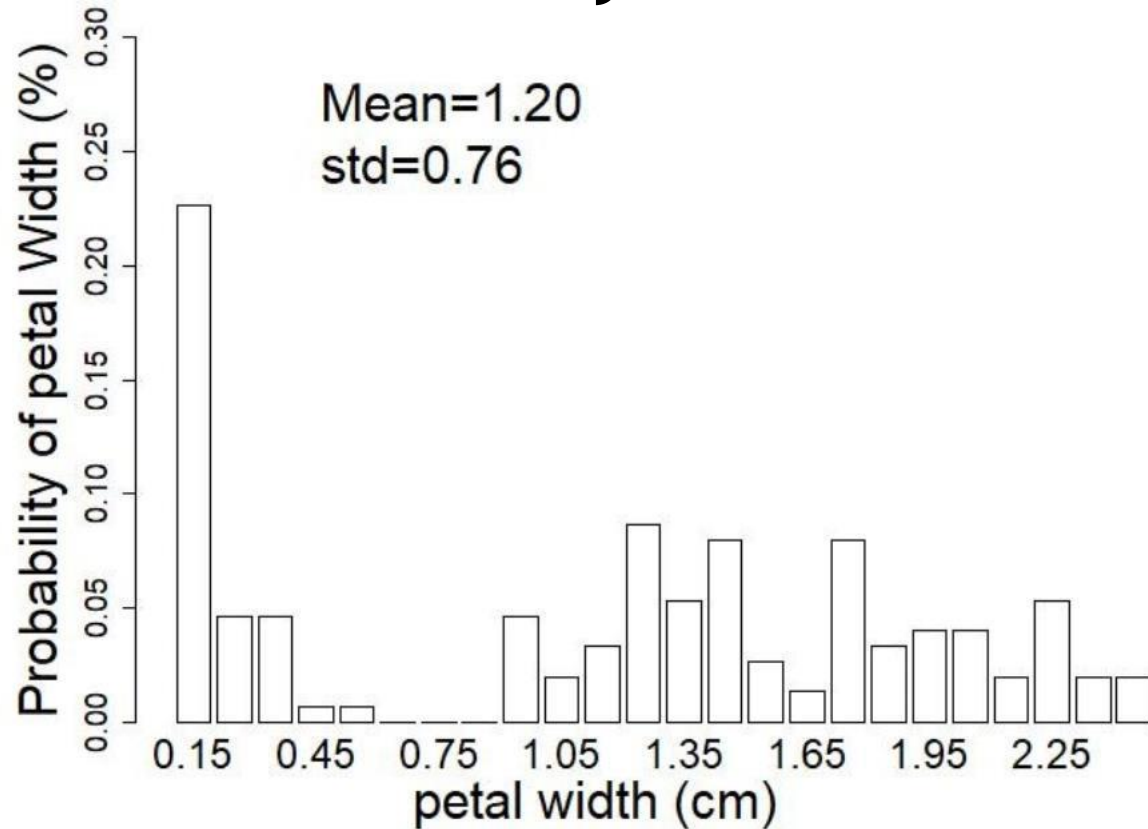
NOT ALL CONTINUOUS VARIABLES ARE NORMALLY DISTRIBUTED

**Distribution of Petal Width (N=150)**



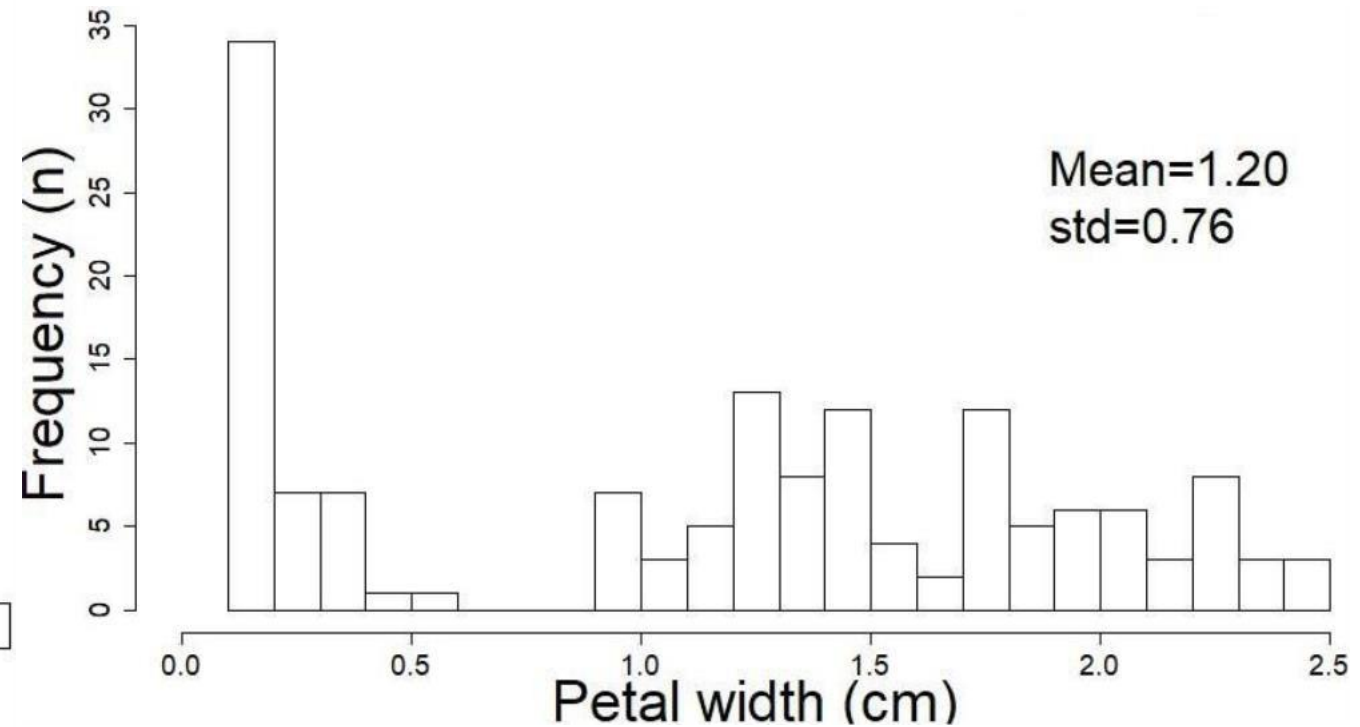
# Frequency vs. Probability Distributions

## Probability Distributions



*# frequency of counts per bin*

## Frequency Distributions



*# counts per bin*

# Boxplots and quantiles

