# 11. Connection between Module 1 (Intro to R) and Module 2 (Intro to Probability)

Principles of Data Science with R

Dr. Uma Ravat

PSTAT 10

## We did:

- Discrete and Continuous Random Variables and their distributions.
- PMF/PDF, plotting pmf/pdf, cdf, probability calculations by hand and using R for Discrete and Continuous uniform, Binomial and Normal distributions.

**binomial distribution** $\text{Binom}(\text{size, prob})$

- dbinom(x, size, prob)
- pbinom(q, size, prob)
- rbinom(n, size, prob)
- qbinom(p, size, prob)

**uniform distribution** $\text{Unif}(\text{min, max})$

- dunif(x, min, max)
- punif(q, min, max)
- runif(n, min, max)
- qunif(p, min, max)

**normal distribution** $N(\text{mean, sd})$

- dnorm(x, mean, sd)
- pnorm(q, mean, sd)
- rnorm(n, mean, sd)
- qnorm(p, mean, sd)

**Next we will see. . .**

- Connection modules and other courses in PSTAT department

- Some extras

    - Creating your own .Rmd file
    - Environment
    - Working directory

## Overall Connection

Module 1(Intro to R ) and Module 2(Introduction to Probability)

# Sample



The penguins dataset is stored in a **data frame** with

- **344 observations/samples/cases/subjects** (rows)
  - each case represents a penguin
- **8 variables** (columns)
  - species, island, bill_length_mm, bill_depth_mm etc
  - each corresponds to some measurement of the penguin

**Where do random varaibles come from?**

Recall, for data(a sample) we said a variable can be

- Numerical - discrete or continuous
- Categorical - ordinal or nominal

Random variables encode all possible data we may ever see!

## Why learn probability?

- Used plots and summary statistics to explore distributions and relationships of different variables in our (observed) data/sample.

- Now, Statistics aims to generalize these findings to the entire population.

## Population



**Random variables**
Population parameters
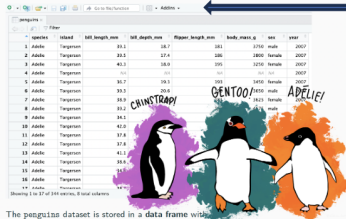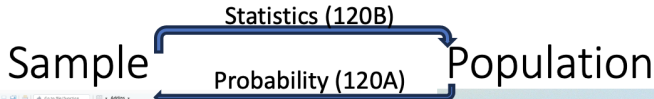- Population mean
- Population variance
Sampling distributions
Central Limit Theorem

**Statistics: Generalizing from a sample to the population**

- There's always some uncertainty about the true distributions and relationships in the population

- Probability is the mathematical tool used to measure and express this uncertainty. (PSTAT 120A)

- We should clearly specify the extent of our uncertainty. (PSTAT 120B)

Statistics (120B)

Sample ⟷ Population

Probability (120A)

The penguins dataset is stored in a **data frame** with:

- 344 observations/samples/cases/subjects (rows)
  - each case represents a penguin
- 8 variables (columns)
  - species, island, bill_length_mm, bill_depth_mm etc
  - each corresponds to some measurement of the penguin

**Variables**
Summary statistics
  - sample mean
  - sample variance
Visualizations

**Random variables**
Population parameters
  - Population mean
  - Population variance
Sampling distributions
Central Limit Theorem

## Courses that build on probability fundamentals

- Measure and express uncertainty in going from sample to population (PSTAT 120B)
- Hypothesis testing (PSTAT 120B)
- Bayesian statistics (PSTAT 115)
- Linear Regression (PSTAT 126)
- Statistical Machine Learning (PSTAT 131)
- Computational statistics (PSTAT 194CS)
  - Monte Carlo methods, Social Network Analysis, AI

**Dangers**

- Theory not used correctly

# Some extras [OPTIONAL]

E. Kale Edmiston PhD
@EKaleEdmiston

**Follow**

A friend/colleague who is an excellent programmer offhandedly told me the other day that coding is 90% googling error messages & 10% writing code. Until this point, I thought that all the time I spent googling error messages meant I was bad at coding. What a perspective change!

8:12 AM - 4 Jan 2019

**151** Retweets **1,069** Likes

💬 27      ↻ 151      ≋           ♡ 1.1K      ✉           ☑

## Debugging

Debugging is the process of getting rid of errors in your code.

3 types of errors:

1. Syntax Errors: code does not follow R's rules
2. Runtime Errors: errors that occur during knitting
3. Logic Errors: code runs but produces unexpected results.

## Know Your RStudio Environment

There are a *lot* of keyboard shortcuts in RStudio. To view all the options, you must engage the keyboard shortcut that rules them all:

- Windows: `Alt` + `Shift` + `K`
- macOS: `Option` + `Shift` + `K`

## Some favorites

1. Autocomplete command.
    - Both: `Tab`
2. Run the current line, selection from the editor.
    - Windows: `Ctrl` + `Enter`
    - macOS: `Cmd` + `Enter`
3. Run the current code chunk from the editor.
    - Windows: `Ctrl` + `Shift` + `Enter`
    - macOS: `Cmd` + `Shift` + `Enter`

## Downloading R

Go to: https://cran.r-project.org/

Chose from:

- Download R for (Mac) OS X
- Download R for Windows

Mac users choose Mac download

Windows users choose Windowns download

## Downloading RStudio

1. Download and install R first.
2. Go to https://rstudio.com/products/rstudio/download/

## We did

- Connection of modules and other courses in PSTAT department

- Some extras

  - Creating your own .Rmd file
  - Environment
  - Working directory