

## **6. Random Variables calculations and wrap up**

Transfer exploration seminar: Statistics and Data Science

---

Dr. Uma Ravat

PSTAT 194TR

## Summary: Random Variables

- Random Variables: Discrete , Continuous
  - distribution function: pmf  $P(X = x)$ , pdf  $f(x)$
  - probabilities : sums, series and integrals
  - Expected value:  $\mu$ ,  $\bar{X}$ ,  $E(X)$
  - Variance:  $\sigma^2$ ,  $Var(X)$

**Become familiar with the notation and concepts, the algebra will follow much more easily**

## Next we will see. . .

- Random variables
  - Using calculus (integrals and series) in probability calculations
- Wrap up with a look at connection of Intro to R and Intro to Probability.
- Feedback about these modules.
- Discuss anything else you want.

**Become familiar with the notation and concepts, the algebra will follow much more easily**

**Pre-reading: Math Review:**

- Sum and series
- Integrals

# Comparison of Discrete and Continuous Random Variables

| Random Variable                   | Discrete  | Continuous  |
|-----------------------------------|---|---|
| Takes on                          | finite or a countable number of distinct values   | any value in a given range/interval                         |
| <b>Example</b>                    | # of correct answers on a 100 question test       | Time taken to hike around the lagoon                        |
| <b>Prob distribution function</b> | probability mass function (PMF) $P(X = x)$        | probability density function (PDF) $f(x)$                   |
| <b>Properties</b>                 | $0 \leq P(X = x) \leq 1$<br>$\sum_x P(X = x) = 1$ | $0 \leq f(x)$ $\int_x f(x) = 1$ (total area under PDF = 1 ) |

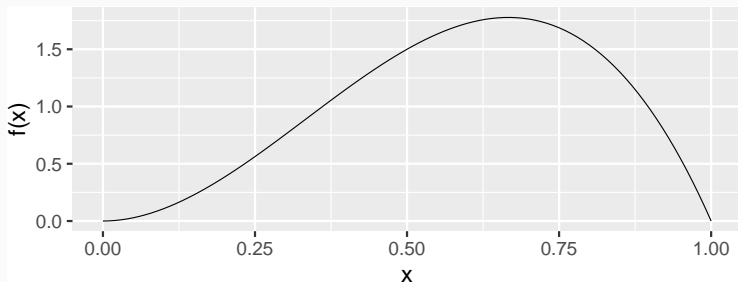
# Comparison of Discrete and Continuous Random Variables

| Random Variable                | Discrete   | Continuous   |
|--------------------------------|--|--|
| <b>Probability Calculation</b> | $P(X = x)$ gives the probability of a specific value | $P(a \leq X \leq b) = \int_a^b f(x) dx$ gives probability over an interval |
| <b>Mean (Expected Value)</b>   | $E(X) = \sum_x x \cdot P(X = x)$                     | $E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx$                           |
| <b>Variance</b>                | $Var(X) = \sum_x (x - E(X))^2 \cdot P(X = x)$        | $Var(X) = \int_{-\infty}^{\infty} (x - E(X))^2 \cdot f(x) dx$              |

## Revisit: Example of continuous random variable

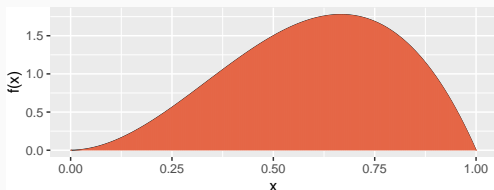
A statistical quality control example.

$$f(x) = \begin{cases} 12 x^2 (1 - x) & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$



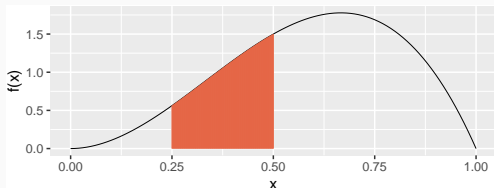
# Area under the curve = 1

$$f(x) = 12(x^2)(1-x), \quad 0 \leq x \leq 1$$



$$\begin{aligned} \int_{x \in S} f(x) dx &= \int_0^1 12(x^2)(1-x) dx \\ &= 12 \int_0^1 (x^2 - x^3) dx = 12 \left[ \frac{x^3}{3} - \frac{x^4}{4} \right]_0^1 = 1 \end{aligned}$$

# Probability is Area Under the Curve



$$\begin{aligned} P(0.25 < X < 0.50) &= \int_{0.25}^{0.50} 12(x^2)(1-x)dx = 12 \int_{0.25}^{0.50} (x^2 - x^3)dx \\ &= 12 \left[ \frac{x^3}{3} - \frac{x^4}{4} \right]_{0.25}^{0.50} = 0.2617188 \end{aligned}$$



## Expected value

$$f(x) = \begin{cases} 12x^2(1-x) & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} E(X) &= \int_{x \in S} xf(x)dx \\ &= \int_0^1 x12(x^2)(1-x)dx \\ &= 12 \int_0^1 (x^3 - x^4)dx \\ &= 12 \left[ \frac{x^4}{4} - \frac{x^5}{5} \right]_0^1 \\ &= \boxed{0.6} \end{aligned}$$

# Variance

$$f(x) = 12 x^2 (1 - x) \text{ if } 0 \leq x \leq 1$$

$$\begin{aligned} \text{Var}(X) &= \int_0^1 (x - E(X))^2 \cdot f(x) dx \\ &= \int_0^1 (x - 0.6)^2 \cdot 12 \cdot x^2 \cdot (1 - x) dx \\ &= \dots \\ &= \dots \\ &= \dots \\ &= \dots \\ &= \dots \\ &= \boxed{0.04} \end{aligned}$$

# Variance

$$f(x) = 12 x^2 (1 - x) \text{ if } 0 \leq x \leq 1$$

$$\begin{aligned} \text{Var}(X) &= \int_0^1 (x - E(X))^2 \cdot f(x) dx \\ &= \int_0^1 (x - 0.6)^2 12 x^2 (1 - x) dx = \int_0^1 (x^2 - 1.2x + 0.36) 12 x^2 (1 - x) dx \\ &= 12 \int_0^1 (x^4 - 1.2x^3 + 0.36x^2)(1 - x) dx = 12 \int_0^1 (x^4 - 1.2x^3 + 0.36x^2) - (x^5 - 1.2x^4 + 0.36x^3) dx \\ &= 12 \int_0^1 (x^4 - 1.2x^3 + 0.36x^2) - x^5 + 1.2x^4 - 0.36x^3 dx = 12 \int_0^1 -x^5 + 2.2x^4 - 1.56x^3 + 0.36x^2 dx \\ &= 12 \left[ \frac{-x^6}{6} + \frac{2.2x^5}{5} - \frac{1.56x^4}{4} + \frac{0.36x^3}{3} \right]_0^1 \\ &= 12 \left[ -16 + \frac{2.2}{5} - \frac{1.56}{4} + \frac{0.36}{3} \right] \\ &= \boxed{0.04} \end{aligned}$$

## Variance: Equivalent definitions

$$\begin{aligned} \text{Var}(X) &= \int_{-\infty}^{\infty} (x - E(X))^2 \cdot f(x) dx \\ &= E[(X - E(X))^2] \\ &= E(X^2) - [E(X)]^2 \end{aligned}$$

$$E(X^2) = ?$$

$$\begin{aligned} E(X^2) &= \int_{-\infty}^{\infty} x^2 f(x) dx = \int_0^1 x^2 12(x^2)(1-x) dx \\ &= 12 \int_0^1 (x^4 - x^5) dx \\ &= 12 \left[ \frac{x^5}{5} - \frac{x^6}{6} \right]_0^1 \\ &= \boxed{0.4} \end{aligned}$$

$$\text{Var}(X) = 0.4 - 0.6^2 = \boxed{0.04}$$

## Example: Flipping two coins

$X$  = number of heads in two independent coin flips

$X$  is a \_\_\_\_\_ random variable.

Sample space is  $S = \{ \text{_____} \}$

PMF of  $X$  is \_\_\_\_\_

## Example: Flipping two coins

$X$  = number of heads in two independent coin flips

$X$  is a discrete random variable.

Sample space is  $S = \{HH, HT, TH, TT\} = \{0, 1, 2\}$

PMF is

| Values: $X = x$                | 0             | 1             | 2             |
|--------------------------------|---------------|---------------|---------------|
| Probability: $P(X = x) = P(x)$ | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{1}{4}$ |

## Expected value and variance

$X$  = number of heads in two independent coin flips,  $S_X = \{0, 1, 2\}$

$$E(X) = \sum_x xP(x)$$

| $x$     | 0                    | 1                              | 2                              |            |
|---------|----------------------|--------------------------------|--------------------------------|------------|
| $P(x)$  | $\frac{1}{4}$        | $\frac{1}{2}$                  | $\frac{1}{4}$                  |            |
| $xP(x)$ | $0 = (0)\frac{1}{4}$ | $\frac{1}{2} = (1)\frac{1}{2}$ | $\frac{1}{2} = (2)\frac{1}{4}$ | $1 = E(X)$ |

$$\text{So, } E(X) = \sum_x xP(x) = 1$$

On average, you expect to see 1 head if you flip two coins.

## Variance

$X$  = number of heads in two independent coin flips,  $S_X = \{0, 1, 2\}$

$$\text{Var}X = E(X^2) - [E(X)]^2 = \sum_{\text{all } x} x^2 P(x) - [E(X)]^2$$

$$E(X) = 1, E(X^2) = ?$$

| $x$       | 0                      | 1                                | 2                              |                |
|-----------|------------------------|----------------------------------|--------------------------------|----------------|
| $P(x)$    | $\frac{1}{4}$          | $\frac{1}{2}$                    | $\frac{1}{4}$                  |                |
| $xP(x)$   | $0 = (0)\frac{1}{4}$   | $\frac{1}{2} = (1)\frac{1}{2}$   | $\frac{1}{2} = (2)\frac{1}{4}$ | $1 = E(X)$     |
| $x^2P(x)$ | $0 = (0^2)\frac{1}{4}$ | $\frac{1}{2} = (1^2)\frac{1}{2}$ | $1 = (2^2)\frac{1}{4}$         | $1.5 = E(X^2)$ |

$$\text{Var}(X) = E(X^2) - [E(X)]^2 = 1.5 - 1^2 = 0.5$$



## Recall: Series and it's sum

If  $b$  is any number, the geometric series and it's sum is given by

$$\sum_{k=0}^{\infty} b^k = (1 + b^1 + b^2 + \dots + b^k + \dots) = \frac{1}{1-b} = \frac{\text{first term}}{1 - \text{base}}, \text{ if } |b| < 1$$

$$k! = 1 \cdot 2 \cdot 3 \cdots (k-1) \cdot k \text{ eg. } 5! = 1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 = 120$$

If  $a$  is any number, the exponential series and it's sum is given by

$$\begin{aligned} \sum_{k=0}^{\infty} \frac{a^k}{k!} &= \frac{a^0}{0!} + \frac{a^1}{1!} + \frac{a^2}{2!} + \dots + \frac{a^k}{k!} + \dots \\ &= 1 + a + \frac{a^2}{2!} + \dots + \frac{a^k}{k!} + \dots \\ &= e^a \end{aligned}$$

See also, series review and math review sheet.

## Example: uncountable sample space

Find  $E(X)$  of a discrete random variable  $X$  that has the following probability distribution:

$$P(X = 0) = P(0) = 2 - \sqrt{e}; P(X = k) = P(k) = \frac{1}{2^k \cdot k!}, \quad k = 1, 2, 3, \dots$$

$$\begin{aligned} E(X) &= \sum_x xP(x) = 0(2 - \sqrt{e}) + \sum_{k=1}^{\infty} k \cdot \frac{1}{2^k k!} \\ &= \sum_{k=1}^{\infty} \frac{\left(\frac{1}{2}\right)^k}{(k-1)!} \\ &= \frac{1}{2} \sum_{k-1=0}^{\infty} \frac{\left(\frac{1}{2}\right)^{(k-1)}}{(k-1)!} \\ &= \frac{1}{2} \sum_{n=0}^{\infty} \frac{\left(\frac{1}{2}\right)^{(n)}}{(n)!} \quad (k-1 = n) \\ &= \frac{1}{2} e^{1/2} = \frac{1}{2} \sqrt{e} \end{aligned}$$

## Summary:

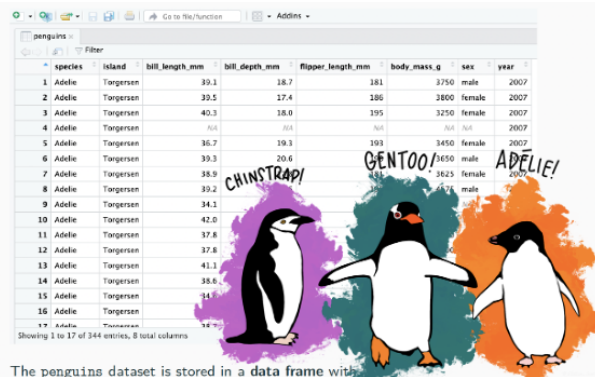
- Probability,  $E(X)$ ,  $V(X)$  calculations using
  - integrals for continuous random variables
  - series for discrete random variables with uncountable sample space.

### Overall summary:

- Study skills for success in PSTAT courses
- Introduction to R (in preparation for PSTAT 10)
- Introduction to Probability (in preparation for PSTAT 120a)
  - Also, review Math review slides provided.
  - Review Double Integrals

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

## Sample



|    | species | island    | bill_length_mm | bill_depth_mm | flipper_length_mm | body_mass_g | sex    | year |
|----|---------|-----------|----------------|---------------|-------------------|-------------|--------|------|
| 1  | Adelie  | Tongersen | 39.1           | 18.7          | 181               | 3750        | male   | 2007 |
| 2  | Adelie  | Tongersen | 39.5           | 17.4          | 186               | 3800        | female | 2007 |
| 3  | Adelie  | Tongersen | 40.3           | 18.0          | 195               | 3250        | female | 2007 |
| 4  | Adelie  | Tongersen | NA             | NA            | NA                | NA          | NA     | 2007 |
| 5  | Adelie  | Tongersen | 36.7           | 19.3          | 193               | 3450        | female | 2007 |
| 6  | Adelie  | Tongersen | 39.3           | 20.6          | 193               | 3650        | male   | 2007 |
| 7  | Adelie  | Tongersen | 38.9           |               | 193               | 3625        | female | 2007 |
| 8  | Adelie  | Tongersen | 39.2           |               | 193               | 4675        | male   | 2007 |
| 9  | Adelie  | Tongersen | 34.1           |               |                   |             |        |      |
| 10 | Adelie  | Tongersen | 42.0           |               |                   |             |        |      |
| 11 | Adelie  | Tongersen | 37.8           |               |                   |             |        |      |
| 12 | Adelie  | Tongersen | 37.8           |               |                   |             |        |      |
| 13 | Adelie  | Tongersen | 41.1           |               |                   |             |        |      |
| 14 | Adelie  | Tongersen | 38.6           |               |                   |             |        |      |
| 15 | Adelie  | Tongersen | 41.1           |               |                   |             |        |      |
| 16 | Adelie  | Tongersen | 38.6           |               |                   |             |        |      |
| 17 | Adelie  | Tongersen | 36.7           |               |                   |             |        |      |

Showing 1 to 17 of 344 entries, 8 total columns

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

## Population



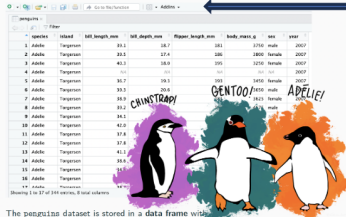
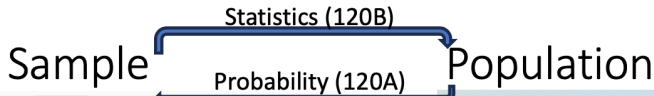
### **Random variables**

Population parameters

- Population mean
- Population variance

Sampling distributions

Central Limit Theorem



The screenshot shows a R console window with the 'penguins' dataset loaded. The dataset has 17 columns: species, island, bill\_length\_mm, bill\_depth\_mm, flipper\_length\_mm, body\_mass\_g, sex, and year. The first 16 rows are displayed, showing data for Chinstrap, Gentoo, and Adelie penguins. Below the table, there are three cartoon penguins with labels: CHINSTRAP! (black and white), GENTOO! (black and white), and ADELIE! (black and white).

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

## Overall summary:

- Study skills for success in PSTAT courses
- Introduction to R (in preparation for PSTAT 10)
- Introduction to Probability (in preparation for PSTAT 120a)
  - Also, review Math review slides provided.
  - Review Double Integrals

Next, Intro to Python