# Inference and Modeling Intro

Spring R '24

Dominic Bordelon, Research Data Librarian, ULS

University of Pittsburgh | Library System

# Agenda

1. Inference and hypothesis testing

    1. Distribution functions

    2. Hypothesis tests

2. Modeling (simple linear regression)

University of Pittsburgh | Library System

# About the trainer

**Dominic Bordelon, Research Data Librarian**

University Library System, University of Pittsburgh

dbordelon@pitt.edu

Services for the Pitt community:

- Consultations
- Training (on-request and via public workshops)
- Talks (on-request and publicly)
- Research collaboration

Support areas and interests:

- Computer programming fundamentals, esp. for data processing and analysis
- Open Science and Data Sharing
- Data stewardship/curation
- Research methods; science and technology studies

# Today's packages

`stats` (part of base R, automatically attaches)

`infer`: tidy inference

`tidymodels`, particularly:

- `broom`: tidy model representation

- `parsnip`: standardized modeling interface

```
1  install.packages(c("tidymodels", "infer"))
```

University of Pittsburgh | Library System

# …and using penguins examples

```r
install.packages("palmerpenguins")
```

```r
library(palmerpenguins)

# load palmerpenguins' data into your environment:
data(penguins)
names(penguins)
```
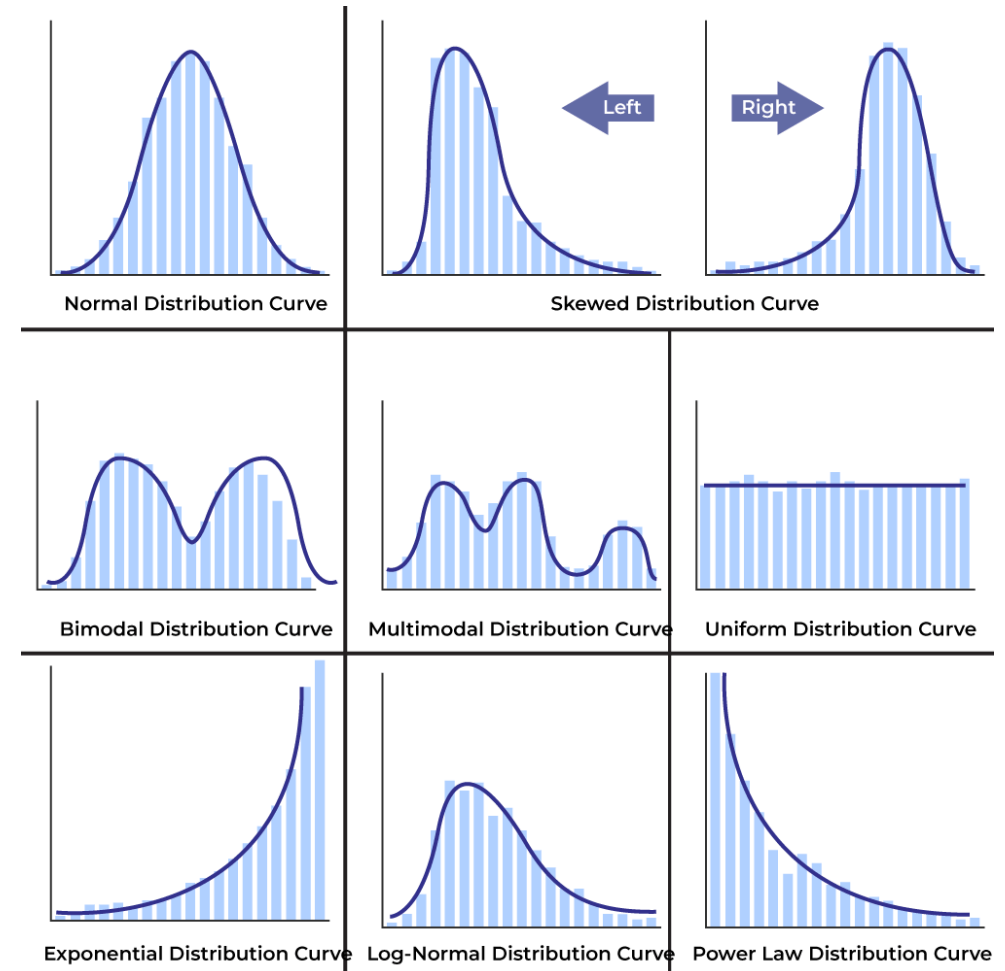
```
[1] "species"         "island"            "bill_length_mm"
[4] "bill_depth_mm"   "flipper_length_mm" "body_mass_g"
[7] "sex"             "year"
```
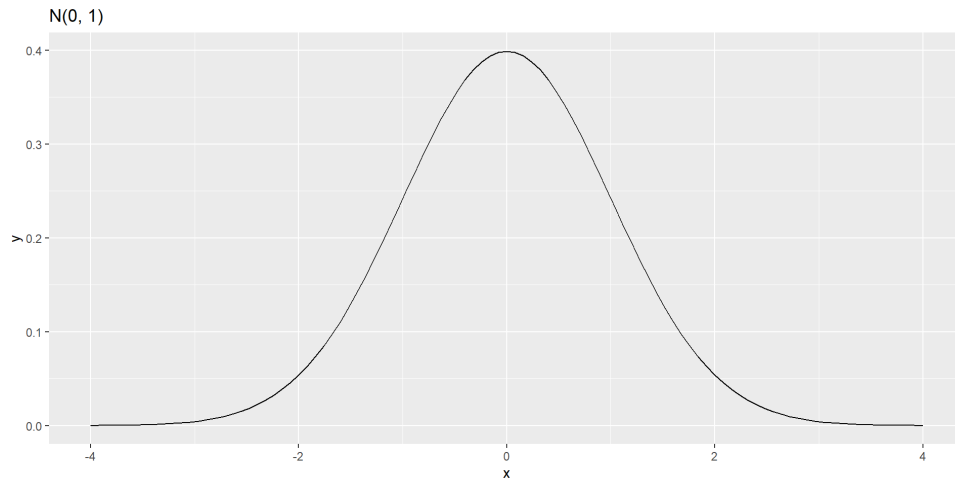
# Inference and hypothesis testing

# Frequency distributions

- In statistics we want to estimate parameters of the *population* using observed *samples*

- In frequentist or parametric statistics, we assume that the population distribution can be approximated with standard forms such as the Gaussian (normal) distribution



A variety of distributions. Image source: Geeks for Geeks

# Distribution functions

N(0, 1)



Probability distributions in R are typically described by four functions. The examples here are for the Normal distribution.

All have `mean` and `sd` arguments.

- `dnorm(x)` for density (height of curve) at $x$

- `pnorm(q)` for getting a probability $p$ value at some $q$, i.e., $P(X \leq q)$

- `qnorm(p)` for finding the quantile of some $p$

- `rnorm(n)` for generating a random sample of size $n$ from $N(\mu, \sigma)$

University of Pittsburgh | Library System

# Distributions

Commonly used distributions in the `stats` package

| Distribution | Name in functions | Example function | Applications |
|---|---|---|---|
| Normal (Gaussian) | norm | pnorm() | Numerous (often assumed) |
| Student's $t$ | t | pt() | Estimating parameters without knowing $\sigma$; $t$-tests |
| Binomial (Bernoulli) | binom | pbinom() | Number of successes in `size` trials; logistic regression |
| Chi-squared ($\chi^2$) | chisq | pchisq() | Chi-squared tests (goodness of fit in a 2-way table) |
| $F$ | f | pf() | ANOVA; $F$-tests (model goodness-of-fit) |

# Hypothesis tests

Hypothesis tests in base R tend to follow this format:

- a function call

- accepts vector or data frame inputs

- returns an object which can be assigned

- can be calculated "by hand" using other R functions

    - example: `t.test()` can also be found by calculating the $t$ test statistic, followed by `pt()` to obtain a $p$

University of Pittsburgh | Library System

# Student's $t$-test

One-sample $t$-test

"Is sample mean $\bar{x}$ different from population mean $\mu$?"

$$t = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}$$

```r
1  mu <- 4147
2  x_bar <- mean(penguins$body_mass_g, na.rm=TRUE)
3  s <- sd(penguins$body_mass_g, na.rm=TRUE)
4  n <- penguins[!is.na(penguins$body_mass_g), "body_
5    pull() %>%
6    length()
7
8  t <- (x_bar - mu) / (s / sqrt(n))
9
10 # find p for this t score, on t distribution with
11 pt(q = t,
12    df = (n - 1),
13    lower.tail = TRUE)
14
15 # or:
16 t.test(x = penguins$body_mass_g,
17        mu = mu)
```

Two-sample $t$-test (independent)

"Is sample mean $\bar{X}_1$ different from sample mean $\bar{X}_2$?"

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_p \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

```r
1  adelie <- penguins %>%
2    filter(species == "Adelie")
3  chinstrap <- penguins %>%
4    filter(species == "Chinstrap")
5
6  t.test(x = adelie$body_mass_g,
7         y = chinstrap$body_mass_g)
```

University of Pittsburgh | Library System

# Chi-squared ($\chi^2$) test

"Is the conditional distribution across two categorical variables equal, or varied?"

`chisq.test(x, y = NULL)` where `x` is a matrix, or `x` and `y` are vectors

`chisq_test(formula)` where `formula` is written in the format `response ~ explanatory`

```
1  # how are penguins distributed across species and islands?
2  chisq_test(penguins, island ~ species)
```

⚠ This is not a good example of a chi-squared test, statistically speaking: our data do not meet the requirements for the test! But the example is left here because R will let us do it! (and it is syntactically correct)

University of Pittsburgh | Library System

# Tidy inference: infer

```r
1  t <- penguins %>%
2    specify(response = body_mass_g) %>%
3    hypothesize(null = "point", mu = 4147) %>%
4    calculate(stat = "t") %>%
5    pull()
6
7  penguins %>%
8    specify(response = body_mass_g) %>%
9    assume(distribution = "t") %>%
10   visualize() +
11   shade_p_value(obs_stat = t, direction = "greater")
12
13 penguins %>%
14   t_test(response = body_mass_g,
15          mu = 4147)
16
17 penguins %>%
18   filter(species %in% c("Adelie", "Chinstrap")) %>%
19   t_test(formula = body_mass_g ~ species)
```

University of Pittsburgh | Library System

# Modeling

University of Pittsburgh | Library System

# Simple linear regression

"What is the relationship between $X$ and $Y$?"

Considering

$$Y = \beta_0 + \beta_1 X + \epsilon$$

where $\beta_0$ is the $y$ intercept, $\beta_1$ is the slope coefficient, $X$ is the explanatory or predictor variable, $\epsilon$ is irreducible error, and $Y$ is the response variable. We want to estimate (or "fit") $\hat{\beta}_0$ and $\hat{\beta}_1$ such that the sum of squared residuals is minimized as much as possible.

`lm(formula, data)` fits a linear model and returns a model object, where formula is written `response ~ explanatory` and `data` is a data frame.

```
1  lm(flipper_length_mm ~ bill_length_mm,
2    data = penguins) %>%
3   summary()
4
5  lm_flipper_bill <- lm(flipper_length_mm ~ bill_length_mm, data = penguins)
```

# predict() outputs from a model

"For some $x$, what $\hat{y}$ does the model predict?"

Models in R implement behavior for the predict() function, which supplies a data frame of input values $X$ and returns a data frame containing corresponding $Y$ values. The input data frame needs variable(s) of the same name(s) as the predictor(s). Or you can omit the data frame and get the fitted values.

💡 predict() functions are found in documentation under article names like predict.lm, even though the function call is predict().

```r
1  # min to max values, incrementing by 0.5:
2  x_range <- data.frame(bill_length_mm = seq(from=min(penguins$bill_length_mm, na.rm=TRUE), to=max(penguins$bil
3
4  predict(lm_flipper_bill, newdata = x_range)
5
6  # fitted values
7  predict(lm_flipper_bill)
```

University of Pittsburgh | Library System

# Tidy modeling: **parsnip**

```
1  linear_reg() %>%
2    set_engine("lm") %>%
3    fit(flipper_length_mm ~ bill_length_mm,
4        data = penguins)
5
6  linear_reg() %>%
7    set_engine("lm") %>%
8    fit(flipper_length_mm ~ bill_length_mm,
9        data = penguins)
```

University of Pittsburgh | Library System

# Model *responsibly*

- Consider the validity and reliability of your measures

- Correlation $\neq$ causation

- Always keep in mind the assumptions made by your chosen method/model

- Validate unexpected results by recomputing "by hand" and/or running a different code implementation (or even in a different software)

# Wrap up

University of Pittsburgh | Library System

# Session in review

Today we learned about:

- how distributions and inference work in R

- linear modeling in R

- old and new ways of doing these things

Join us next week for machine learning intro!

University of Pittsburgh | Library System