# STA130 Study Guide

| Data type | Symbol | Description | Examples |
|---|---|---|---|
| Integer | int | Numbers (w/ decimals) | 1, - 2, 3, - 4 |
| Double | dbl | Numbers (w/ or w/o decimals) | 2, 2.02, 22222 |
| Logical | lgl | TRUE or FALSE | TRUE, FALSE |
| Character | chr | Words, surrounded by quotation marks | "I", "love", "stats" |
| Factor | fct | Looks like "character" type, but only take values from a prespecified list | If continents list →can't take "blue" |

| Operator | Syntax |
|---|---|
| equal | == |
| not equal | ! = |
| less than (less than or equal to) | < (<=) |
| greater than (greater than or equal to) | > (>=) |
| not | ! |
| and | & |
| or | \| |

## Vects

- c() function combines single elements into a vector
- is. functions to check the data type of a vector (e.g. is.numeric(), is.character())

## Coercion

- R switches between data types automatically for certain operations
- ex) sum(c(TRUE, FALSE)) becomes sum(c(1, 0))
  - Counts the number of values of TRUE in a vector

## Data frames

| rows | Individual observations / records |
|---|---|
| columns | Variables |

| read_csv() | <ul><li>Load data in R<ul><li>the resulting object type is called a "tibble"</li></ul></li></ul> |
|---|---|

| glimpse() | <ul><li>Tells how many rows and columns there are</li><li>Listing out the column names</li><li>Tells what their data type is</li><li>Giving a peak at the first few values</li></ul> |
|---|---|
| head() | <ul><li>Shows what the top couple rows of the data look like</li></ul> |
| Pipes %>% | <ul><li>Make it easy to apply functions to our data, step-by-step</li></ul> |

| | Nominal variable | Unordered descriptions (ex: turtle, butterfly, snail) |
|---|---|---|
| **Categorical Variables** | Ordinal variable | Ordered descriptions (ex: unhappy, ok, awesome) |
| | Binary variable | Only 2 mutually exclusive outcomes (ex: yes, no) |
| **Quantitative (numerical) Variables** | Continuous variable | Measured data, can have infinite values within possible range (ex: 3.2 inch, 34.16g) |
| | Discrete variable | Observations can only exist at limited values, often counts (ex: 7, 5, 9) |

**[Visualizing and describing the distribution of a quantitative(numerical) variable]**

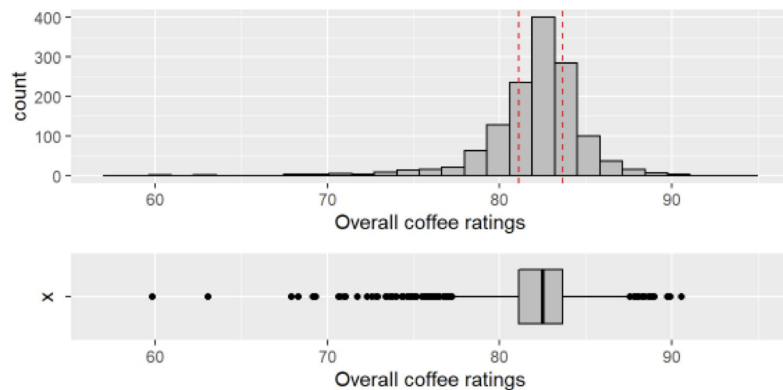| Histogram | |
|---|---|
| | <ul><li>**Height** of each **bar** counts the # of values of the variable which fall within the corresponding **bin**</li><li>**Horizontal axis**: numerical (no gaps between bins)</li><li>**Vertical axis**: the number of values which fall within each bin</li></ul> |
| **Features** | <ul><li>**Shape**: overall pattern of the values of the variable</li><li>**Center**: describes a 'typical' value of the variable</li><li>**Spread**: describes how concentrated the values of the variables are</li></ul> |
| **Shape** | **Skewness**<ul><li>Symmetric (not skewed)</li><li>Left-skewed</li><li>Right-skewed</li></ul>**Modality**<ul><li>Unimodal</li><li>Bimodal</li><li>Multimodal</li><li>Uniform</li></ul> |
| **Creating a histogram in R** | <ul><li>Replace **blue part**</li></ul>ggplot(data = **data_name**, aes(x = **numerical variable**)) + geom_histogram(color = "**black**", |

| | |
|---|---|
| | fill = "**gray**",<br>bins = **#**) +<br>labs(x = "**horizontal axis name**") |
| **mean / median / modality** | <ul><li>Left skewed<ul><li>mean < median < mode</li></ul></li><li>Right skewed<ul><li>mode < median < mean</li></ul></li><li>Symmetric<ul><li>mode $\approx$ median $\approx$ mean</li></ul></li></ul> |

| Center | | Spread | |
|---|---|---|---|
| Mean | $$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} = \frac{1}{n}\sum_{i=1}^{n} x_i$$ | Variance | $$s^2 = \frac{(x_1 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n-1} = \frac{1}{n-1}\sum_{i=1}^{n}(x_i - \bar{x})^2$$<br>Average squared distance between the values and their mean |
| Median | Half of the values are smaller/larger than the median | Standard Deviation | $$s = \sqrt{s^2}$$ |

| Boxplot | |
|---|---|
| | <ul><li>Summarizes the distribution of a quantitative variable using five statistics</li><li>Plot unusual observations (outliers)</li></ul> |
| **Features of the plot** | <ul><li>**Median**: line in the middle of the box</li><li>**First quartile (Q$_1$)**: one quarter (25%) of the data values are smaller than it</li><li>**Third quartile (Q$_3$)**: three quarters (75%) of the data values are smaller than it</li><li>**Inter-quartile range (IQR)**: $Q_3$ - $Q_1$, 50% of the values</li><li>Whiskers extend within 1.5 x IQR</li><li>**Outliers**: points beyond the whiskers</li></ul> |
| **Creating a boxplot in R** | <ul><li>Replace **blue part**</li></ul>ggplot(data = **data_name**, aes(x = "**"**,<br><br>    y = **numerical variable**)) +<br><br> geom_boxplot(color = "**black**", fill = "**gray**") +<br><br> labs(y = "**vertical axis name**")<br><br><ul><li>When **comparing distributions** across different categories</li></ul>ggplot(data = **data_name**, aes(x = **categorical variable,**<br><br>    y = **numerical variable**)) +<br><br> geom_boxplot(color = "**black**", fill = "**gray**") +<br><br>labs(y = "**vertical axis name**",<br><br> X = "**horizontal axis name**") |

**[Histogram vs. boxplot for the same distribution]**

| Barplots | |
|---|---|
| | • One bar for each category<br>• Height of a bar →# of values of the variable which fall within the corresponding category<br>• Arbitrary widths (should all be the same)<br>• Gap between each bar<br>• Arbitrary order of the bars<br>• **Shape or center don't make sense** |
| **Creating a barplot in R** | ggplot(data = **data_name**, aes(x = **categorical variable**)) +<br>  geom_bar(color = "**black**",<br>        fill = "**gray**") +<br>labs(x = "**horizontal axis name**")<br><br>• **Flipped version**<br>ggplot(data = **data_name**, aes(x = **categorical variable**)) +<br>  geom_bar(color = "**black**",<br>        fill = "**gray**") +<br>labs(x = "**horizontal axis name**") +<br>coord_flip() |

**[Tidy Data]**

Dataset is tidy because:

☐ Each value has its own cell

☐ Each observation has its own row

☐ Each variable has its own column

| Tidy data examples | Not Tidy data examples |
|---|---|