

Week 2: *Data Wrangling*

III EMSE 6035: Marketing Analytics for Design Decisions

👤 John Paul Helveston

📅 September 08, 2020

Required Packages (check `notes.R` file)

Make sure you have these libraries installed:

```
install.packages(c("tidyverse", "here"))
```

Remember: you only need to install packages once!

Once installed, you'll need to *load* the libraries every time you open RStudio:

```
library(tidyverse)  
library(here)
```

Week 2: *Data Wrangling*

1. Working with data frames
2. Data wrangling with the *tidyverse*

BREAK

3. Project proposals

Week 2: *Data Wrangling*

1. Working with data frames
2. Data wrangling with the *tidyverse*

BREAK

3. Project proposals

The data frame...in Excel

A screenshot of Microsoft Excel showing a data frame for the Beatles. The data is organized into columns A through H. Column A contains row numbers from 1 to 15. Columns B through E contain specific data points for each member: firstName, lastName, instrument, yearOfBirth, and deceased status.

	A	B	C	D	E	F	G	H
1	firstName	lastName	instrument	yearOfBirth	deceased			
2	John	Lennon	guitar	1940	TRUE			
3	Paul	McCartney	bass	1942	FALSE			
4	Ringo	Starr	drums	1940	FALSE			
5	George	Harrison	guitar	1943	TRUE			
6								
7								
8								
9								
10								
11								
12								
13								
14								
15								

The data frame...in R

```
beatles <- tibble(  
  firstName = c("John", "Paul", "Ringo", "George"),  
  lastName = c("Lennon", "McCartney", "Starr", "Harrison"),  
  instrument = c("guitar", "bass", "drums", "guitar"),  
  yearOfBirth = c(1940, 1942, 1940, 1943),  
  deceased = c(TRUE, FALSE, FALSE, TRUE)  
)  
  
beatles
```

```
#> # A tibble: 4 x 5  
#>   firstName lastName instrument yearOfBirth deceased  
#>   <chr>     <chr>      <chr>          <dbl> <lgl>  
#> 1 John       Lennon      guitar        1940  TRUE  
#> 2 Paul       McCartney bass           1942 FALSE  
#> 3 Ringo      Starr       drums        1940 FALSE  
#> 4 George     Harrison    guitar        1943  TRUE
```

Columns: Vectors of values (must be same data type)

```
beatles
```

```
#> # A tibble: 4 x 5
#>   firstName lastName instrument yearOfBirth deceased
#>   <chr>     <chr>    <chr>        <dbl> <lgl>
#> 1 John       Lennon    guitar        1940  TRUE
#> 2 Paul       McCartney bass         1942 FALSE
#> 3 Ringo      Starr     drums        1940 FALSE
#> 4 George     Harrison  guitar        1943  TRUE
```

Extract a column using \$

```
beatles$firstName
```

```
#> [1] "John"    "Paul"    "Ringo"   "George"
```

Rows: Information about individual observations

Information about *John Lennon* is in the first row:

```
beatles[1,]
```

```
#> # A tibble: 1 x 5
#>   firstName lastName instrument yearOfBirth deceased
#>   <chr>     <chr>      <chr>        <dbl> <lgl>
#> 1 John       Lennon     guitar        1940  TRUE
```

Information about *Paul McCartney* is in the second row:

```
beatles[2,]
```

```
#> # A tibble: 1 x 5
#>   firstName lastName instrument yearOfBirth deceased
#>   <chr>     <chr>      <chr>        <dbl> <lgl>
#> 1 Paul       McCartney bass        1942 FALSE
```

Take a look at the `beatles` data frame in `notes.R`

Getting data into R

1. Load external packages
2. Read in external files (usually a `.csv*` file)

*csv = "comma-separated values"

Data from an R package

```
library(ggplot2)
```

See which data frames are available in a package:

```
data(package = "ggplot2")
```

Find out more about a package data set:

```
?msleep
```

Back to **notes.R**

Importing an external data file

Note the `data.csv` file in your `data` folder.

- **DO NOT** double-click it!
- **DO NOT** open it in Excel!

Excel can **corrupt** your data!

If you **must** open it in Excel:

- Make a copy
- Open the copy

Steps to importing external data files

1. Create a path to the data

```
library(here)
path_to_data <- here('data', 'data.csv')
path_to_data
```

```
#> [1] "C:/Users/jph/Downloads/gh/2021-Fall/class/2-data-wrangling/data/data.csv"
```

2. Import the data

```
library(tidyverse)
data <- read_csv(path_to_data)
```

Using the **here** package to make file paths

The `here()` function builds the path to your **root** to your *working directory* (this is where your `.Rproj` file lives!)

```
here()
```

```
#> [1] "C:/Users/jph/Downloads/gh/2021-Fall/class/2-data-wrangling"
```

The `here()` function builds the path to files *inside* your working directory

```
path_to_data <- here('data', 'data.csv')  
path_to_data
```

```
#> [1] "C:/Users/jph/Downloads/gh/2021-Fall/class/2-data-wrangling/data/data.csv"
```

Avoid hard-coding file paths!

(they can break on different computers)

```
path_to_data <- 'data/data.csv'  
path_to_data
```

```
#> [1] "data/data.csv"
```



Back to reading in data

```
path_to_data <- here('data', 'data.csv')  
data <- read_csv(path_to_data)
```

Important: Use `read_csv()` instead of `read.csv()`

10:00

Think-Pair-Share

- 1) Use the `here()` and `read_csv()` functions to load the `data.csv` file that is in the `data` folder. Name the data frame object `data`.
- 2) Use the `data` object to answer the following questions:
 - How many rows and columns are in the data frame?
 - What type of data is each column? (Just look, don't need to type out the answer)
 - Preview the different columns - what do you think this data is about? What might one row represent?
 - How many unique airlines are in the data frame?
 - What is the earliest and latest observation in the data frame?
 - What is the shortest and longest air time for any one flight in the data frame?

Week 2: *Data Wrangling*

1. Working with data frames
2. Data wrangling with the *tidyverse*

BREAK

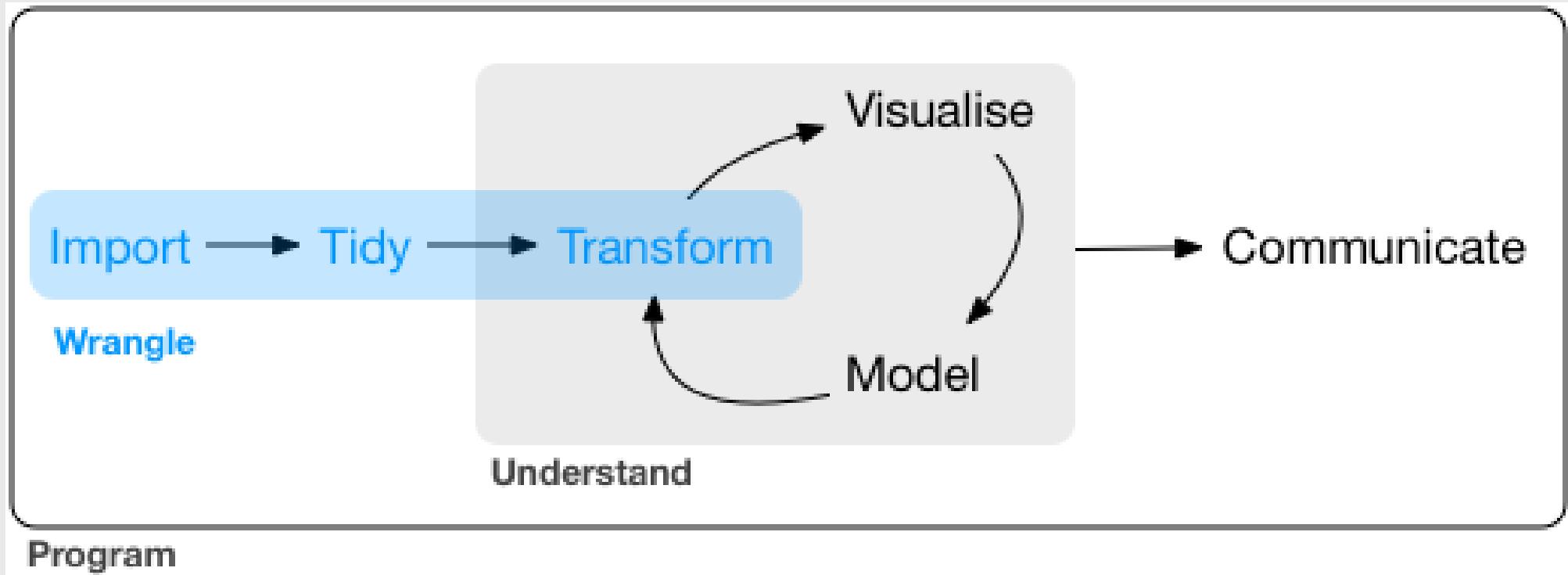
3. Project proposals

The tidyverse: `stringr` + `dplyr` + `readr` + `ggplot2` + ...



Art by Allison Horst

80% of the job is data wrangling



Today: data wrangling with **dplyr**



Art by [Allison Horst](#)

The main `dplyr` "verbs"

"Verb"	What it does
<code>select()</code>	Select columns by name
<code>filter()</code>	Keep rows that match criteria
<code>arrange()</code>	Sort rows based on column(s)
<code>mutate()</code>	Create new columns
<code>summarize()</code>	Create summary values

Core **tidyverse** concept: **Chain functions together with "pipes"**

%>%

Think of the words "...and then..."

```
data %>%  
  do_something() %>%  
  do_something_else()
```

Think of %>% as the words "...and then..."

Without Pipes (read from inside-out):

```
leave_house(get_dressed(get_out_of_bed(wake_up(me))))
```

With Pipes:

```
me %>%  
  wake_up %>%  
  get_out_of_bed %>%  
  get_dressed %>%  
  leave_house
```

Select columns with `select()`

Subset Variables (Columns)



Select columns with `select()`

```
beatles <- tibble(  
  firstName = c("John", "Paul", "Ringo", "George"),  
  lastName = c("Lennon", "McCartney", "Starr", "Harrison"),  
  instrument = c("guitar", "bass", "drums", "guitar"),  
  yearOfBirth = c(1940, 1942, 1940, 1943),  
  deceased = c(TRUE, FALSE, FALSE, TRUE)  
)  
  
beatles
```

```
#> # A tibble: 4 x 5  
#>   firstName lastName instrument yearOfBirth deceased  
#>   <chr>     <chr>      <chr>          <dbl> <lgl>  
#> 1 John       Lennon      guitar        1940  TRUE  
#> 2 Paul       McCartney  bass         1942  FALSE  
#> 3 Ringo      Starr       drums        1940  FALSE  
#> 4 George     Harrison    guitar        1943  TRUE
```

Select columns with `select()`

Select the columns `firstName` & `lastName`

```
beatles %>%  
  select(firstName, lastName)
```

```
#> # A tibble: 4 x 2  
#>   firstName lastName  
#>   <chr>     <chr>  
#> 1 John      Lennon  
#> 2 Paul      McCartney  
#> 3 Ringo    Starr  
#> 4 George   Harrison
```

Select columns with `select()`

Use the `-` sign to drop columns

```
beatles %>%  
  select(-firstName, -lastName)
```

```
#> # A tibble: 4 x 3  
#>   instrument yearOfBirth deceased  
#>   <chr>        <dbl> <lgl>  
#> 1 guitar       1940  TRUE  
#> 2 bass         1942 FALSE  
#> 3 drums        1940 FALSE  
#> 4 guitar       1943 TRUE
```

Select columns with `select()`

Select columns based on name criteria:

- `ends_with()` = Select columns that end with a character string
- `contains()` = Select columns that contain a character string
- `matches()` = Select columns that match a regular expression
- `one_of()` = Select column names that are from a group of names

Select columns with `select()`

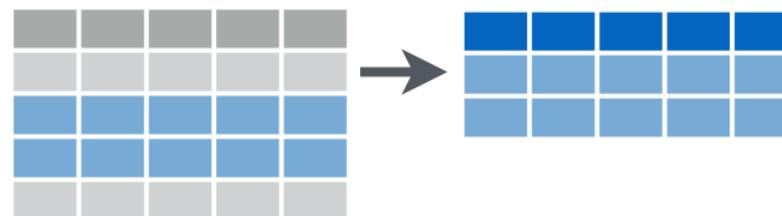
Select the columns that end with "Name":

```
beatles %>%  
  select(ends_with("Name"))
```

```
#> # A tibble: 4 x 2  
#>   firstName lastName  
#>   <chr>     <chr>  
#> 1 John      Lennon  
#> 2 Paul      McCartney  
#> 3 Ringo    Starr  
#> 4 George   Harrison
```

Select rows with `filter()`

Subset Observations (Rows)



Select rows with filter()

Select the band members born after 1941

```
#> # A tibble: 4 × 5
#>   firstName lastName instrument yearOfBirth deceased
#>   <chr>     <chr>    <chr>        <dbl> <lgl>
#> 1 John       Lennon    guitar      1940  TRUE
#> 2 Paul       McCartney bass       1942 FALSE
#> 3 Ringo     Starr     drums      1940 FALSE
#> 4 George    Harrison  guitar      1943  TRUE
```

Select rows with filter()

Select the band members born after 1941

```
beatles %>%  
  filter(yearOfBirth > 1941)
```

```
#> # A tibble: 2 x 5  
#>   firstName lastName instrument yearOfBirth deceased  
#>   <chr>     <chr>    <chr>          <dbl> <lgl>  
#> 1 Paul      McCartney bass            1942 FALSE  
#> 2 George    Harrison  guitar          1943 TRUE
```

Select rows with filter()

Select the band members born after 1941 & **are still living**

```
beatles %>%  
  filter(yearOfBirth > 1941, deceased == FALSE)
```

```
#> # A tibble: 1 x 5  
#>   firstName lastName instrument yearOfBirth deceased  
#>   <chr>     <chr>    <chr>          <dbl> <lgl>  
#> 1 Paul      McCartney bass            1942 FALSE
```

```
beatles %>%  
  filter(yearOfBirth > 1941 & deceased == FALSE)
```

```
#> # A tibble: 1 x 5  
#>   firstName lastName instrument yearOfBirth deceased  
#>   <chr>     <chr>    <chr>          <dbl> <lgl>  
#> 1 Paul      McCartney bass            1942 FALSE
```

Logic operators for `filter()`

Description	Example
Values greater than 1	<code>value > 1</code>
Values greater than or equal to 1	<code>value >= 1</code>
Values less than 1	<code>value < 1</code>
Values less than or equal to 1	<code>value <= 1</code>
Values equal to 1	<code>value == 1</code>
Values not equal to 1	<code>value != 1</code>
Values in the set c(1, 4)	<code>value %in% c(1, 4)</code>

Removing missing values

Drop all rows where `variable` is `NA`

```
data %>%  
  filter(!is.na(variable))
```

Combine `filter()` and `select()`

Get the **first & last name** of members born after 1941 & are still living

```
beatles %>%  
  filter(yearOfBirth > 1941, deceased == FALSE) %>%  
  select(firstName, lastName)
```

```
#> # A tibble: 1 x 2  
#>   firstName lastName  
#>   <chr>     <chr>  
#> 1 Paul       McCartney
```

10:00

Think-Pair-Share

- 1) Use the `here()` and `read_csv()` functions to load the `data.csv` file that is in the `data` folder. Name the data frame object `data`.
- 2) Use the `data` object and the `select()` and `filter()` functions to answer the following questions:
 - Create a new data frame, `flights_fall`, that contains only flights that departed in the fall semester.
 - Create a new data frame, `flights_dc`, that contains only flights that flew to DC airports (Reagan or Dulles).
 - Create a new data frame, `flights_dc_carrier`, that contains only flights that flew to DC airports (Reagan or Dulles) and only the columns about the month and airline.
 - How many unique airlines were flying to DC airports in July?

Create new variables with `mutate()`

Make New Variables





Create new variables with `mutate()`

Use the `yearOfBirth` variable to compute the age of each band member

```
beatles %>%  
  mutate(age = 2021 - yearOfBirth)
```

```
#> # A tibble: 4 x 6  
#>   firstName lastName instrument yearOfBirth deceased    age  
#>   <chr>     <chr>    <chr>          <dbl> <lgl>      <dbl>  
#> 1 John       Lennon    guitar        1940  TRUE       81  
#> 2 Paul       McCartney bass          1942 FALSE       79  
#> 3 Ringo      Starr     drums         1940 FALSE       81  
#> 4 George     Harrison  guitar        1943 TRUE       78
```

You can *immediately* use new variables

```
beatles %>%
  mutate(
    age = 2021 - yearOfBirth,
    meanAge = mean(age))
```

```
#> # A tibble: 4 x 7
#>   firstName lastName instrument yearOfBirth deceased   age meanAge
#>   <chr>     <chr>      <chr>        <dbl> <lgl>     <dbl>    <dbl>
#> 1 John       Lennon     guitar        1940  TRUE      81     79.8
#> 2 Paul       McCartney bass          1942 FALSE      79     79.8
#> 3 Ringo      Starr      drums         1940 FALSE      81     79.8
#> 4 George     Harrison   guitar        1943 TRUE      78     79.8
```

Handling if/else conditions

`ifelse(<condition>, <if TRUE>, <else>)`

```
beatles %>%  
  mutate(playsGuitar = ifelse(instrument == "guitar", TRUE, FALSE))
```

```
#> # A tibble: 4 x 6  
#>   firstName lastName instrument yearOfBirth deceased playsGuitar  
#>   <chr>     <chr>      <chr>        <dbl> <lgl>    <lgl>  
#> 1 John       Lennon     guitar        1940  TRUE     TRUE  
#> 2 Paul       McCartney bass          1942  FALSE    FALSE  
#> 3 Ringo      Starr      drums         1940  FALSE    FALSE  
#> 4 George     Harrison   guitar        1943  TRUE     TRUE
```

Sort data frame with `arrange()`

Sort `beatles` data frame by year of birth

```
beatles %>%  
  arrange(yearOfBirth)
```

```
#> # A tibble: 4 x 5  
#>   firstName lastName instrument yearOfBirth deceased  
#>   <chr>     <chr>    <chr>          <dbl> <lgl>  
#> 1 John       Lennon    guitar        1940  TRUE  
#> 2 Ringo      Starr     drums        1940 FALSE  
#> 3 Paul       McCartney bass        1942 FALSE  
#> 4 George     Harrison  guitar        1943  TRUE
```

Sort data frame with `arrange()`

Use the `desc()` function to sort in descending order

```
beatles %>%  
  arrange(desc(yearOfBirth))
```

```
#> # A tibble: 4 x 5  
#>   firstName lastName instrument yearOfBirth deceased  
#>   <chr>     <chr>    <chr>          <dbl> <lgl>  
#> 1 George    Harrison  guitar        1943  TRUE  
#> 2 Paul      McCartney bass         1942 FALSE  
#> 3 John      Lennon   guitar        1940  TRUE  
#> 4 Ringo    Starr    drums        1940 FALSE
```

Sort rows with `arrange()`

Compute the band member age, then sort based on the youngest:

```
beatles %>%  
  mutate(age = 2021 - yearOfBirth) %>%  
  arrange(age)
```

```
#> # A tibble: 4 x 6  
#>   firstName lastName instrument yearOfBirth deceased    age  
#>   <chr>     <chr>      <chr>        <dbl> <lgl>      <dbl>  
#> 1 George     Harrison    guitar       1943 TRUE        78  
#> 2 Paul       McCartney  bass        1942 FALSE       79  
#> 3 John       Lennon     guitar       1940 TRUE        81  
#> 4 Ringo     Starr      drums       1940 FALSE       81
```

10:00

Think pair share

- 1) Use the `here()` and `read_csv()` functions to load the `data.csv` file that is in the `data` folder. Name the data frame object `data`.
- 2) Using the `data` object, create the following new variables:
 - `speed`: The speed in mph, computed from the time (in minutes) and distance (in miles) variables.
 - `dep_delay_hour`: Is `TRUE` if the departure delay is greater or equal to one hour, `FALSE` otherwise.
- 3) Which flight flew the fastest?
- 4) Remove rows that have `NA` for `air_time` and re-arrange the resulting data frame based on the longest air time and longest flight distance.

Break

05 : 00

Week 2: *Data Wrangling*

1. Working with data frames
2. Data wrangling with the *tidyverse*

BREAK

3. Project proposals

Project Proposal Guidelines

Proposal Items

Item	Description
Abstract	Product / technology in just a few sentences
Introduction	Description, picture, background
Market Opportunity	Identify your customer, competitors, and market size
Product Attributes & Decision Variables	2-4 key variables related to product's design and performance
Questions	Major outstanding questions to be resolved

Today

Market Opportunity

- Identify customer
- Identify competitors
- Identify market size

Product Attributes

Features your *customer* cares about

Decision Variables

Features that *the designer* cares about

Example: **Folding solar panels**



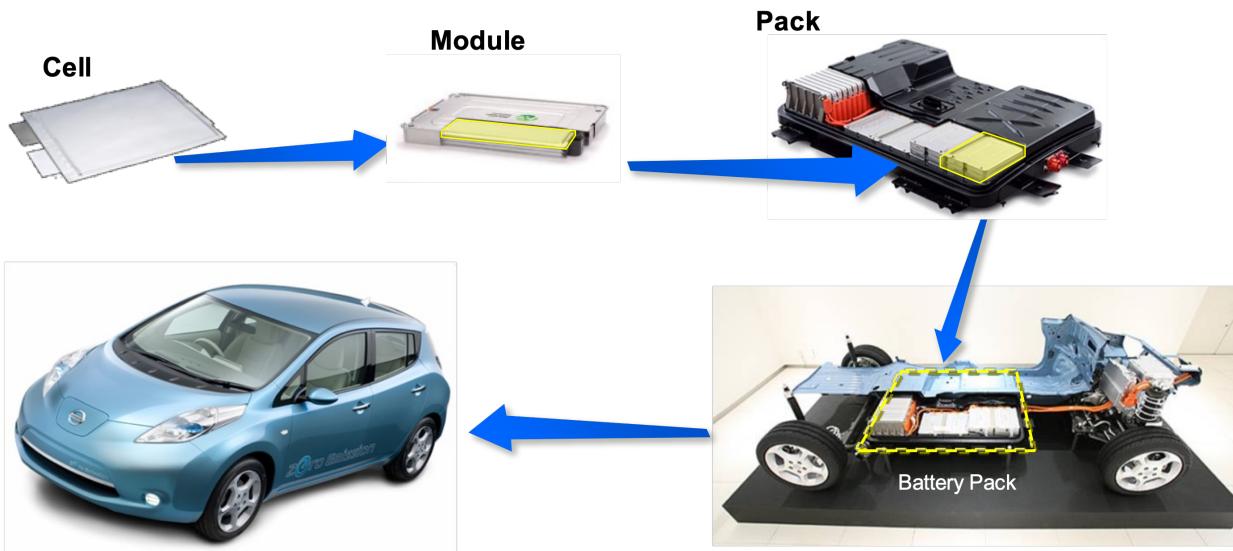
Who is your customer?

- General public?
- Outdoor enthusiasts?
- Emergency gear?

Competitors?

- Similar folding panels
- Batteries?

Example: Electric vehicle battery



Who is your customer?

- Car buyers

Competitors?

- Hybrid vehicles?
- Efficient gasoline vehicles?

Product Attributes

Features your *customer* cares about

Decision Variables

Features that *the designer* cares about

Product Diagram

Attribute Units

Price - USD

Weight - Kg

Power Output - Watts

Durability - Months

Portability - LxWxH

Decision Variable Units

Power Density - W/Kg

Degradation Rate - Hours

Packing Design - Cm³



Durability

Degradation Protections



Image Sources:

1. <https://www.deviantart.com/kata3214/art/Rainy-Sun-2444779413>
2. <https://www.pennscale.com/products/intec-certified-scales/label-printing-scales>
3. <https://www.thisiswhyimbroke.com/portfolio/foldable-solar-panel-charger/>

Model Relationships Table (example)

	Decision Variables			Demand	Competitors			
	Power Density	Degradation Rate	Packing Design		Aims Solar Panel	SUAOKI Solar Charger	Units	
Product Attributes	Price	-	-	+	-	225	160	USD
	Weight	-	-	+	-	2.6	2.06	kg
	Power Output	+	+	+	+	120	60	W
	Durability	-	+	-	+	60	12	Months
	Portability	-	-/+	+	+	20.6"x11"x 1.2"	11.5"x7.1"x2.9"	L"xW"xH"
Domain		[2.5, 60]	[24,1000]	[200, 2800]				
Units		W/kg	Hours	cm^3				

15 : 00

Team Proposals

1. Re-arrange tables to sit with your team
2. Discuss & identify your customer & potential competitors
3. Discuss & identify key *Product Attributes & Decision Variables*
4. Start building out your model relationships table (copy from [this example](#))

Suggestions

- You may want to start with simple bullet lists
- Start with more items rather than fewer (can always cut back later)