



THE OHIO STATE UNIVERSITY

FISHER COLLEGE OF BUSINESS

**BUSMGT 7331: Descriptive Analytics and Visualization**

Week 2

# More Data Wrangling, Covariation, and Customizing Visualization

Hyunwoo Park

Fisher College of Business

The Ohio State University

# Merge, String, Factor, Date

## Reading

- R4DS Chapter 10. Relational Data with dplyr
- R4DS Chapter 11. Strings with stringr
- R4DS Chapter 12. Factors withforcats
- R4DS Chapter 13. Dates and Times with lubridate

# Dataset for this week

- <https://www.kaggle.com/usdot/flight-delays>

kaggle Search  Competitions Datasets Kernels Discussion Learn ...  

  
Dataset  
2015 Flight Delays and Cancellations  
Which airline should you fly on to avoid significant delays?  
Department of Transportation • updated 2 years ago (Version 1)

Data Overview Kernels (92) Discussion (12) Activity Download (192 MB) New Kernel

Data (192 MB) API `kaggle datasets download -d usdot/flight-delays` ?  

Data Sources	About this file	Columns
<ul style="list-style-type: none"><li>airlines.csv 2 columns</li><li>airports.csv 7 columns</li><li>flights.csv 31 columns</li></ul>	IATA airline codes and names	 <ul style="list-style-type: none"><li>YEAR Year of the Flight Trip</li><li>MONTH Month of the Flight Trip</li><li>DAY Day of the Flight Trip</li></ul>

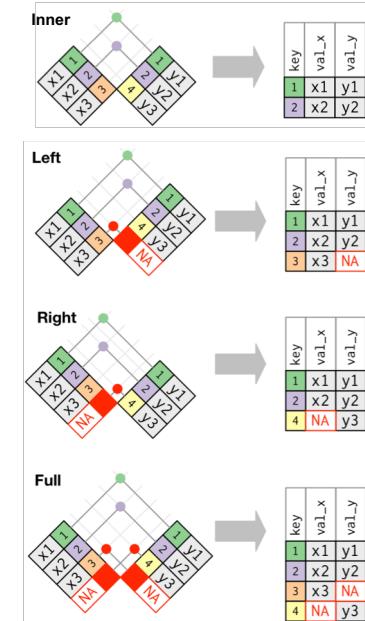
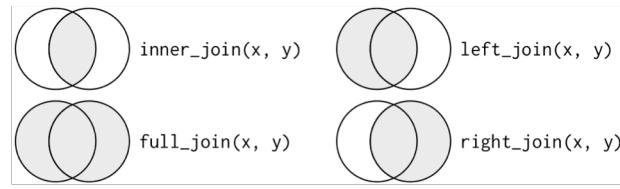
# Data import

```
1  airlines <- read_csv("data/airlines.csv")
2  glimpse(airlines)
3  airports <- read_csv("data/airports.csv")
4  glimpse(airports)
5  flights <- read_csv("data/flights.csv.zip")
6  glimpse(flights)
```

```
> glimpse(airlines)
Observations: 14
Variables: 2
$ IATA_CODE <chr> "UA", "AA", "US", "F9", "B6", "OO", "AS", "NK", "WN", "DL", "EV", "HA", "MQ"...
$ AIRLINE <chr> "United Air Lines Inc.", "American Airlines Inc.", "US Airways Inc.", "Front...
> glimpse(airports)
Observations: 322
Variables: 7
$ IATA_CODE <chr> "ABE", "ABI", "ABQ", "ABY", "ACK", "ACT", "ACV", "ACY", "ADK", "ADQ"...
$ AIRPORT <chr> "Lehigh Valley International Airport", "Abilene Regional Airport", "Albuer...
$ CITY <chr> "Allentown", "Abilene", "Albuquerque", "Aberdeen", "Albany", "Nantucket", "W...
$ STATE <chr> "PA", "TX", "NM", "SD", "GA", "MA", "TX", "CA", "NJ", "AK", "AK", "LA", "GA"...
$ COUNTRY <chr> "USA", "USA"...
$ LATITUDE <dbl> 40.65236, 32.41132, 35.04022, 45.44906, 31.53552, 41.25305, 31.61129, 40.978...
$ LONGITUDE <dbl> -75.44040, -99.68190, -106.60919, -98.42183, -84.19447, -70.06018, -97.23052...
> glimpse(flights)
Observations: 5,819,079
Variables: 31
$ YEAR           <int> 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, ...
$ MONTH          <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ...
$ DAY            <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ...
$ DAY_OF_WEEK    <int> 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, ...
$ AIRLINE        <chr> "AS", "AA", "US", "AA", "AS", "DL", "NK", "US", "AA", "DL", "DL", ...
$ FLIGHT_NUMBER  <int> 98, 2336, 840, 258, 135, 806, 612, 2013, 1112, 1173, 2336, 1674, 1...
$ TAIL_NUMBER    <chr> "N407AS", "N3KUAA", "N171US", "N3HYAA", "N527AS", "N373OB", "N635N...
$ ORIGIN_AIRPORT <chr> "ANC", "LAX", "LAX", "SEA", "SFO", "LAS", "LAX", "SFO", "LAX", ...
$ DESTINATION_AIRPORT <chr> "SEA", "PBI", "CLT", "MIA", "ANC", "MSP", "MSP", "CLT", "DFW", "AT...
$ SCHEDULED_DEPARTURE <time> 00:05:00, 00:10:00, 00:20:00, 00:25:00, 00:25:00, 00:25:...
$ DEPARTURE_TIME   <time> 23:54:00, 00:02:00, 00:18:00, 00:15:00, 00:24:00, 00:20:00, 00:19:...
$ DEPARTURE_DELAY  <int> -11, -8, -2, -5, -1, -5, -6, 14, -11, 3, -6, -8, 0, -6, -1, -4, -1...
$ TAXI_OUT         <int> 21, 12, 16, 15, 11, 18, 11, 13, 17, 12, 12, 21, 18, 18, 28, 17, 25...
$ WHEELS_OFF       <time> 00:15:00, 00:14:00, 00:34:00, 00:30:00, 00:35:00, 00:38:00, 00:30...
```

# Join two tables

- “join” means “merge” in tidyverse. The join command family merges two tibbles.
- Suppose we are joining two tables, x and y.
- Types of join
  - Mutating join* combines the variables from two tables.
    - [1] inner\_join(x, y)
    - Outer joins
      - [2] left\_join(x, y)
      - [3] right\_join(x, y)
      - [4] full\_join(x, y)
  - Filtering join* only affects the observations based on the match.
    - [1] semi\_join(x, y)
    - [2] anti\_join(x, y)



# Mutating joins and SQL syntax

---

dplyr join syntax	SQL syntax
<code>inner_join(x, y, by = "z")</code>	<code>SELECT * FROM x INNER JOIN y USING (z)</code>
<code>left_join(x, y, by = "z")</code>	<code>SELECT * FROM x LEFT OUTER JOIN y USING (z)</code>
<code>right_join(x, y, by = "z")</code>	<code>SELECT * FROM x RIGHT OUTER JOIN y USING (z)</code>
<code>full_join(x, y, by = "z")</code>	<code>SELECT * FROM x FULL OUTER JOIN y USING (z)</code>

---

- **by** argument
  - If you don't provide any **by** argument, dplyr will use the combination of all common variables (i.e., having the same variable name) as key.
  - If you supply a set of strings (e.g., `by="AIRLINE"`), dplyr use only those variables to merge.
  - If you want to match based on variables with different names in two data, use `by=c("COL_IN_X"="COL_IN_Y")`.

# Keys

- Key is a variable that uniquely identifies an observation.
  - You can create a key by combining multiple variables when they collectively identifies a unique observation.
- 
- *Primary key* is a key in its own table.
  - *Foreign key* is a key in another table.
  - *Surrogate key* is an artificially generated key with `mutate()` and `row_number()`.

# Let's join flights with airlines and airports.

- To save time, first count flights by airline and airport.
- We can do this multiple ways.

---

```
1 flights %>% count(AIRLINE, ORIGIN_AIRPORT)
2 flights %>% count(AIRLINE, ORIGIN_AIRPORT, sort=T)
3 flights %>% group_by(AIRLINE, ORIGIN_AIRPORT) %>% tally() %>% ungroup()
4 flights %>% group_by(AIRLINE, ORIGIN_AIRPORT) %>% tally(sort=T) %>% ungroup()
5 crosstab <- flights %>% count(AIRLINE, ORIGIN_AIRPORT, sort=T)
```

---

- The dplyr function `count` is equivalent to `group_by %>% tally %>% ungroup`.

```
# A tibble: 2,386 x 3
  AIRLINE ORIGIN_AIRPORT     n
  <chr>   <chr>       <int>
1 DL      ATL        221705
2 AA      DFW        134270
3 WN      MDW        76350
4 WN      LAS        68520
5 WN      BWI        64063
6 MQ      ORD        63130
7 UA      ORD        59538
8 EV      IAH        58330
9 EV      ATL        56959
10 DL     MSP        55929
# ... with 2,376 more rows
```

# Mutating join: inner\_join

```
1 crosstab %>%
2   inner_join(airlines, by=c("AIRLINE"="IATA_CODE")) %>%
3   inner_join(airports, by=c("ORIGIN_AIRPORT"="IATA_CODE"))
```

	AIRLINE	ORIGIN_AIRPORT	n AIRLINE.y
	<chr>	<chr>	<int> <chr>
1	DL	ATL	221705 Delta Air Lines Inc.
2	AA	DFW	134270 American Airlines Inc.
3	WN	MDW	76350 Southwest Airlines Co.
4	WN	LAS	68520 Southwest Airlines Co.
5	WN	BWI	64063 Southwest Airlines Co.
6	MQ	ORD	63130 American Eagle Airlines Inc.
7	UA	ORD	59538 United Air Lines Inc.
8	EV	IAH	58330 Atlantic Southeast Airlines
9	EV	ATL	56959 Atlantic Southeast Airlines
10	DL	MSP	55929 Delta Air Lines Inc.
# ... with 2,376 more rows			

	AIRLINE	ORIGIN_AIRPORT	n AIRLINE.y	AIRPORT	CITY	STATE	COUNTRY	LATITUDE	LONGITUDE
	<chr>	<chr>	<int> <chr>	<chr>	<chr>	<chr>	<chr>	<dbl>	<dbl>
1	DL	ATL	221705 Delta Air L... Hartsfield-... Atlan... GA	USA	33.6	-84.4			
2	AA	DFW	134270 American Ai... Dallas/Fort... Dalla... TX	USA	32.9	-97.0			
3	WN	MDW	76350 Southwest A... Chicago Mid... Chica... IL	USA	41.8	-87.8			
4	WN	LAS	68520 Southwest A... McCarran In... Las V... NV	USA	36.1	-115.			
5	WN	BWI	64063 Southwest A... Baltimore-W... Balti... MD	USA	39.2	-76.7			
6	MQ	ORD	63130 American Ea... Chicago O'H... Chica... IL	USA	42.0	-87.9			
7	UA	ORD	59538 United Air ... Chicago O'H... Chica... IL	USA	42.0	-87.9			
8	EV	IAH	58330 Atlantic So... George Bush... Houst... TX	USA	30.0	-95.3			
9	EV	ATL	56959 Atlantic So... Hartsfield-... Atlan... GA	USA	33.6	-84.4			
10	DL	MSP	55929 Delta Air L... Minneapolis... Minne... MN	USA	44.9	-93.2			
# ... with 1,294 more rows									

- Try other mutating joins (left\_join, right\_join, full\_join) and see if you get what you expect to get. How are they different? Why?

# Filtering join: semi\_join

```
1 crosstab %>%
2   semi_join(airlines, by=c("AIRLINE"="IATA_CODE")) %>%
3   semi_join(airports, by=c("ORIGIN_AIRPORT"="IATA_CODE"))
```

# A tibble: 2,386 x 3			
	AIRLINE	ORIGIN_AIRPORT	n
# <chr>	<chr>	<int>	
1	DL	ATL	221705
2	AA	DFW	134270
3	WN	MDW	76350
4	WN	LAS	68520
5	WN	BWI	64063
6	MQ	ORD	63130
7	UA	ORD	59538
8	EV	IAH	58330
9	EV	ATL	56959
10	DL	MSP	55929
# ... with 2,376 more rows			

# A tibble: 1,304 x 3			
	AIRLINE	ORIGIN_AIRPORT	n
# <chr>	<chr>	<int>	
1	DL	ATL	221705
2	AA	DFW	134270
3	WN	MDW	76350
4	WN	LAS	68520
5	WN	BWI	64063
6	MQ	ORD	63130
7	UA	ORD	59538
8	EV	IAH	58330
9	EV	ATL	56959
10	DL	MSP	55929
# ... with 1,294 more rows			

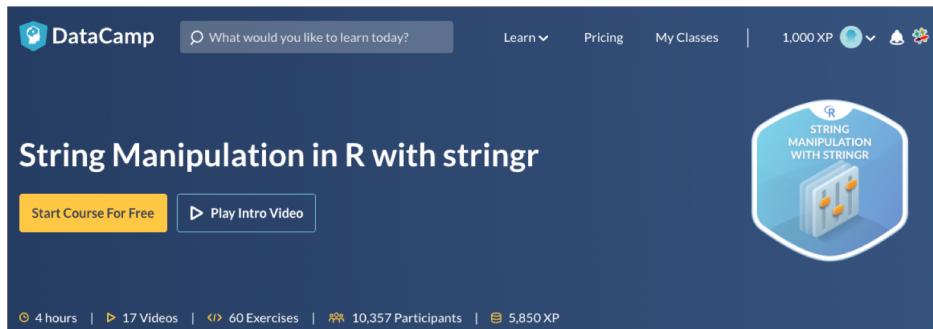
- Try the other filtering join (anti\_join) and see if you get what you expect to get. How are they different? Why?

# Basic string manipulation with stringr

- Make sure to load stringr package by `library(stringr)`.  
It's not part of core tidyverse because not all data analysis involves textual data.
- All stringr functions start with `str_`.
- Common (easier) functions
  - `str_length` counts the number of characters in a string.
  - `str_c` combines multiple strings into one.
  - `str_sub` extracts a substring from a string.
  - `str_trim` removes whitespaces from start and end of a string.
  - `str_squish` replaces two or more consecutive whitespaces into one.
  - `str_to_lower` lowercases a string. (This is cool => this is cool)
  - `str_to_upper` uppercases a string. (This is cool => THIS IS COOL)
  - `str_to_title` titlecases a string. (This is cool => This Is Cool)

# More string manipulation with stringr

- To learn more advanced use of stringr functions, if involving regular expression in particular,
  - Read R4DS Chapter 11: Strings with stringr.  
<https://r4ds.had.co.nz/strings.html>
  - Complete DataCamp course: String Manipulation in R with stringr.  
<https://www.datacamp.com/courses/string-manipulation-in-r-with-stringr>



The screenshot shows the DataCamp website interface. At the top, there's a navigation bar with the DataCamp logo, a search bar, and links for 'Learn', 'Pricing', 'My Classes', and a user account section showing '1,000 XP'. Below the header, the course title 'String Manipulation in R with stringr' is displayed in large white text. Underneath the title are two buttons: 'Start Course For Free' and 'Play Intro Video'. To the right of the title is a circular icon containing a stylized 'R' and the text 'STRING MANIPULATION WITH STRINGR' over an illustration of a server or storage unit. At the bottom of the main course area, there are statistics: '4 hours', '17 Videos', '60 Exercises', '10,357 Participants', and '5,850 XP'. A 'Course Description' section follows, which is partially cut off at the bottom.

## Course Description

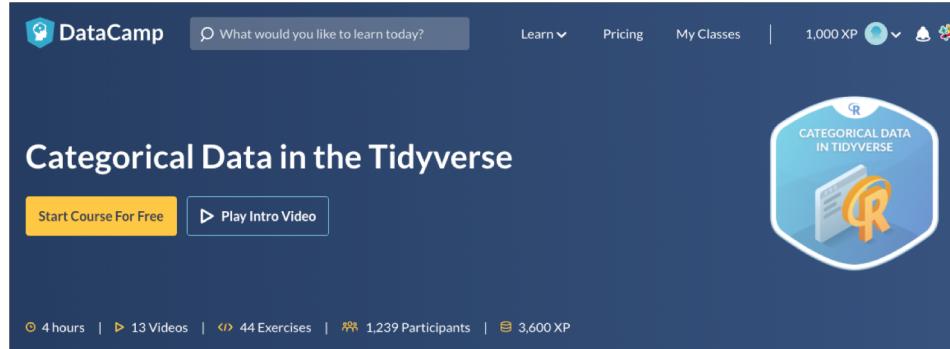
Character strings can turn up in all stages of a data science project. You might have to clean messy string input before analysis, extract data that is embedded in text or automatically turn numeric results into a sentence to include in a report. Perhaps the strings themselves are the data of interest, and you need to detect and match patterns within them. This course will help you master these tasks by teaching you how to pull strings apart, put them back together and use stringr to detect, extract, match and split strings using regular expressions, a powerful way to express patterns.

This course is part of these tracks:

R Programmer  
Text Mining with R

# Factor

- Make sure to loadforcats package by `library(forcats)`.
- Factor can be considered implicitly ordered categorical variable.
- It's essential to understand factor for customizing a visualization involving a categorical variable.
- To learn more,
  - Read R4DS Chapter 12: Factors withforcats. <https://r4ds.had.co.nz/factors.html>
  - Complete DataCamp course: Categorical Data in the Tidyverse  
<https://www.datacamp.com/courses/categorical-data-in-the-tidyverse>



The screenshot shows the DataCamp learning platform. At the top, there is a navigation bar with the DataCamp logo, a search bar containing "What would you like to learn today?", and links for "Learn", "Pricing", and "My Classes". A user icon indicates 1,000 XP and shows a notification bell with three dots. Below the navigation, the course title "Categorical Data in the Tidyverse" is displayed in large white text on a dark blue background. Underneath the title are two buttons: "Start Course For Free" and "Play Intro Video". To the right of the title is a hexagonal badge featuring the R logo and the text "CATEGORICAL DATA IN TIDYVERSE". At the bottom of the page, there is footer information: "4 hours | 13 Videos | 44 Exercises | 1,239 Participants | 3,600 XP".

# Date, time, datetime

- Make sure to load lubridate package by `library(lubridate)`.
- Date and time are numerical variables in nature, but their notations are rather unique and different from usual numerical variables.
- Understanding how to handle date and time is important for analyzing data involving time dimension.
- To learn more,
  - Read R4DS Chapter 13. Dates and Times with lubridate <https://r4ds.had.co.nz/dates-and-times.html>
  - Complete DataCamp course: Working with Dates and Times in R <https://www.datacamp.com/courses/working-with-dates-and-times-in-r>

The screenshot shows the DataCamp website interface. At the top, there's a navigation bar with the DataCamp logo, a search bar asking 'What would you like to learn today?', and links for 'Learn', 'Pricing', and 'My Classes'. A user icon indicates 1,000 XP. Below the header, the main title 'Working with Dates and Times in R' is displayed, along with two buttons: 'Start Course For Free' and 'Play Intro Video'. To the right of the title is a circular badge featuring an R logo and the text 'WORKING WITH DATES AND TIMES' over an illustration of a stack of books.

# Covariation

# Relationship between two variables

- Last week we learned how to describe the distribution of a single categorical/discrete/continuous variable.
  - A natural next question is about relationship between two variables.
- 
- *Variation* is how much measurements of a single variable vary across observations.
  - *Covariation* is how the variation of one variable is associated with the variation of another variable.

# Measures of variation and covariation

	Variation (single variable)	Covariation (two variables)
Unscaled	Variance $s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$	Covariance $q_{x,y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$
Normalized		Correlation $r_{x,y} = \frac{q_{x,y}}{s_x s_y}$ $= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$

# Hand calculation of variance, covariance, and correlation

```
1 a <- tibble(x=1:3, y=3:1)
2 a
3 c(var(a$x), var(a$y), cov(a$x, a$y), cor(a$x, a$y))
```

```
> a <- tibble(x=1:3, y=3:1)
> a
# A tibble: 3 x 2
  x     y
  <int> <int>
1 1     3
2 2     2
3 3     1
> c(var(a$x), var(a$y), cov(a$x, a$y), cor(a$x, a$y))
[1]  1  1 -1 -1
```

$$s_x^2 = s_y^2 = \frac{(1-2)^2 + (2-2)^2 + (3-1)^2}{3-1} = 1$$

$$q_{x,y} = \frac{(1-2)(3-2) + (2-2)(2-2) + (3-2)(2-3)}{3-1} = -1$$

$$r_{x,y} = \frac{q_{x,y}}{s_x s_y} = \frac{-1}{1 \times 1} = -1$$

# Different levels of correlation

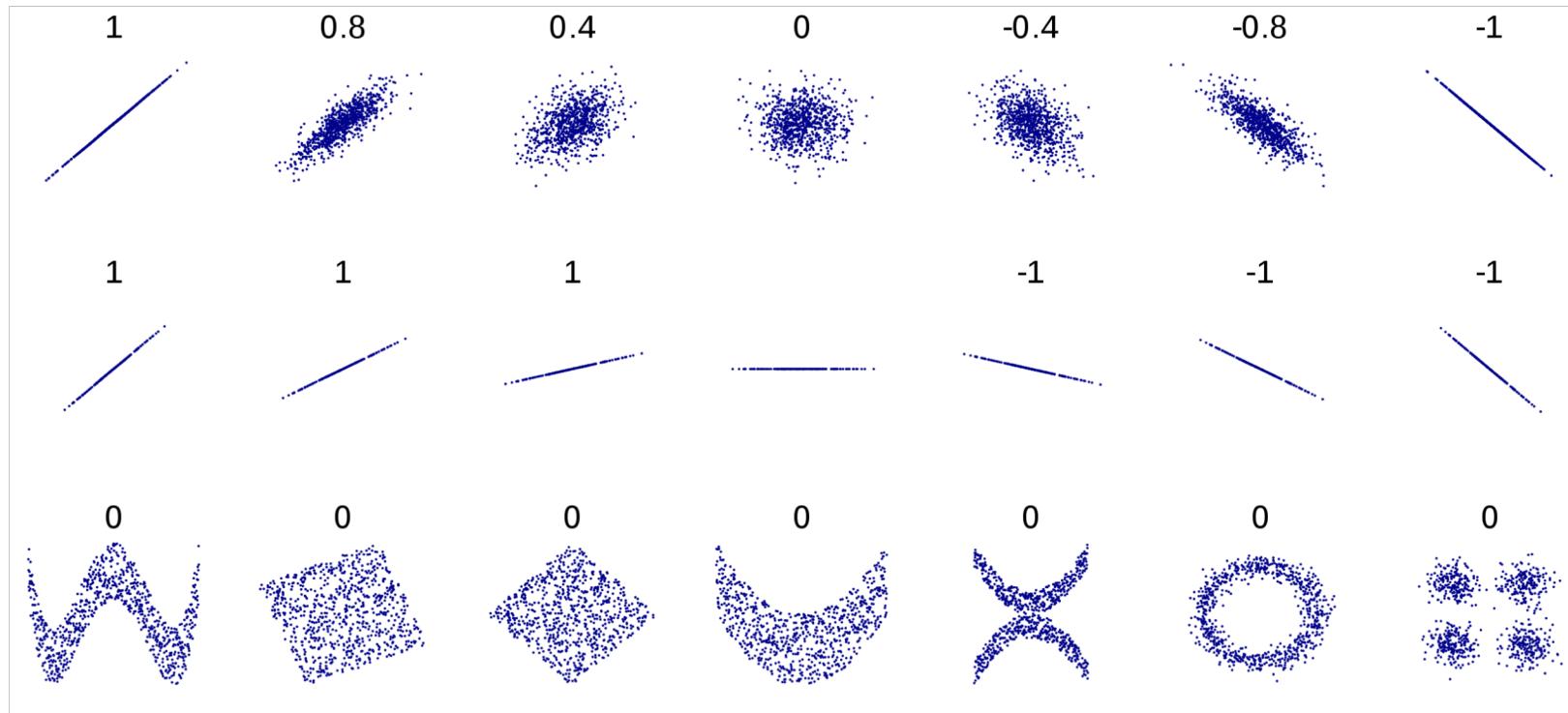


Figure from Wikipedia ([https://en.wikipedia.org/wiki/Correlation\\_and\\_dependence](https://en.wikipedia.org/wiki/Correlation_and_dependence))

# Anscombe's quartet

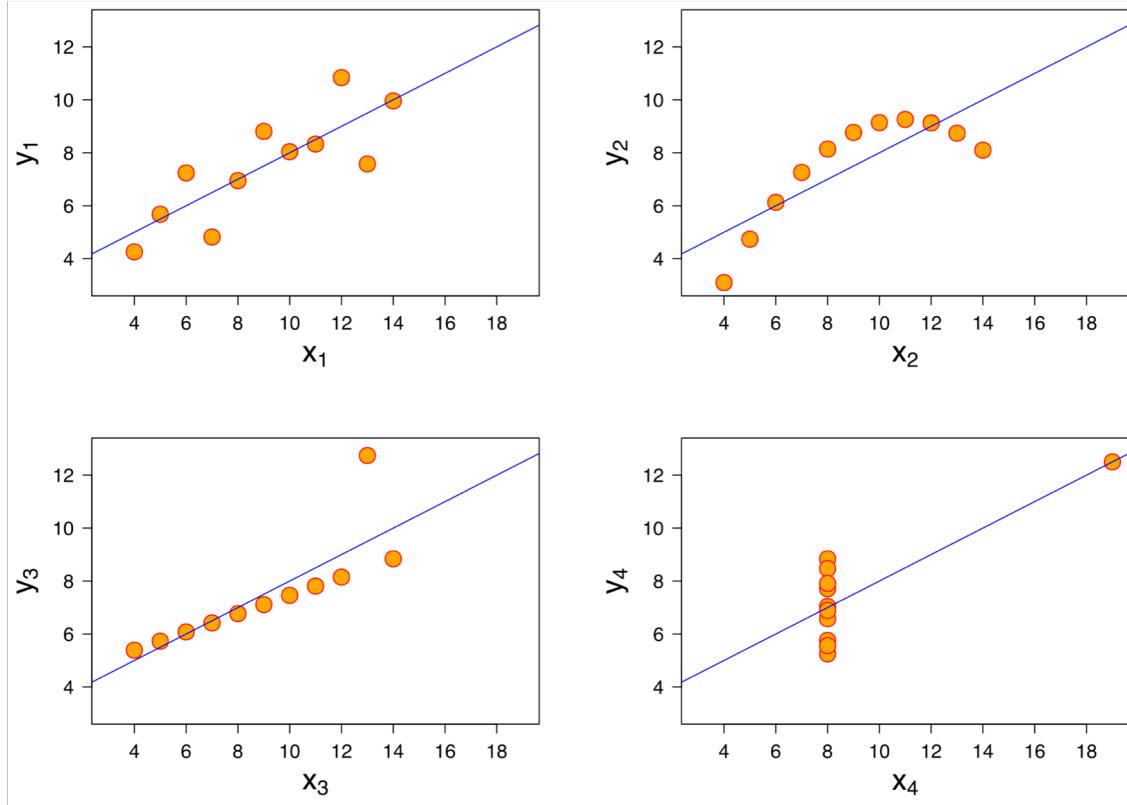


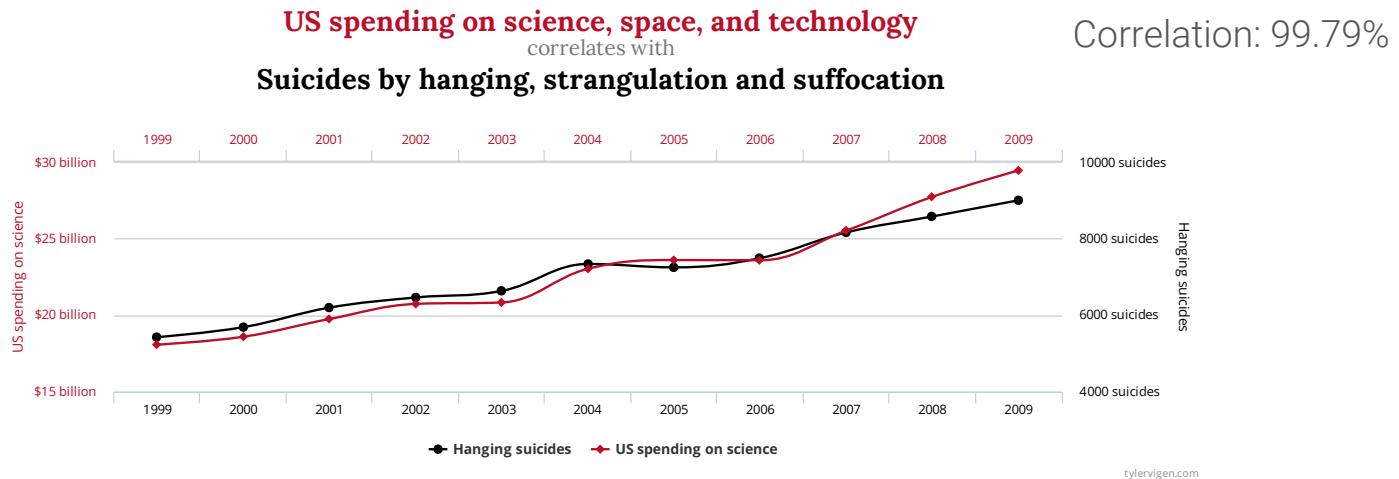
Figure from Wikipedia ([https://en.wikipedia.org/wiki/Anscombe%27s\\_quartet](https://en.wikipedia.org/wiki/Anscombe%27s_quartet))

# Correlation vs. causation

- Short reading on Wikipedia: [Correlation does not imply causation.](#)
- When one observes a strong correlation between A and B, one is strongly tempted to conclude “[1] A causes B”. It’s certainly a possibility. But, there are a few alternative scenarios that may give rise to a strong correlation between A and B.
  - [2] B causes A. (Reverse causality)
  - [3] C causes both A and B. (Common causal variable; omitted variable bias)
  - [4] A and B both cause C, and observations were only collected conditioned on C. (Selection bias)
  - [5] A causes B, and B also causes A. (Bidirectional / cyclic causation; imprecise magnitude of effect)
  - [6] A causes C which cases B. (Indirect causation; unknown mechanism)
  - [7] A and B have no connection. (Coincidence)

# An example of misled causal conclusion from correlation

- A fun collection of spurious correlations  
<http://tylervigen.com/spurious-correlations>

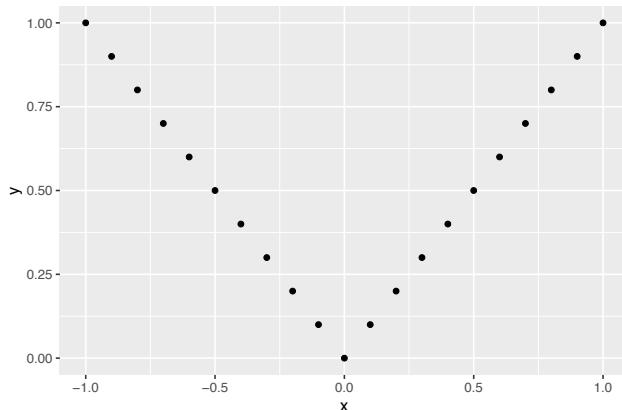


- Is US spending on SST "make" people suicide by hanging? Hell no! Probably not...

# Uncorrelated doesn't mean independence, either.

```
1 a <- tibble(x=seq(-1,1,.1)) %>% mutate(y=abs(x))
2 ggplot(a) + geom_point(aes(x,y))
3 a %>% summarize(cor(x,y))
4 cor(a$x,a$y)
```

Correlation: 0.0000000000000001081586



```
> a %>% summarize(cor(x,y))
# A tibble: 1 × 1
`cor(x, y)`
<dbl>
1 1.08e-16
> cor(a$x,a$y)
[1] 1.081586e-16
```

True relationship between x and y:  $y = |x|$   
In this case, x causes y.

# How to establish a causal relationship?

- It's really a deep (philosophical, really) problem, especially for social sciences, fundamentally because we cannot rewind history and relive the identically same world with only changes in variables of our interest.
- In order to establish causality, we need counterfactual (i.e., what if A were B and everything else remained the same).
  - Well-designed scientific experiments can achieve this counterfactual with randomization by randomly selecting groups from a population and giving a treatment to only one group.
  - In the business world, A/B testing is such an experiment.
  - Many interesting and important questions in social sciences do not lend themselves to such randomized experiments due to ethical and practical issues.
- Then, what to do with observational data?
  - To me, establishing causality from observational data analysis seems a matter of argumentation with sophisticated slicing-and-dicing of the data.
  - Also, I believe, establishing causality is not the only good goal of data analysis.
  - Descriptive analytics is important at least because it allows us to better understand what's actually happening in the world we live in.
  - Visualization is a great tool for communicating our understanding with each other.

# Combination of types of two variables

	Categorical	Discrete	Continuous
Categorical	[1] Count observations for each combination of levels	[2] Summarize the continuous variable by each category	
Discrete			
Continuous			[3] After binning, count observations for each combination of levels

# [1] Cross tabulation (two categorical variables)

- Count observations for each combination of levels in two categorical variables.
- If you put one variable in rows and the other variable in columns, it's called contingency table.

```
1 crosstab %>%
2   spread(AIRLINE, n) %>%
3   replace(is.na(.), 0) %>%
4   semi_join(airports, by=c("ORIGIN_AIRPORT"="IATA_CODE")) %>%
5   mutate(n=rowSums(.[2:ncol(.)])) %>%
6   arrange(-n)
```

	ORIGIN_AIRPORT	AA	AS	B6	DL	EV	F9	HA	MQ	NK	OO	UA	US	VX	WN	n
1	ATL	7498	429	0	221705	56959	4494	0	1295	4369	4230	3405	2939	0	39513	346836
2	ORD	50741	1748	2256	6962	44243	4611	0	63130	9965	37858	59538	3452	1380	0	285884
3	DFW	134270	1348	598	5095	24112	1288	0	53603	8519	4131	3505	3082	0	0	239551
4	DEN	7107	1589	958	7080	5507	21175	0	0	3691	44633	46218	2329	0	55768	196055
5	LAX	32738	8071	4108	25171	0	1382	1555	184	4344	36691	27429	3688	11801	37511	194673
6	SFO	12216	5177	4928	9701	0	1946	668	0	0	34994	45587	2615	15940	14236	148008
7	PHX	28214	2275	663	6918	8	1870	334	0	966	18552	5469	27781	0	53765	146815
8	IAH	6501	381	0	2927	58330	1529	0	27	5096	15381	53985	2465	0	0	146622
9	LAS	11359	4073	4064	10626	0	4327	811	0	9037	2421	11153	3024	3766	68520	133181
10	MSP	5487	587	0	55929	6840	1513	0	187	3495	25513	3208	2086	0	7272	112117
	# ... with 312 more rows															

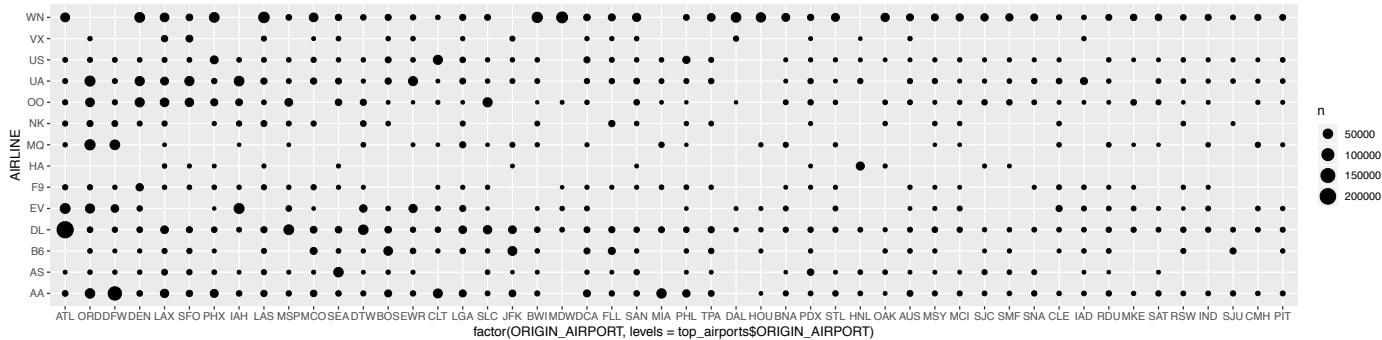
# [1] Visualizing crosstab with tile chart (heatmap)

```
1 top_airports <- crosstab %>%
2   semi_join(airports, by=c("ORIGIN_AIRPORT"="IATA_CODE")) %>%
3   group_by(ORIGIN_AIRPORT) %>%
4   summarize(n=sum(n)) %>% arrange(-n) %>% head(50)
5
6 filtered_crosstab <- crosstab %>%
7   semi_join(top_airports, by="ORIGIN_AIRPORT")
8
9 ggplot(filtered_crosstab) +
10   geom_tile(aes(x=AIRLINE, y=factor(ORIGIN_AIRPORT,
11                 levels=top_airports$ORIGIN_AIRPORT), fill=n)) +
12   coord_flip()
```



# [1] Visualizing crosstab with bubble chart

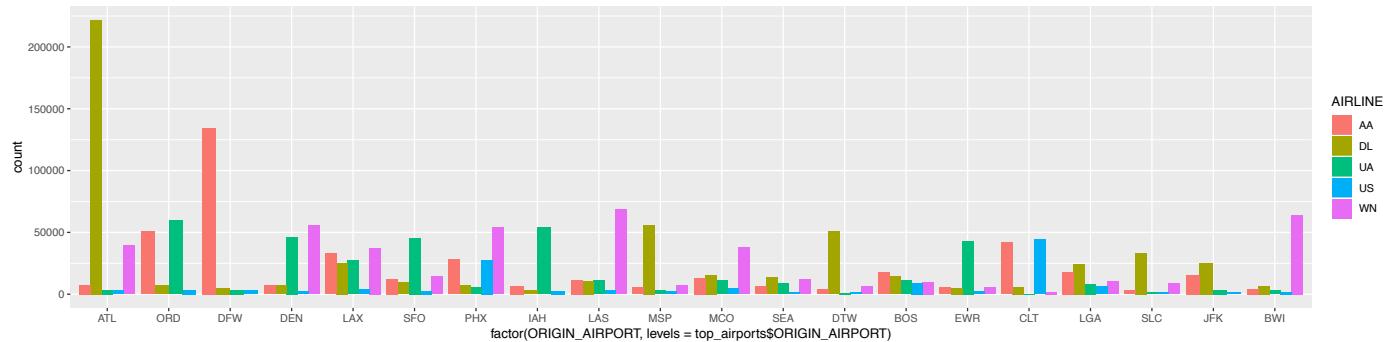
```
1 ggplot(filtered_crosstab) +  
2   geom_point(aes(x=AIRLINE, y=factor(ORIGIN_AIRPORT,  
3 levels=top_airports$ORIGIN_AIRPORT), size=n)) +  
4   coord_flip()
```



- Just like bar chart shows distribution of a categorical variable, these charts are like a 3-D bar chart looked down from above.

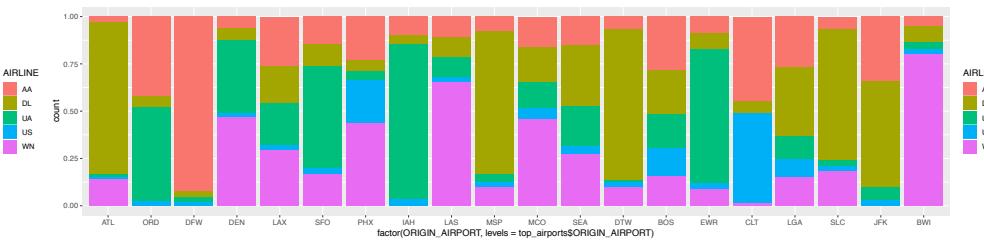
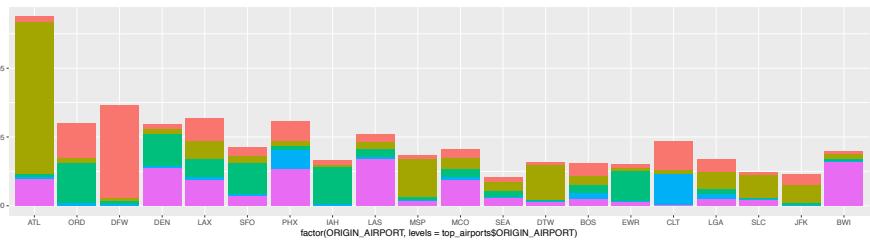
# [1] Grouped bar chart

```
1 top_airports <- crosstab %>%
2   semi_join(airports, by=c("ORIGIN_AIRPORT"="IATA_CODE")) %>%
3   group_by(ORIGIN_AIRPORT) %>%
4   summarize(n=sum(n)) %>% arrange(-n) %>% head(20)
5 filtered_flights <- flights %>%
6   semi_join(top_airports, by="ORIGIN_AIRPORT") %>%
7   filter(grepl("AA|DL|WN|UA|US", AIRLINE))
8 ggplot(filtered_flights) +
9   geom_bar(aes(
10     x=factor(ORIGIN_AIRPORT, levels=top_airports$ORIGIN_AIRPORT),
11     fill=AIRLINE),
12     position="dodge")
```



# [1] Stacked bar chart & proportion stacked bar chart

```
1 ggplot(filtered_flights) +  
2   geom_bar(aes(  
3     x=factor(ORIGIN_AIRPORT, levels=top_airports$ORIGIN_AIRPORT),  
4     fill=AIRLINE), position="stack")  
5  
6 ggplot(filtered_flights) +  
7   geom_bar(aes(  
8     x=factor(ORIGIN_AIRPORT, levels=top_airports$ORIGIN_AIRPORT),  
9     fill=AIRLINE), position="fill")
```



## [2] Summarize by group (categorical + continuous variables)

- The true power of `summarize` is realized when used with the `group_by`.
- The typical syntax is:
  - `df %>%  
 group_by(<VARIABLES>) %>%  
 summarize(<NEW_SUMMARY_VARIABLE> = <STAT_FUNCTION>(), count = n())`

---

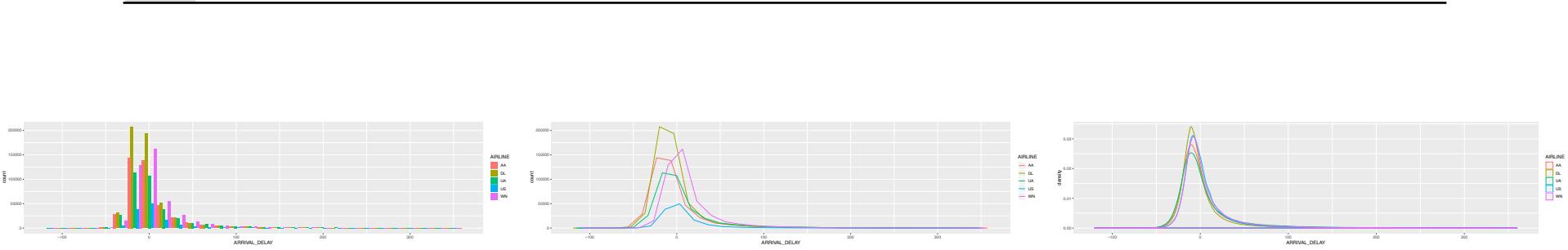
```
1 top_airports <- crosstab %>%  
2   semi_join(airports, by=c("ORIGIN_AIRPORT"="IATA_CODE")) %>%  
3   group_by(ORIGIN_AIRPORT) %>%  
4   summarize(n=sum(n)) %>% arrange(-n) %>% head(50)
```

---

```
> top_airports  
# A tibble: 50 x 2  
  ORIGIN_AIRPORT     n  
  <chr>        <int>  
1 ATL            346836  
2 ORD            285884  
3 DFW            239551  
4 DEN            196055  
5 LAX            194673  
6 SFO            148008  
7 PHX            146815  
8 IAH            146622  
9 LAS            133181  
10 MSP           112117  
# ... with 40 more rows
```

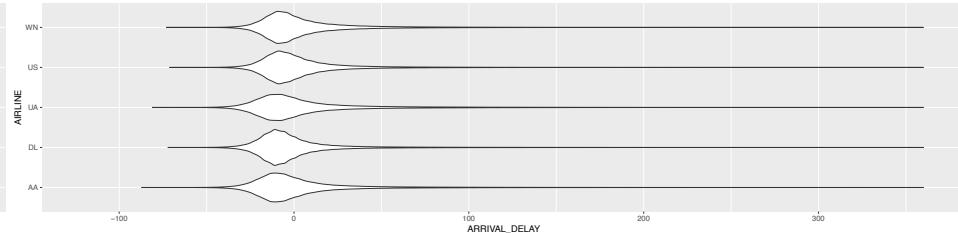
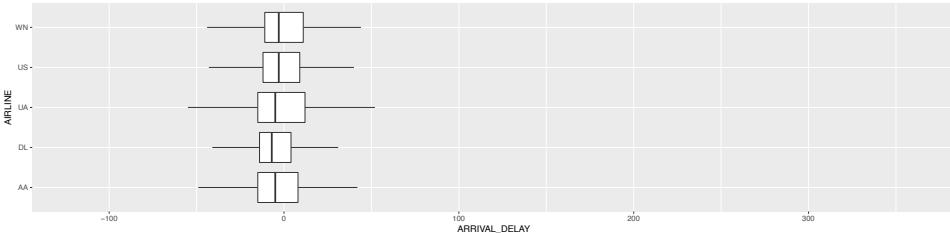
## [2] Grouped histogram & line plot & density plot

```
1 ggplot(filtered_flights) +  
2   geom_histogram(aes(x=ARRIVAL_DELAY, fill=AIRLINE), position="dodge") +  
3   xlim(c(-120,360))  
4 ggplot(filtered_flights) +  
5   geom_freqpoly(aes(x=ARRIVAL_DELAY, color=AIRLINE), position="dodge") +  
6   xlim(c(-120,360))  
7 ggplot(filtered_flights) +  
8   geom_density(aes(x=ARRIVAL_DELAY, color=AIRLINE), position="dodge") +  
9   xlim(c(-120,360))
```



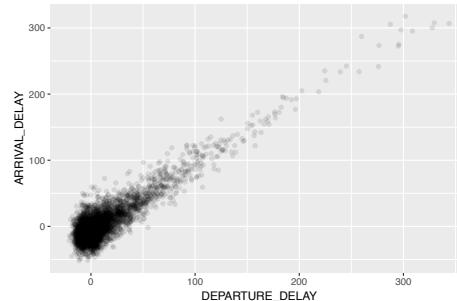
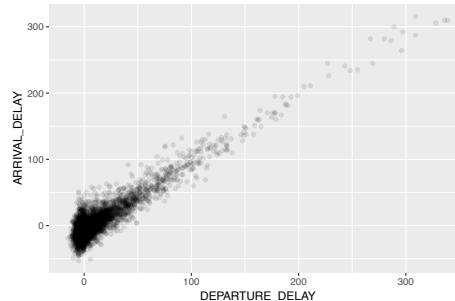
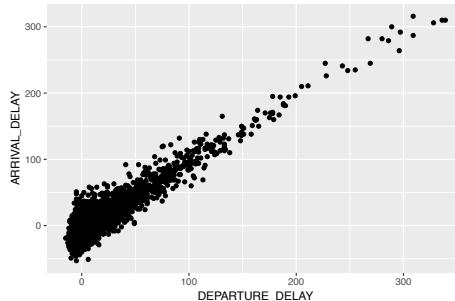
## [2] Grouped box plot & violin plot

```
1 ggplot(filtered_flights) +  
2   geom_boxplot(aes(x=AIRLINE, y=ARRIVAL_DELAY), outlier.shape=NA) +  
3   ylim(c(-120,360)) + coord_flip()  
4  
5 ggplot(filtered_flights) +  
6   geom_violin(aes(x=AIRLINE, y=ARRIVAL_DELAY)) +  
7   ylim(c(-120,360)) + coord_flip()
```



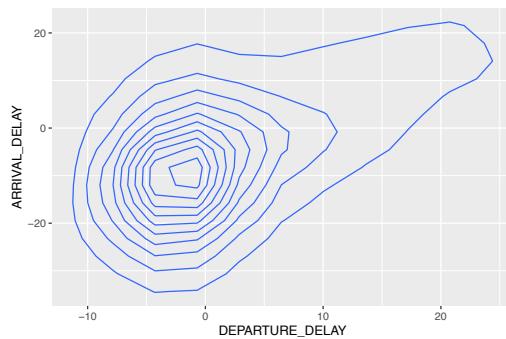
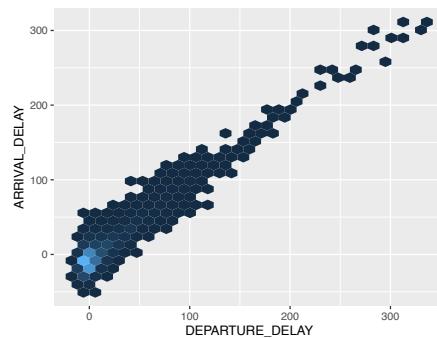
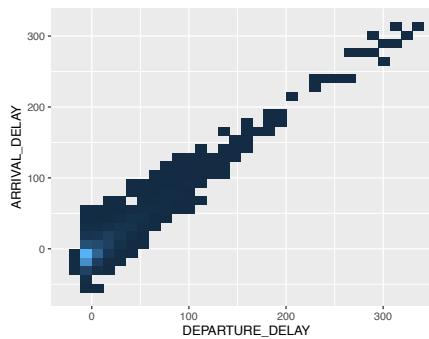
### [3] Scatter plot & jittering

```
1 more_filtered_flights <- filtered_flights %>%
2   filter(MONTH==1, DAY==1, !is.na(DEPARTURE_DELAY), !is.na(ARRIVAL_DELAY),
3   DEPARTURE_DELAY<400, ARRIVAL_DELAY<400)
4 ggplot(more_filtered_flights) +
5   geom_point(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY))
6 ggplot(more_filtered_flights) +
7   geom_point(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY), alpha=.1)
8 ggplot(more_filtered_flights) +
9   geom_jitter(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY), alpha=.1, width=10,
height=10)
```



### [3] Counting after binning: bin2d, hex, density\_2d

```
1 ggplot(more_filtered_flights) +  
2   geom_bin2d(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY))  
3 ggplot(more_filtered_flights) +  
4   geom_hex(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY))  
5 ggplot(more_filtered_flights) +  
6   geom_density_2d(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY))
```

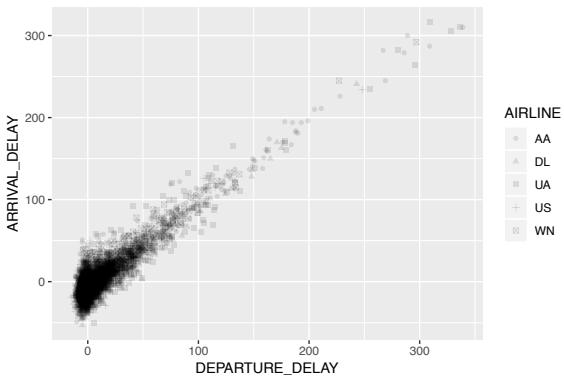
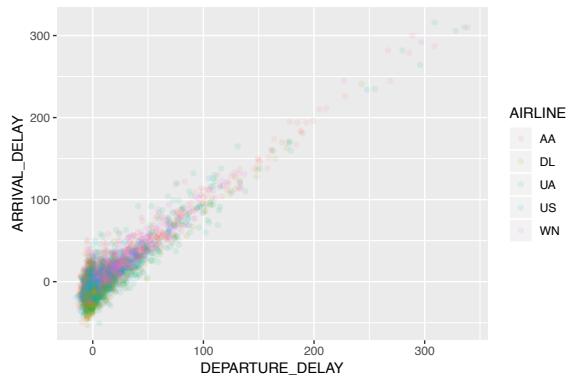


# Relationship among three or more variables

- It's all about visually encoding (=mapping) variables onto aesthetic elements.
- What do we have?
  - Size
  - Color
  - Fill
  - Shape
  - Alpha
  - Line type
  - Weight
  - Facets (small multiples)

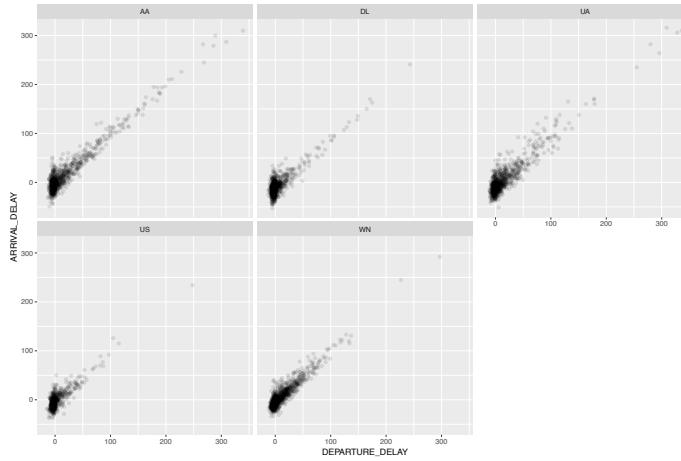
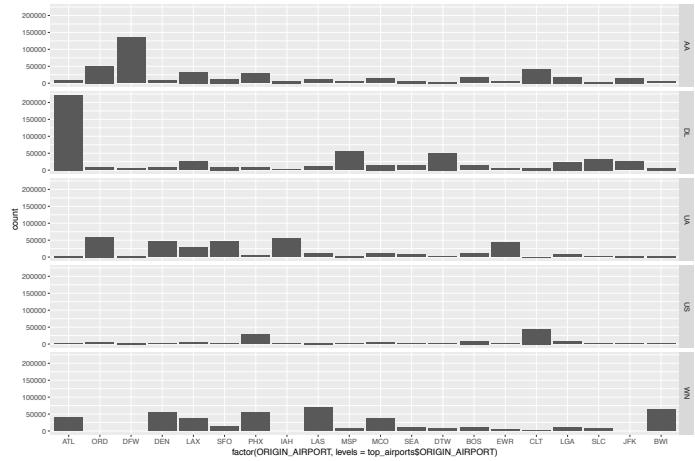
# Examples encoding more than two variables

```
1 ggplot(more_filtered_flights) +  
2   geom_point(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY, color=AIRLINE),  
3   alpha=.1)  
4 ggplot(more_filtered_flights) +  
5   geom_point(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY, shape=AIRLINE),  
6   alpha=.1)
```



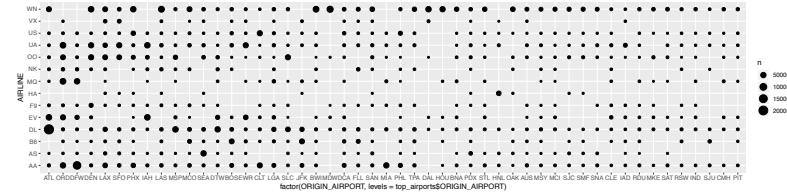
# Facets (small multiples)

```
1 ggplot(filtered_flights) +
2   geom_bar(aes(x=factor(ORIGIN_AIRPORT, levels=top_airports$ORIGIN_AIRPORT)))+
3   facet_grid(rows=vars(AIRLINE))
4 ggplot(more_filtered_flights) +
5   geom_point(aes(x=DEPARTURE_DELAY, y=ARRIVAL_DELAY), alpha=.1) +
6   facet_wrap(~AIRLINE)
```



# Visualizing grouped summary stats

- So far, we learned how to visualize distribution of a single variable, two variables, and three variables.
- To do that, we just “counted” the number of observations.
- We can plot other summary statistics such as mean instead of count.
- Imagine that color or size in the visualizations below represents average departure delay instead of the number of observations.



# Shortcut to create a plot: qplot

- qplot is a convenient wrapper for ggplot. The following example is from [?qplot](#).

```
1 # Use data from data.frame
2 qplot(mpg, wt, data = mtcars)
3 qplot(mpg, wt, data = mtcars, colour = cyl)
4 qplot(mpg, wt, data = mtcars, size = cyl)
5 qplot(mpg, wt, data = mtcars, facets = vs ~ am)
6
7 # qplot will attempt to guess what geom you want depending on the input
8 # both x and y supplied = scatterplot
9 qplot(mpg, wt, data = mtcars)
10 # just x supplied = histogram
11 qplot(mpg, data = mtcars)
12 # just y supplied = scatterplot, with x = seq_along(y)
13 qplot(y = mpg, data = mtcars)
14
15 # Use different geoms
16 qplot(mpg, wt, data = mtcars, geom = "path")
17 qplot(factor(cyl), wt, data = mtcars, geom = c("boxplot", "jitter"))
18 qplot(mpg, data = mtcars, geom = "dotplot")
```

# Visualization Typology

# Dr. Abela's chart chooser

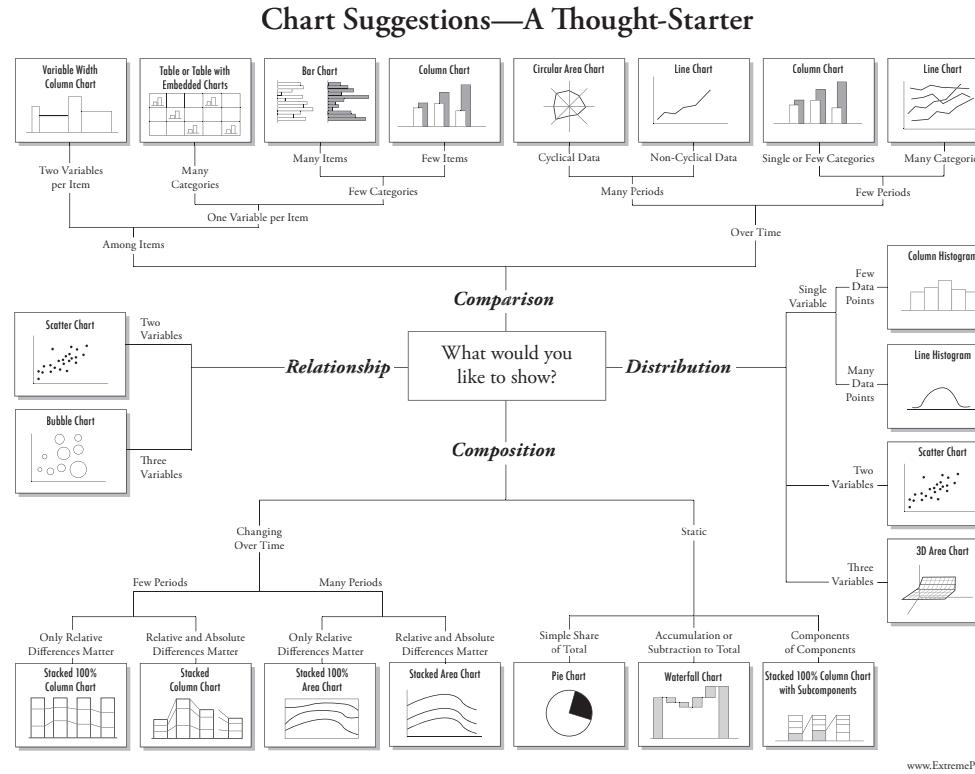


Figure from <https://extremepresentation.com/design/7-charts/>

# A periodic table of visualization methods

A PERIODIC TABLE OF VISUALIZATION METHODS																																													
> ⊖ < <b>C</b> continuum		> ⊖ < <b>Tb</b> table		> ⊖ < <b>Ga</b> cartesian coordinates		> ⊖ < <b>Pi</b> pie chart		> ⊖ < <b>L</b> line chart		> ⊖ < <b>B</b> bar chart		> ⊖ < <b>Ac</b> area chart		> ⊖ < <b>R</b> radar chart		> ⊖ < <b>Pa</b> parallel coordinates		> ⊖ < <b>Hy</b> hyperolic tree		> ⊖ < <b>Cy</b> cycle diagram		> ⊖ < <b>T</b> timeline		> ⊖ < <b>Ve</b> venn diagram		> ⊖ > <b>Mi</b> mindmap		> ⊖ > <b>Sq</b> square of oppositions		> ⊖ < <b>Cc</b> concentric circles		> ⊖ < <b>Ar</b> argument slide		> ⊖ < <b>Sw</b> swim lane diagram		> ⊖ < <b>Gc</b> gantt chart		> ⊖ > <b>Pm</b> perspectives diagram		> ⊖ < <b>D</b> dilemma diagram		> ⊖ > <b>Pr</b> parameter ruler		> ⊖ < <b>Kn</b> knowledge map	
> ⊖ < <b>Hi</b> histogram		> ⊖ < <b>Sc</b> scatterplot		> ⊖ < <b>Sa</b> sankey diagram		> ⊖ < <b>In</b> information lens		> ⊖ < <b>E</b> entity relationship diagram		> ⊖ < <b>Pt</b> petri net		> ⊖ < <b>Fl</b> flow chart		> ⊖ < <b>Cl</b> clustering		> ⊖ < <b>Lc</b> layer chart		> ⊖ < <b>Py</b> minto pyramid technique		> ⊖ < <b>Ce</b> cause-effect chains		> ⊖ < <b>Tl</b> toulmin map		> ⊖ < <b>Dt</b> decision tree		> ⊖ < <b>Cp</b> cpm critical path method		> ⊖ > <b>Cf</b> concept fan		> ⊖ < <b>Co</b> concept map		> ⊖ < <b>Ic</b> iceberg		> ⊖ < <b>Lm</b> learning map											
> ⊖ < <b>Tk</b> tukey box plot		> ⊖ < <b>Sp</b> spectrogram		> ⊖ < <b>Da</b> data map		> ⊖ < <b>Tp</b> tremap		> ⊖ < <b>Cn</b> cone tree		> ⊖ < <b>Sy</b> system dyn./simulation		> ⊖ < <b>Df</b> data flow diagram		> ⊖ < <b>Se</b> semantic network		> ⊖ < <b>So</b> soft system modeling		> ⊖ < <b>Sn</b> synergy map		> ⊖ > <b>Fo</b> force field diagram		> ⊖ < <b>Ib</b> ibni argumentation map		> ⊖ < <b>Pr</b> process event chains		> ⊖ < <b>Pe</b> pert chart		> ⊖ < <b>Ev</b> evocative knowledge map		> ⊖ < <b>V</b> Vee diagram		> ⊖ > <b>Hh</b> heaven 'n' hell chart		> ⊖ < <b>I</b> infomural											
<b>Cy</b> <b>Process Visualization</b>																	Note: Depending on your location and connection speed it can take some time to load a pop-up picture. © Ralph Lengler & Martin J. Eppler, www.visual-literacy.org																		version 1.5										
<b>Hy</b> <b>Structure Visualization</b>		<b>Overview Detail</b>		<b>Detail AND Overview</b>		<b>Divergent thinking</b>		<b>Convergent thinking</b>		> ⊖ < <b>Su</b> supply demand curve		> ⊖ < <b>Pc</b> performance charting		> ⊖ < <b>St</b> strategy map		> ⊖ < <b>Oc</b> organisation chart		> ⊖ < <b>Ho</b> house of quality		> ⊖ < <b>Fd</b> feedback diagram		> ⊖ < <b>Fr</b> failure tree		> ⊖ < <b>Mq</b> magic quadrant		> ⊖ < <b>Ld</b> life-cycle diagram		> ⊖ < <b>Po</b> porter's five forces		> ⊖ < <b>S</b> s-cycle		> ⊖ < <b>Sm</b> stakeholder map		> ⊖ < <b>Is</b> ishikawa diagram		> ⊖ < <b>Tc</b> technology roadmaps									
<b>Ed</b> edgeworth box		<b>Pf</b> portfolio diagram		<b>Sg</b> strategic game board		<b>Mz</b> mintberg's organigraph		<b>Z</b> zwickly's morphological box		<b>Ad</b> affinity diagram		<b>De</b> decision discovery diagram		<b>Bm</b> bg matrix		<b>Stc</b> strategy canvas		<b>Vc</b> value chain		<b>Hy</b> hyper-cycle		<b>Sc</b> stakeholder rating map		<b>Sr</b> stakeholder rating map		<b>Ta</b> tags		<b>Sd</b> spray diagram																	

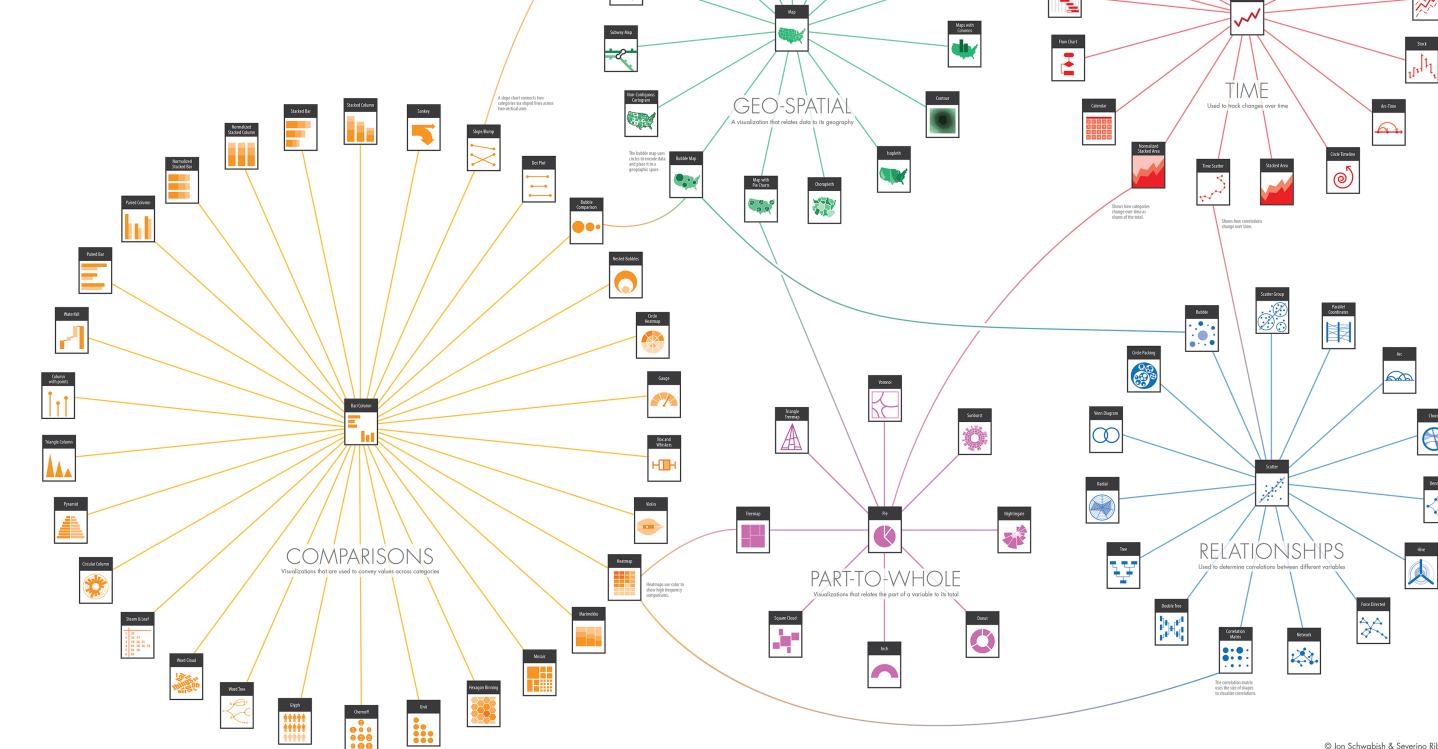
Figure from [http://www.visual-literacy.org/periodic\\_table/periodic\\_table.html](http://www.visual-literacy.org/periodic_table/periodic_table.html)

## THE GRAPHIC CONTINUUM

Graphs use different shapes, markings, and colors to show patterns and relationships within and across variables. The Graphic Continuum is a representation – a map, if you will – of the different kinds of graphs one might use to plot data and show different relationships. The relationships between graph types is inherently non-linear: a bar chart, for example, can be used to show comparisons and also changes over time, while variations on a scatterplot can be used to show correlations within a single time period or over time.

The graphs shown are not meant to be an exhaustive list of all types of graphs, nor do the connections represent all of the possible links. The Continuum shows some representative graphs and the links that bind them together. The five categories shown here: Comparisons, Time, Geo-spatial, Partio-Whole, and Relationships are also not meant to be exhaustive as there exist overlap between them as there exist overlap between data bases and visualization techniques.

You can use this map as a guide to possible graphic choices, but your imagination and your data will help you determine other graphs you can use to present your data.



© Jan Schumach & Sonja Röhrs

Figure from <https://visual.ly/blog/graphic-continuum/>

# The Data Visualization Catalog

The Data Visualisation Catalogue

About · Blog · Shop · Resources

What do you want to show?

Here you can find a list of charts categorised by their data visualization functions or by what you want a chart to communicate to an audience. While the allocation of each chart into specific functions isn't a perfect system, it still works as a useful guide for selecting chart based on your analysis or communication needs.

Comparisons   Proportions   Relationships   Hierarchy

Concepts   Location   Part-to-a-whole   Distribution

How things work   Processes & methods   Movement or flow   Patterns

Range   Data over time   Analysing text   Reference tool

Search by Function   View by List

中文   Español   Русский   Türkçe

Search by Function	View by List				
Arc Diagram	Area Graph	Bar Chart	Box & Whisker Plot	Brainstorm	Bubble Chart
Bubble Map	Bullet Graph	Calendar	Candlestick Chart	Chord Diagram	Choropleth Map
Circle Packing	Connection Map	Density Plot	Donut Chart	Dot Map	Dot Matrix Chart
Error Bars	Flow Chart	Flow Map	Gantt Chart	Heatmap	Histogram
Illustration Diagram	Kagi Chart	Line Graph	Marimekko Chart	Multi-set Bar Chart	Network Diagram
Nightingale Rose Chart	Non-ribbon Chord Diagram	Open-high-low-close Chart	Parallel Coordinates Plot	Parallel Sets	Pictogram Chart

Figure from <https://datavizcatalogue.com/>

# Data Visualization with ggplot2 :: CHEAT SHEET



## Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and geoms—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot(data = <DATA>) +
  <GEOM_FUNCTION>(<mapping> = aes(<POSITIONS>),
  stat = <STAT>, position = <POSITION>) +
  <COORDINATE_FUNCTION> +
  <FACET_FUNCTION> +
  <SCALE_FUNCTION> +
  <THEME_FUNCTION>
```

required  
Not required, sensible defaults supplied

`ggplot(data = mpg, aes(x = cyl, y = hwy))` Begins a plot that you finish by adding layers to. Add one geom function per layer.

**aesthetic mappings**    **data**    **geom**  
`qplot(x = cyl, y = hwy, data = mpg, geom = "point")`  
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

`last_plot()` Returns the last plot

`ggsave("plot.png", width = 5, height = 5)` Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.

## Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemployed))
b <- ggplot(seals, aes(x = long, y = lat))

# + geom_blank()
# (Useful for expanding limits)

# + geom_curve(aes(yend = lat + 1,
#   xend = long + 1, curvature = -1), x = xend, y = yend,
#   alpha, angle, color, curvature, linetype, size)

# + geom_path(lineend = "butt", linejoin = "round",
#   x, y, alpha, color, group, linetype, size)

# + geom_polygon(aes(group = group))
# x, y, alpha, color, fill, group, linetype, size

# + geom_rect(aes(xmin = long, ymin = lat, xmax =
#   long + 1, ymax = lat + 1)) -> xmax, ymin,
#   alpha, color, fill, linetype, size

# + geom_ribbon(aes(ymin = unemployed - 900,
#   ymax = unemployed + 900)) -> ymax, ymin,
#   alpha, color, fill, group, linetype, size
```

### LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

```
b + geom_abline(aes(intercept=0, slope=1))
b + geom_hline(aes(intercept = lat))
b + geom_vline(aes(intercept = long))

b + geom_segment(aes(yend=lat+1, vend=long+1))
b + geom_spoke(aes(angle = 1:115, radius = 1))
```

### ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

# + geom_area(stat = "bin")
# x, y, alpha, color, fill, linetype, size, weight

# + geom_density(kernel = "gaussian")
# x, y, alpha, color, fill, group, linetype, size, weight

# + geom_dotplot()
# x, y, alpha, color, fill

# + geom_freqpoly() x, y, alpha, color, group,
#   linetype, size

# + geom_histogram(binwidth = 5) x, y, alpha,
#   color, fill, linetype, size, weight

# + geom_qq(aes(sample = hwy)) x, y, alpha,
#   color, fill, linetype, size, weight
```

### discrete

```
d <- ggplot(mpg, aes(f))
# + geom_bar()
```

### TWO VARIABLES

**continuous x, continuous y**

```
e <- ggplot(mpg, aes(cty, hwy))

# + geom_label(aes(label = cyl), nudge_x = 1,
#   nudge_y = 1, check_overlap = TRUE) x, y, label,
#   alpha, angle, color, family, fontface, hjust,
#   lineheight, size, vjust

# + geom_jitter(height = 2, width = 2)
# x, y, alpha, color, fill, shape, size

# + geom_point() x, y, alpha, color, fill, shape,
#   size, stroke

# + geom_quantile() x, y, alpha, color, group,
#   linetype, size, weight

# + geom_rug(sides = "bl") x, y, alpha, color,
#   linetype, size

# + geom_smooth(method = lm) x, y, alpha,
#   color, fill, group, linetype, size, weight

# + geom_text(aes(label = cyl), nudge_x = 1,
#   nudge_y = 1, check_overlap = TRUE) x, y, label,
#   alpha, angle, color, family, fontface, hjust,
#   lineheight, size, vjust
```

### discrete x , continuous y

```
f <- ggplot(mpg, aes(class, hwy))

# + geom_col() x, y, alpha, color, fill, group,
#   linetype, size

# + geom_boxplot() x, y, lower, middle, upper,
#   ymax, ymin, alpha, color, fill, group, linetype,
#   shape, size, weight

# + geom_dotplot(binaxis = "y", stackdir =
#   "center") x, y, alpha, color, fill, group

# + geom_violin(scale = "area") x, y, alpha, color,
#   fill, group, linetype, size, weight
```

### discrete x , discrete y

```
g <- ggplot(diamonds, aes(cut, color))

# + geom_count() x, y, alpha, color, fill, shape,
#   size, stroke
```

### THREE VARIABLES

```
sealsSz <- with(seals, sqrt(delta_long^2 + delta_lat^2)) l <- ggplot(seals, aes(long, lat))

# + geom_contour(aes(z = z))
# x, y, z, alpha, colour, group, linetype, size, weight

# + geom_raster(aes(fill = z), hijust = 0.5, vjust = 0.5,
#   interpolate = FALSE)
# x, y, alpha, fill

# + geom_tile(aes(fill = z)) x, y, alpha, color, fill,
#   linetype, size, width
```

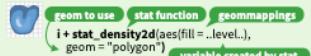
## Stats

An alternative way to build a layer

A stat builds new variables to plot (e.g., count, prop).



Visualize a stat by changing the default stat of a geom function, `geom_bar(stat="count")` or by using a stat function, `stat_count(geom="bar")`, which calls a default geom to make a layer (equivalent to a geom function). Use `.name..` syntax to map stat variables to aesthetics. Use `.variable` to map stat variables to aesthetics.



```

c + stat_bin(binwidth = 1, origin = 10)
x, y | ..count.., ..ncount.., ..density.., ..ndensity..
c + stat_count(width = 1) x, y | ..count.., ..prop..
c + stat_density(adjust = 1, kernel = "gaussian")
x, y, | ..count.., ..density.., ..scaled..
e + stat_bin_2d(bins = 30, drop = T)
x, y, fill | ..count.., ..density..
e + stat_hex(bins = 30) x, y, fill | ..count.., ..density..
e + stat_density_2d(contour = TRUE, n = 100)
x, y, color, size | ..level..
e + stat_ellipse(level = 0.95, segments = 51, type = "t")
l + stat_contour(aes(z = z)) x, y, z, order | ..level..
l + stat_summary_hex(aes(z = z), bins = 30, fun = max)
x, y, z, fill | ..value..
l + stat_summary_2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..value..
f + stat_boxplot(coef = 1.5) x, y | ..lower...
..middle..., ..upper..., ..width..., ..ymin..., ..ymax..
f + stat_ydensity(kernel = "gaussian", scale = "area") x, y |
..density.., ..scaled.., ..count.., ..n.., ..violinwidth.., ..width..
e + stat_ecdf(n = 40) x, y | ..x.., ..
e + stat_quantile(quartiles = c(0.1, 0.9), formula = y ~ log(x), method = "rq") x, y | ..quantile..
e + stat_smooth(method = "lm", formula = y ~ x, se = T, level = 0.95) x, y | ..se.., ..x.., ..ymin.., ..ymax..
ggplot() + stat_function(aes(x = -3:3), n = 99, fun = dnorm, args = list(sd = 0.5)) x, y, ..
e + stat_identity(harm = TRUE)
ggplot() + stat_qq(aes(sample = 1:100), dist = qt, dparam = list(df = 5)) sample, x, y | ..sample.., ..theoretical..
e + stat_sum(x, y, size | ..n.., ..prop..
e + stat_summary(fun.data = "mean_cl_boot")
h + stat_summary_bin(fun.y = "mean", geom = "bar")
e + stat_unique()

```

## Scales

Scales map data values to the visual values of an aesthetic. To change a mapping, add a new scale.



### GENERAL PURPOSE SCALES

Use with most aesthetics

`scale_*` continuous - map cont' values to visual ones  
`scale_*` discrete - map discrete values to visual ones  
`scale_*` identity - use data values as visual ones  
`scale_*` manual(values = c...) - map discrete values to manually chosen visual ones  
`scale_*` date(date, labels = "%m/%d%Y") - date values as dates.  
`scale_*` datetime() - treat data x values as date times. Use same arguments as `scale_x_date()`. See ?strptime for formats.

### X & Y LOCATION SCALES

Use with x or y aesthetics (x shown here)

`scale_x_log10()` - Plot x on log10 scale  
`scale_x_reverse()` - Reverse direction of x axis  
`scale_x_sqrt()` - Plot x on square root scale

### COLOR AND FILL SCALES (DISCRETE)

`n <- d + geom_bar(aes(fill = fl))`  
`n + scale_fill_brewer(palette = "Blues")`  
For palette choices:  
RColorBrewer::display.brewer.all()  
`n + scale_fill_distiller(palette = "Blues")`

`n + scale_fill_gradient(low = "red", high = "yellow")`  
`n + scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)`

`n + scale_fill_hexbin(hex.colors = brewer.pal(6, "Reds"))`  
Also: rainbow(), heat.colors(), terrain.colors(), cm.colors(), RColorBrewer::brewer.pal()

### SHAPE AND SIZE SCALES

`p <- e + geom_point(aes(shape = fl, size = cyl))`  
`p + scale_shape() + scale_size()`  
`p + scale_shape_manual(values = c(3:7))`  
`p + scale_shape_discrete()`  
`p + scale_size_continuous()`  
`p + scale_size_area(max_size = 6)`

## Coordinate Systems

`r <- d + geom_bar()`

`r + coord_cartesian(xlim = c(0, 5))`  
The default cartesian coordinate system  
`r + coord_fixed(ratio = 1/2)`  
ratio, xlim, ylim  
Cartesian coordinates with fixed aspect ratio  
`r + coord_flip()`  
Flipped Cartesian coordinates  
`r + coord_polar(theta = "x", direction = 1)`  
theta, start, direction  
Polar coordinates  
`r + coord_trans(xtrans = "sqrt")`  
xtrans, ytrans, xlim, ylim  
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

`r + coord_quickmap()`  
`r + coord_map(projection = "ortho", orientation = c(45, -74, 0))` projection, orientation, xlim, ylim  
Map projections from the mapproj package  
(mercator (default), aezimuthala, lagrange, etc.)

## Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

`s <- ggplot(mpg, aes(f1, fill = drv))`  
`s + geom_bar(position = "dodge")`  
Arrange elements side by side.  
`s + geom_bar(position = "fill")`  
Stack elements on top of one another, normalize height.  
`e + geom_point(position = "jitter")`  
Add random noise to x and y position of each element to avoid overplotting.  
`e + geom_label(position = "nudge")`  
Nudge labels away from points.  
`s + geom_bar(position = "stack")`  
Stack elements on top of one another.

Each position adjustment can be recast as a function with `manual` width and `height` arguments  
`s + geom_bar(position = position_dodge(width = 1))`

`r + theme_bw()` White background with gray lines  
`r + theme_light()` Gray background  
`r + theme_linedraw()` Minimal themes  
`r + theme_dark()` dark for contrast  
`r + theme_void()` Empty theme

## Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

`t <- ggplot(mpg, aes(cty, hwy)) + geom_point()`

`t + facet_grid(cols = vars(f1))`  
facet into columns based on f1  
`t + facet_grid(rows = vars(year))`  
facet into rows based on year  
`t + facet_grid(rows = vars(year), cols = vars(f1))`  
facet into both rows and columns  
`t + facet_wrap(vars(f1))`  
wrap facets into a rectangular layout

`Set scales to let axis limits vary across facets`

`t + facet_grid(rows = vars(drv), cols = vars(f1), scales = "free")`  
x and y axis limits adjust to individual facets  
`"free", "x"` - x axis limits adjust  
`"free", "y"` - y axis limits adjust

`Set labeller to adjust facet labels`

`t + facet_grid(cols = vars(f1), labeller = label_both)`  
`f1: c f1: d f1: e f1: p f1: r`  
`t + facet_grid(rows = vars(f1), labeller = label_bquote(alpha ^ .(f1)))`  
`alpha^c alpha^d alpha^e alpha^p alpha^r`

## Labels

`t + labs( x = "New x axis label", y = "New y axis label", title = "Add a title above the plot", subtitle = "Add a subtitle below title", caption = "Add a caption below plot", `aes` = "New AES legend title")`  
Use scale functions to update legend labels

`t + annotate(geom = "text", x = 8, y = 9, label = "A")`  
geom to place manual values for geom's aesthetics

## Legends

`n + theme(legend.position = "bottom")`  
Place legend at "bottom", "top", "left", or "right"  
`n + guides(fill = "none")`  
Set legend type for each aesthetic: colorbar, legend, or none in the legend  
`n + scale_fill_discrete(name = "Title", labels = c("C", "B", "A", "D", "E"))`  
Set legend title and labels with a scale function.

## Zooming

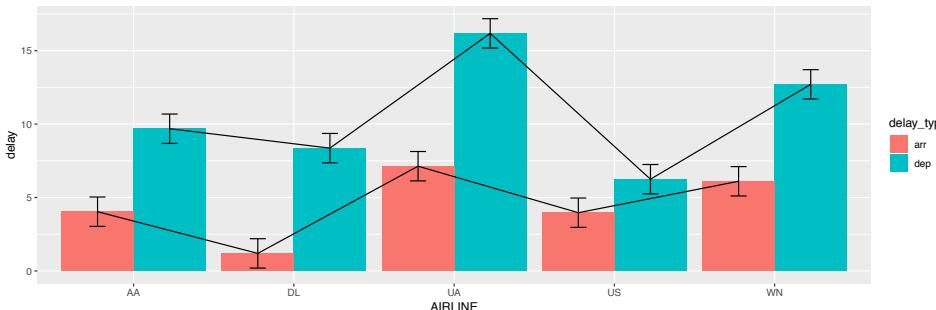
Without clipping (preferred)  
`t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))`  
With clipping (removes unseen data points)  
`t + xlim(0, 100) + ylim(10, 20)`  
`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`

# Customizing Visualization

# Layering multiple geoms

- The power of the grammar of graphics is that you can stack up visualizations as long as they share the same core mappings.

```
1 delay_summary <- filtered_flights %>% group_by(AIRLINE) %>%
2   summarize(dep=mean(DEPARTURE_DELAY, na.rm=T), arr=mean(ARRIVAL_DELAY,
3     na.rm=T)) %>%
4   gather(delay_type, delay, -AIRLINE)
5
6 ggplot(delay_summary, aes(x=AIRLINE, y=delay, group=delay_type)) +
7   geom_col(aes(fill=delay_type), position="dodge") +
8   geom_line(position=position_dodge(width=0.9)) +
9   geom_errorbar(aes(ymin=delay-1, ymax=delay+1), width=0.2,
10   position=position_dodge(width=0.9))
```



# Annotation

- For plotting data-driven labels, use geom\_label or geom\_text.
- You can also annotate your plot with text, line, etc.
- Review the following two URLs.
  - geom\_text: [https://ggplot2.tidyverse.org/reference/geom\\_text.html](https://ggplot2.tidyverse.org/reference/geom_text.html)
  - annotate: <https://ggplot2.tidyverse.org/reference/annotate.html>

# Legends and labels

- Legend customization
  - [http://www.cookbook-r.com/Graphs/Legends\\_\(ggplot2\)/](http://www.cookbook-r.com/Graphs/Legends_(ggplot2)/)
- Title, subtitle, axis labels: ggtitle(), xlab(), ylab()
  - <https://ggplot2.tidyverse.org/reference/labs.html>
- Tick label style: theme(axis.text.x = ..., axis.text.y = ...)
  - <http://www.sthda.com/english/wiki/ggplot2-axis-ticks-a-guide-to-customize-tick-marks-and-labels>

# Axis scaling & reordering

**Scales** map data values to the visual values of an aesthetic. To change a mapping, add a new scale.

(n <- d + geom\_bar(aes(fill = fl)))

scale\_      aesthetic to adjust      prepackaged scale to use      scale-specific arguments

n + scale\_fill\_manual()

values = c("skyblue", "royalblue", "blue", "navy"),  
 limits = c("D", "E", "P", "R"), breaks = c("D", "E", "P", "R"),  
 name = "fuel", labels = c("D", "E", "P", "R"))

range of values to include in mapping  
 title to use in legend/axis  
 labels to use in legend/axis  
 breaks to use in legend/axis

## GENERAL PURPOSE SCALES

Use with most aesthetics

**scale\_\*** \_continuous() - map cont' values to visual ones  
**scale\_\*** \_discrete() - map discrete values to visual ones  
**scale\_\*** \_identity() - use data values as visual ones  
**scale\_\*** \_manual(values = c()) - map discrete values to manually chosen visual ones

**scale\_\***(date\_labels = "%m/%d"), date\_breaks = "2 weeks") - treat data values as dates

**scale\_\*\*\_datetime()** - treat data x values as date times.  
Use same arguments as scale\_x\_date(). See ?strftime for  
label formats.

## X & Y LOCATION SCALES

Use with x or y aesthetics (x shown here)

**scale\_x\_log10()** - Plot x on log10 scale

**scale\_x\_reverse()** - Reverse direction of x axis

**scale\_x\_sqrt()** - Plot x on square root scale

## COLOR AND FILL SCALES (DISCRETE)

```
n <- d + geom_bar(aes(fill = fl))
```

```
p + scale_fill_brewer(palette = "Blues")
```

For palette choices:

RColorBrewer::display.brewer.all()

## COLOR AND FILL SCALES (CONTINUOUS)

```
g <- c + geom_dotplot(aes(fill = ..x..))
```

**at scale fill distiller(palette = "Blues")**

Journal of CHI, 2009, Vol 6(1), 1-100

```
o + scale_fill_gradient2(low="red", high="blue",  
mid="white", midpoint=25)
```

**o + scale\_fill\_gradientn(colours=topo.colors(6))**  
Also: rainbow(), heat.colors(), terrain.colors(), cm.colors(), RColorBrewer::brewer.pal()

## SHAPE AND SIZE SCALES

```
p <- e + geom_point(aes(shape = fl, size = cyl))
```

**p + scale shape() + scale size()**

**p + scale shape manual**(values = c(3:7))

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

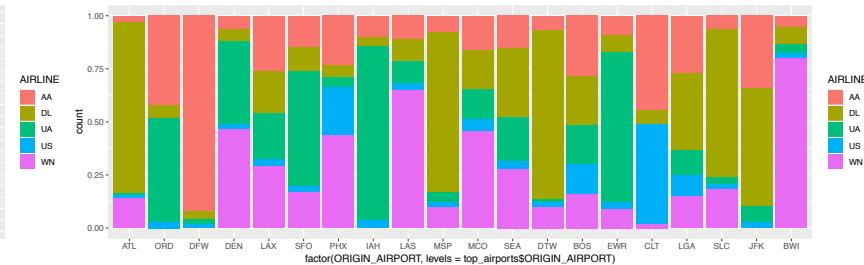
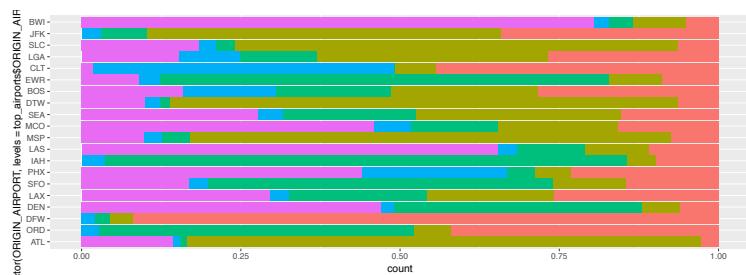
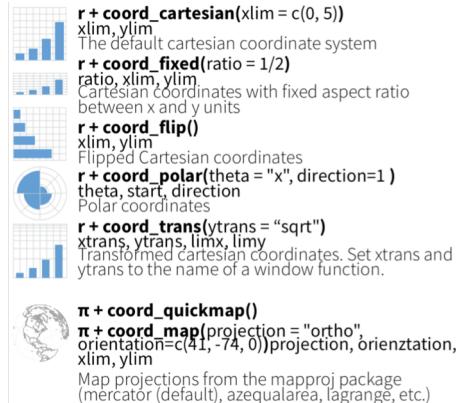
□○△+×◊▽⊗\*⊕◊田⊗田□○

**p + scale\_radius(range = c(1,6))**

**p + scale size area(max\_size = 6)**

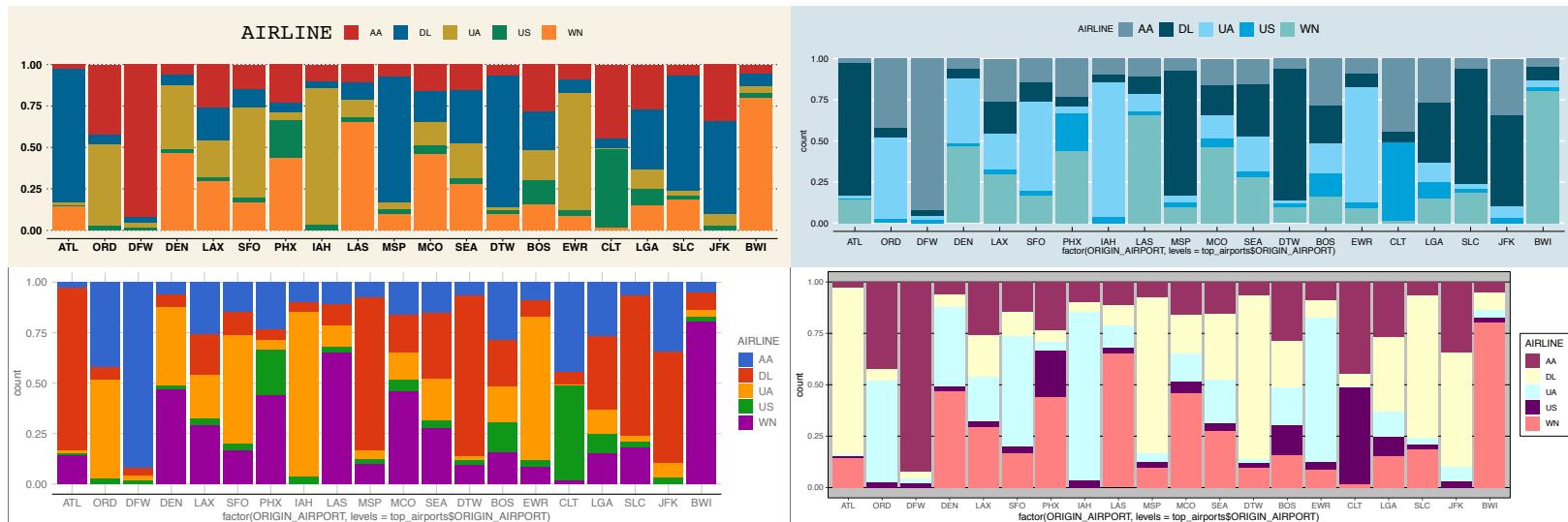
# Coordinate

- Mainly used functions
  - coord\_cartesian(): default coordinate system
  - coord\_polar(): one we used for creating pie chart
  - coord\_map(): for geospatial visualization
  - coord\_flip(): to interchange x and y axes



# Themes

- Default themes: theme\_bw(), theme\_linedraw(), theme\_light(), theme\_dark(), theme\_minimal(), theme\_classic(), theme\_void()
- Additional themes via ggthemes: theme\_tufte(), theme\_solarized(), theme\_excel(), ...



# Weekly Recap

# Things we covered this week

- More data wrangling
  - Join (merge)
  - String
  - Factor
  - Datetime
- Covariation
  - Covariance / Correlation
  - Visualizing two variables
  - Visualizing three or more variables
- Visualization typology
- Customizing visualization
  - Layering multiple geometries
  - Annotation
  - Labels
  - Legends
  - Axes scaling & reordering
  - Coordinate
  - Theme

# Things to do this week

- 4 Concept Checker Quizzes (Due: Thursday, January 17, 11:59pm)
  - Week 2.1. Merge, String, Factor, Date
  - Week 2.2. Covariation
  - Week 2.3. Visualization Topology
  - Week 2.4. Customizing Visualization
- 1 Weekly Problem (Due: Friday, January 18, 11:59pm)
- 3 DataCamp Courses (Due: Friday, January 18, 11:59pm)
  - [Text Mining with R](#)
    - [\[1\] String Manipulation in R with stringr](#)
  - Individual Courses
    - [Categorical Data in the Tidyverse](#)
    - [Working with Dates and Times in R](#)
- 1 In-class Activity (Due: Saturday, January 19, 11:59pm)