

Intro to SYS-6018

Data Mining

SYS 6018 | Fall 2020

00-intro.pdf

Contents

1	Requirements	2
2	About us	2
2.1	About the Instructor	2
2.2	About you	2
2.3	About our TA	2
3	The course	2
3.1	Becoming Data Scientists	2
3.2	Topics	2
3.3	Examples	3
3.4	How this course fits within Data Science	3
4	Syllabus	4
4.1	Course Webpage	4
4.2	Course Prereqs	4
4.3	Exercise 1	5
4.4	Other Syllabus Material	5
4.5	Succeeding in this course	6
5	Introduction to R	7
5.1	Getting Help	7
5.2	RStudio	7
5.3	Using R Packages	7
5.4	Graphics with the <code>ggplot2</code> package	8
5.5	Data Transformation with the <code>dplyr</code> package	8
5.6	Groupwise operations	9
5.7	Relational Data and Joins	11
5.8	Data Importing	12
5.9	Tidy Data with the <code>tidyr</code> package	12
6	RMarkdown	12

1 Requirements

- Working version of R and RStudio
- Install R package: `tidyverse` and `nycflights13`
- Course Webpage: <https://mdporter.github.io/SYS6018>

2 About us

2.1 About the Instructor

2.1.1 My Job

- Faculty Webpage <http://www.faculty.virginia.edu/mdporter/index.html>
- GitHub <https://github.com/mdporter>
 - Blog <https://mdporter.github.io/blog/>

2.1.2 Family

2.1.3 Hobbies

2.2 About you

Fill out a notecard with the following information:

1. Your name (with pronunciation hints)
2. Hometown (include country/region if far away)
3. Degrees
4. What type of job to hope to land on graduation (industry, title)
5. 2 things you hope to learn in this course
6. 2 interesting things about you (to help me remember you)

2.3 About our TA

Mojtaba Heidarysafa

3 The course

3.1 Becoming Data Scientists

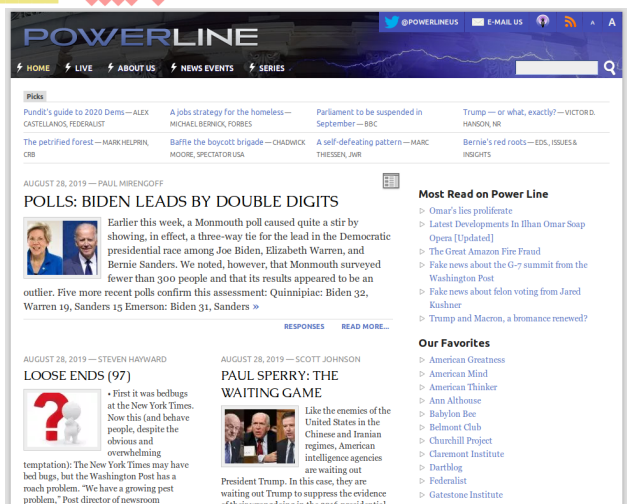
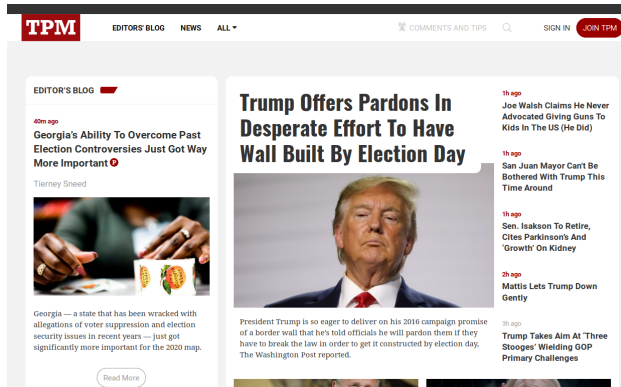
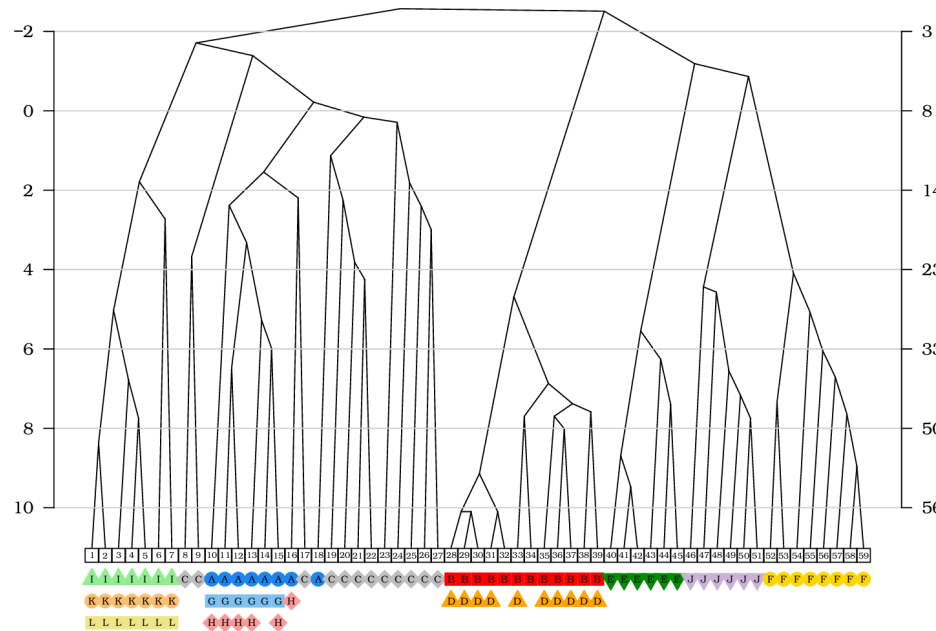
- How do you learn in the DS domain?
- What is knowledge?

3.2 Topics

- See website: <https://mdporter.github.io/SYS6018>
- Data scientists are expected to have *fluency* in:
 - data analysis, modeling, stats, ML, coding, algorithms, probability, etc.
 - This class will be a mix of all
- Data scientists are expected to be problem solvers

- doing good on structured homework sets isn't sufficient

3.3 Examples



3.4 How this course fits within Data Science

1. Science of Data Collection/Acquisition (Systems + Value)
The science of “how”, “when”, “where” data is collected.
2. Science of Data Storage and Representation (Design + Systems)
The science of storing and representing data
3. Science of Data Manipulation/Transformation

The science of transforming/preparing data for analysis

4. Science of Computing with Data (Systems + Analytics)

The science of computing for data analysis (algorithms and performance)

5. Science of Data Analytics (Analytics)

Science of Modelling with Data (or Machine Learning)

6. Science of Summarizing/Communicating Data and Models (Design + Value)

The science of summarizing data for human understanding Another view: *The science of extracting and communicating the information in data*

7. Science of Practicing Data Science (Practice)

The Science of the Data Science System

8. Disciplinary Data Science

The science of applying data science to your discipline

4 Syllabus

4.1 Course Webpage

- We have a course webpage <https://mdporter.github.io/SYS6018/>
 - lectures
 - R scripts
 - data sets
- We will use the Collab site for homework submission, solutions, etc.
- We will use the Slack channel: UVA MSDS 2021/[sys-6018](#) for other group communication
 - Change name
 - Add Mojtaba

4.2 Course Prereqs

- Linear Regression
 - Multiple Linear Regression
 - Logistic Regression
 - Categorical Predictors (dummy coding)
 - Implementation in R (lm(), predict(), coef(), etc.)
 - Estimation / Model Fitting
 - Cross-validation
- Probability and Statistics
 - Bayes Theorem
 - CDF/PDF/PMF
 - Maximum Likelihood Estimation
 - Distributions: normal, binomial, hypergeometric, etc.
 - Expected value, variance, median, quantiles

- Mean Square Error
 - Confidence Intervals
 - Hypothesis Testing
- Math
 - Calculus
 - Matrix Calculations
 - PCA, SVD
- Computing
 - data types: vector, matrix, array, list, etc.
 - writing simple functions
 - flow control: loops, if/else, etc.
 - data wrangling
 - generating random variables
 - RMarkdown [*Note: practice HW will cover RMarkdown*]

4.3 Exercise 1

Your Turn #1

Let X_1, X_2, \dots, X_n be the yearly number of crashes at an intersection (X_i is number of crashes in year i).

- What is an estimate of the probability that there are 100 crashes in year $n + 1$?

4.4 Other Syllabus Material

- Office Hours
- Textbooks
- R, RStudio
- Course Assessment
 - *The exam dates are posted in the Class Schedule (on the course website). Note these now so there are no conflicts.*
 - RMarkdown (See HW0)
 - Class participation. Expect you come prepared with questions. Don't be afraid to ask questions in class. Now is your time to learn.
- Course Management
- **Honor Code**
- Read all syllabus and ask me questions

4.5 Succeeding in this course

- Schedule time now (10 hours out-of-class time)
- Basic Flow
 - Thursdays: New material introduced
 - Thursdays: Homework due before class
 - Changes Nov 3rd
 - Example
 - a. Thursday morning: Read intro material (for Thurs class); finish/review homework
 - b. Thursday day: Participate in Lecture
 - c. Friday: Start homework, Read secondary material (from Thurs class)
 - d. Monday: Read intro material (for Tues class), continue homework
 - e. Tuesday: Participate in Lecture
 - f. Wed: Read secondary material (from Tue class); finish homework; Read intro material (for Thur class)
- Attend office hours!
- I find that most students struggle with coding. This really hinders your ability to get your mind about the concepts and slows down your learning. The course will use R, but all examples will use the tidyverse dialect. There is no better tool for interactive data analysis and both exploratory and confirmatory modeling. Tidyverse is a major improvement over base R, but it can look a bit different and take some time becoming familiar with. The free online book [R for Data Science](#) and [website](#) provide a good introduction and reference. UVA Library has created some resources: [Data Science Essentials in R](#). RStudio [cheatsheets](#). While I encourage tidyverse, you are free to use anything for homework (<https://rstudio.github.io/reticulate/>).
- I find students understand the least about statistical concepts. This is so fundamental to all of ML and Data Mining; a strong grasp of statistics will enable the connections between topics to pop out. If you already feel comfortable coding, I suggest you go a quick stat review. Here is one introductory resource: <https://www.openintro.org/book/isrs/>
- The students who gain the most from the program will embrace mathematical equations. As they say “an equation is worth a thousand words”. While we won’t do any proofs in this class, we will judiciously use equations to clarify concepts. Spend time to become intimate with math notation – it is worth the investment.
- **Use Trustworthy Material**
 - The assigned readings are trustworthy
 - Blogs and videos you find on the web are not
 - Please don’t trust: Toward Data Science, Analytics Vidha, Machine Learning Mastery, Medium

5 Introduction to R

5.1 Getting Help

- A good source of basic data analysis using R is found in the free book [R for Data Science](#).
- Web search, especially *stackoverflow.com* and *stats.stackexchange.com*
- Troubleshooting/Debugging.
 - Check one line of code at a time.
 - Use scripts
 - Make sure it works in plain R before incorporating into Rmd

5.2 RStudio

- Install R and RStudio
- Make use of *Projects* in RStudio

5.3 Using R Packages

It takes two steps to use the functions and data in an R package

1. Install the package
 - i.e., download the package to your computer
 - this only needs to be done one time
 - `install.packages()`
2. Load the package
 - i.e., tell R to look for the package functions and/or data
 - this needs to be done every time R is started (and you want to use the package)
 - `library()`

5.3.1 Note on tidyverse package

- The tidyverse package <https://www.tidyverse.org/packages/> is really just a wrapper to load several related R packages
 - `ggplot2` for graphics
 - `dplyr` for data manipulation
 - `tidyr` for getting data into tidy form
 - `readr` for loading in data
 - `tibble` for improved data frames
 - `purrr` for functional programming
 - `stringr` for string manipulation
 - `forcats` for categorical/factor data
- This provides a nice shortcut to load all of these packages with `library(tidyverse)` instead of each separately:

```
#- the hard way  
library(ggplot2)  
library(dplyr)  
library(tidyr)
```

```
library(readr)
library(tibble)
library(purrr)
library(stringr)
library(forcats)
```

```
#- the easy way
library(tidyverse)
```

5.4 Graphics with the `ggplot2` package

The `ggplot2` package is an approach to creating graphics for data analysis.

- See <https://ggplot2.tidyverse.org/>
- Keep the [ggplot2 cheatsheet](#) handy

5.5 Data Transformation with the `dplyr` package

- See <https://dplyr.tidyverse.org/>
- Keep the [dplyr cheatsheet](#) handy

5.5.1 single table verbs

1. `filter()`: find/keep certain rows
 - alternative to `base::subset()`
 - `slice()` to keep by row number
 - helper functions: `between()`: numeric values in a range
2. `arrange()`: reorder rows
 - alternative to `base::order()`
 - helper functions: `desc()` to use descending order
3. `select()`: find/keep certain columns
 - helper functions: `starts_with()`, `ends_with()`, `matches()`, `contains()`, `?select`
4. `mutate()`: add/create new variables
 - alternative to `base::transform()`
 - `transmute()`: only return new variables
5. `summarize()`: produce summary statistics
 - don't confuse with `summary()`
 - most useful when data is *grouped*

5.5.2 Chaining/Pipes

- Multiple operations can be chained together with the *pipe* operator, `%>%`, (pronounced as *then*). Technically, it performs `x %>% f(y) -> f(x, y)`. This lets you focus on the verbs, or actions you are performing.

```
x = c(1:5, NA)
x %>% mean(na.rm=TRUE)
#> [1] 3
```



```
mean(x, na.rm=TRUE)
#> [1] 3
```

Your Turn #2

1. Load the `nycflights13` package, which contains airline on-time data for all flights departing NYC in 2013. Also includes useful 'metadata' on airlines, airports, weather, and planes.
2. Load the `tidyverse` package
3. Using the `flights` data,
 - find all flights that were less than 1000 miles (`distance`)
 - Keep only the columns: `dep_delay`, `arr_delay`, `origin`, `dest`, `air_time`, and `distance`
 - Add the Z-score for departure delays
 - Convert the departure and arrival delays into hours
 - Calculate the average flight speed (in mph)
 - order by average flight speed (fastest to slowest)
 - return the first 12 rows

5.5.3 Other useful `dplyr` functions

- `distinct()`: retain unique/distinct rows
- `sample_n()` and `sample_frac()`: randomly sample rows
- `top_n()`: selects and orders the top n rows according to wt
- `add_column()` add new column in particular position
- `add_row()` adds new row(s) to the table
- `na_if(x, y)` converts the y valued elements in x to NA

```
x = c(1, 2, -99, 5, 5, -99)
na_if(x, -99)           # replace -99 with NA
#> [1] 1 2 NA 5 5 NA
```

- `coalesce(x, y)` replaces the NA in x with y

```
x = c(1, 2, NA, 5, 5, NA)
coalesce(x, 0)          # replace NA with 0
#> [1] 1 2 0 5 5 0
```

5.6 Groupwise operations

5.6.1 Split - Apply - Combine

The `dplyr` operations are more powerful when they can be used with grouping variables. Split - Apply - Combine.

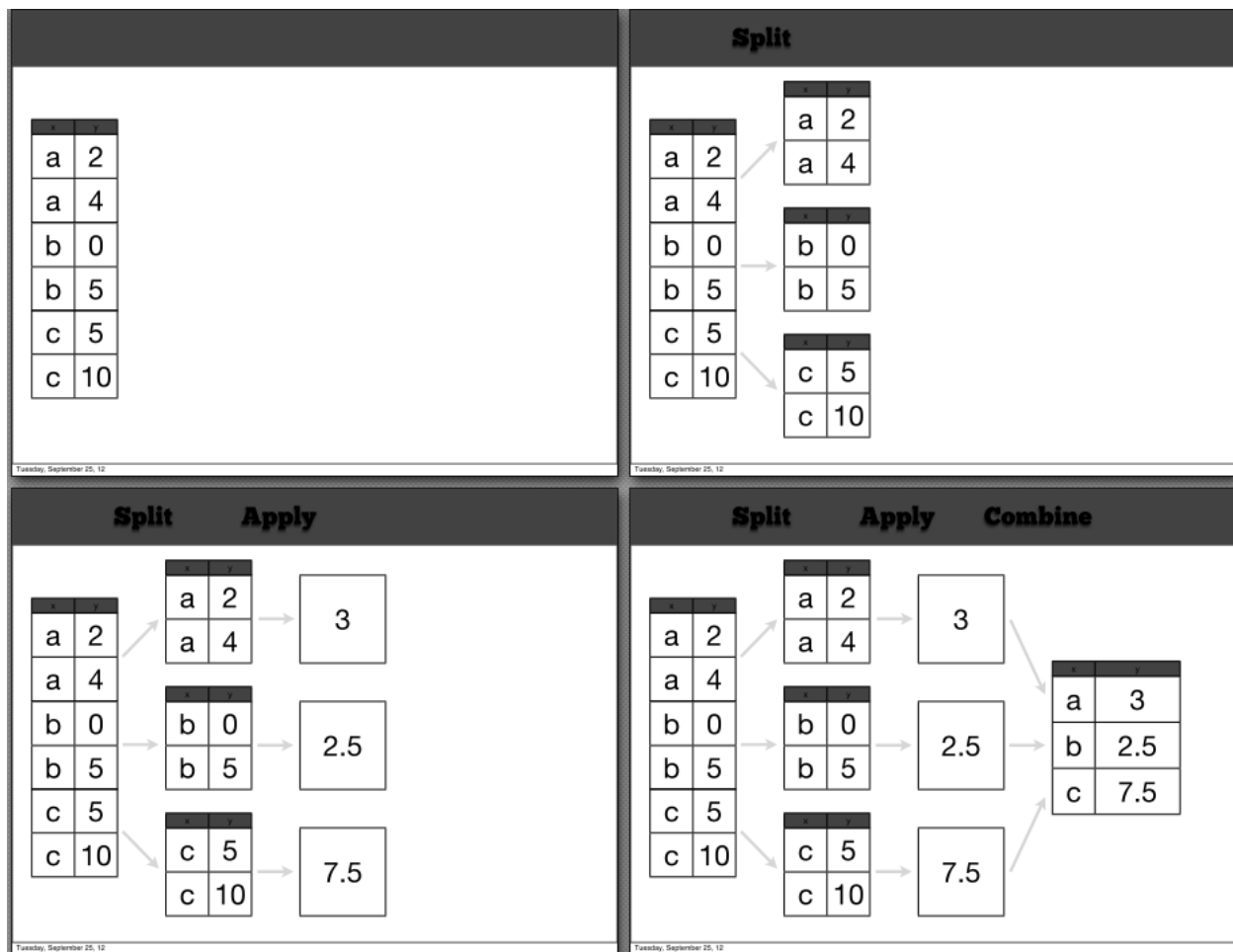


Image from Hadley Wickham UseR tutorial June 2014 <http://www.dropbox.com/sh/i8qnlwmuieicxc/AAAgT9tIKoIm7WZKIyK25lh6a>

5.6.2 group_by()

First use the `group_by()` function to group the data (determines how to split), then apply function(s) to each group using the `summarise()` function. Note: grouping should to be applied on discrete variables (categorical, factor, or maybe integer valued columns).

```
flights %>%
  group_by(origin, dest) %>% # group by both origin and dest
  summarise(max.delay = max(arr_delay, na.rm=TRUE),
            avg.delay = mean(arr_delay, na.rm=TRUE),
            min.delay = min(arr_delay, na.rm=TRUE),
            count = n() )    # n() gives the group count

#> # A tibble: 224 x 6
#> # Groups:   origin [3]
#>   origin dest max.delay avg.delay min.delay count
#>   <chr> <chr>   <dbl>   <dbl>   <dbl> <int>
#> 1 EWR   ALB     328    14.4    -34   439
#> 2 EWR   ANC      39    -2.5    -47     8
#> 3 EWR   ATL    796    13.2    -39  5022
#> 4 EWR   AUS    349   -0.474   -59   968
#> 5 EWR   AVL    228     8.80   -26   265
#> 6 EWR   BDL    266     7.05   -43   443
```

```
#> # ... with 218 more rows
```

- `count(...)` is a shortcut for `group_by(...) %>% summarize(n=n())`
- `ungroup()` removes the grouping

5.6.3 Grouped Mutate and Filter

- When data is *grouped*, `mutate()` and `filter()` operate on each group independently

```
#- proportion of carrier at each dest
flights %>%
  count(dest, carrier) %>%
  group_by(dest) %>% # group by dest
  mutate(total=sum(n), p=n/sum(n)) %>% # grouped mutate sum(n) is by group
  arrange(desc(total), -p) # arrange by most freq dest and prop
#> # A tibble: 314 x 5
#> # Groups:   dest [105]
#>   dest carrier      n total      p
#>   <chr> <chr>   <int> <int>   <dbl>
#> 1 ORD   UA      6984 17283 0.404
#> 2 ORD   AA      6059 17283 0.351
#> 3 ORD   MQ      2276 17283 0.132
#> 4 ORD   9E      1056 17283 0.0611
#> 5 ORD   B6       905 17283 0.0524
#> 6 ORD   EV        2 17283 0.000116
#> # ... with 308 more rows
```

5.7 Relational Data and Joins

Joins are used to combine or merge two datasets. This is a major aspect of SQL.

5.7.1 Mutating Joins

See 13.4 of R4DS

- `inner_join(x, y)` only includes observations that having matching x and y key values. Rows of x can be dropped/filtered.
- `left_join(x, y)` includes all observations in x, regardless of whether they match or not. This is the most commonly used join because it ensures that you don't lose observations from your primary table.
- `right_join(x, y)` includes all observations in y. It's equivalent to `left_join(y, x)`, but the columns will be ordered differently.
- `full_join()` includes all observations from x and y.
- The left, right and full joins are collectively know as **outer joins**. When a row doesn't match in an outer join, the new variables are filled in with missing values.
 - **outer joins** will fill any missing values with NA
- If there are duplicate keys, all combinations are returned.
- Missing values are given NA.

5.8 Data Importing

5.8.1 readr package

- See <https://readr.tidyverse.org/>
- Keep the [data import cheatsheet](#) handy

5.8.2 readxl package

- See <https://readxl.tidyverse.org/> for importing excel files

5.9 Tidy Data with the tidyr package

- <https://tidyr.tidyverse.org/>
- Keep the [data import cheatsheet](#) handy. Page two describes the `tidyr` functionality

5.9.1 Why Tidy Data?

- Tidy data (in form of a data frame) is usually the best form for analysis
 - some exceptions are for modeling (e.g., matrix manipulations and algorithms)
- For presentation of data (e.g., in tables), non-tidy form can often do better
- the functions in `tidyr` usually allow us to covert from non-tidy to tidy for analysis and also from tidy to non-tidy for presentation

5.9.2 Main tidyr functions

function	description
<code>spread()</code>	Spreads a pair of key:value columns into a set of tidy columns
<code>gather()</code>	Gather takes multiple columns and collapses into key-value pairs, duplicating all other columns as needed. You use <code>gather()</code> when you notice that you have columns that are not variables
<code>separate()</code>	turns a single character column into multiple columns
<code>unite()</code>	paste together multiple columns into one (reverse of <code>separate()</code>)

6 RMarkdown

- Homework will be submitted in Rmd and (pdf/html) format
- When you *knit* a Rmd, it:
 1. starts a new instance of R (clean environment)
 2. in the current directory
- Any data or code must first be put into the Rmd file
 - The Rmd won't know about anything in another script or in your R environment
 - Any `source()` or data paths are relative to the current directory of the Rmd

-
- A homework template will be provided for each homework