

ASSIGNMENT FRONT SHEET

Course Name: ALY6040 Data Mining Applications

Professor Name: Nagadeepa Shanmuganathan

Student Name: Dong Quoc Tuong (Lukas)

Student Class: Fall 2019 CPS Term: Winter 2021

Module 3: Data Mangling and Data Wrangling & Logistic Regression, Decision Trees, and Random Forests

Completion Date: February 7th Due Time:12:00am

Statement of Authorship

I confirm that this work is my own. Additionally, I confirm that no part of this coursework, except where clearly quoted and referenced, has been copied from material belonging to any other person e.g. from a book, handout, another student. I am aware that it is a breach of Northeastern University's regulations to copy the work of another without clear acknowledgement and that attempting to do so renders me liable to disciplinary procedures. To this effect, I have uploaded my work onto Turnitin and have ensured that I have made any relevant corrections to my work prior to submission.

X Tick here to confirm that your paper version is identical to the version submitted through Turnitin

The dataset we are working on called Gapminder.org, detailing the values for life expectancy, GDP per capita, and population, and every five years from 1952 -2007 for approximately 150 countries. In this assignment, we will be using 2 packages: dplyr and tidyr

1/ dplyr

dplyr is a grammar of data manipulation, giving the coders a set of verbs that allow them to solve the most common data manipulation problems ("A Grammar of Data Manipulation • dplyr," 2020)

First, we have a look at "filter" in dplyr. It allows you to only iinlcude information that satisfy the condition you set up in the beginning. For example, you can see that we only include any inputs that has "lifeExp" < 29 or "country"= "Mexico" or "Mexico, Afghanistan"

```
filter(gapminder, lifeExp < 29)</p>
+ A tibble: 2 x 6
                                                pop gdpPercap
               continent year lifeExp
  country
                <fct>
                           <int>
                                     <db1>
 Afghanistan Asia
                            <u>1</u>952
                                      28.8 8425333
                                                           779.
 Rwanda
               Africa
                            <u>1</u>992
                                      23.6 7290203
 filter(gapminder, country == "Mexico")
A tibble: 12 x 6
  country continent year lifeExp
                                               pop gdpPercap
                        <int>
                                 <db1>
                                             <int>
                                                         \langle db 1 \rangle
   <fct>
            <fct>
            Americas
                                  50.8 30144317
                                                         3478.
                         1952
1 Mexico
                                  55.2 35<u>015</u>548
58.3 41<u>121</u>485
                         <u>1</u>957
            Americas
 2 Mexico
                                                         4132.
  Mexico
            Americas
                         1962
                                                         4582.
4 Mexico
            Americas
                         <u>1</u>967
                                  60.1 47995559
                                                         5755.
 5 Mexico
            Americas
                         <u>1</u>972
                                   62.4
                                          55<u>984</u>294
                                                         <u>6</u>809.
6 Mexico
            Americas
                         1977
                                   65.0
                                          63759976
                                                         <u>7</u>675.
            Americas
                         <u>1</u>982
                                   67.4
                                         71640904
                                                         9611.
  Mexico
  Mexico
                         1987
                                   69.5
                                          80122492
                                                         8688.
            Americas
                                                         9472.
  Mexico
            Americas
                         <u>1</u>992
                                   71.5
                                         88111030
                                                         <u>9</u>767.
.0 Mexico
            Americas
                         1997
                                  73.7 95895146
  Mexico
            Americas
                         <u>2</u>002
                                  74.9 102479927
                                                        <u>10</u>742.
            Americas
                         2007
                                  76.2 108700891
                                                        <u>11</u>978.
  Mexico
· filter(gapminder, country %in% c("Mexico",
                                                    "Afghanistan"))
A tibble: 24 x 6
  country
                 continent year lifeExp
                                                   pop gdpPercap
                 <fct>
                             <int>
                                      \langle dh 1 \rangle
                                                 <int>
1 Afghanistan Asia
                             1952
                                       28.8 8425333
                                                              779
 2 Afghanistan Asia
                              <u>1</u>957
                                       30.3
                                              9240934
                                                              821.
  Afghanistan Asia
                              <u>1</u>962
                                       32.0 10267083
                                                              853.
  Afghanistan Asia
                              <u>1</u>967
                                       34.0 11537966
                                                              836.
5 Afghanistan Asia
                              1972
                                       36.1 13079460
                                                              740.
```

There are two other ways do it. The first method is to identify the column inside the dataset with "\$" and the second is using "subset" while retaining the same code structure like "filter"

```
gapminder[gapminder$lifeExp < 29, ]</pre>
                                                              ## repeat
j] indexing is distracting
A tibble: 2 x 6
                 continent year lifeExp
                                                      pop gdpPercap
 country
 \langle fct \rangle
                 <fct>
                               <int>
                                         \langle db 1 \rangle
                                                    <int>
                                                                 <db1>
 Afghanistan Asia
                               <u>1</u>952
                                          28.8 8<u>425</u>333
                                                                  779.
 Rwanda
                Africa
                               <u>1</u>992
                                          23.6 7290203
                                                                  737.
 subset(gapminder, country == "Mexico")
... but wait ...
A tibble: 12 x 6
  country continent year lifeExp
                                                     pop gdpPercap
  <fct>
             <fct>
                           <int>
                                     \langle db 1 \rangle
                                                  <int>
                                                                \langle db 1 \rangle
            Americas
                           <u>1</u>952
                                      50.8 30144317
                                                               <u>3</u>478.
L Mexico
 Mexico
             Americas
                           1957
                                      55.2
                                              35<u>015</u>548
                                                               <u>4</u>132.
 Mexico
             Americas
                           <u>1</u>962
                                      58.3 41<u>121</u>485
                                                               <u>4</u>582.
                           1967
                                              47995559
 Mexico
             Americas
                                      60.1
                                                               <u>5</u>755.
 Mexico
            Americas
                           <u>1</u>972
                                      62.4
                                              55<u>984</u>294
                                                               <u>6</u>809.
            Americas
                           <u>1</u>977
 Mexico
                                      65.0
                                              63<u>759</u>976
                                                               <u>7</u>675.
```

Second, we start to dig a little bit deeper in Pipe "%>%" Operator. This operator forwards a value or a result from the expression into the next function call/ expression. Consequently, removing duplication and make your code looks cleaner, readable and more efficient to the readers. For example, codes like this

gapminder %>% head

will have the same result like this

head(gapminder)

```
gapminder %>% head # this...
A tibble: 6 x 6
               continent year lifeExp
country
                                                   pop gdpPercap
               <fct>
                           <int>
                                      \langle dh 1 \rangle
                                                              < db7 >
                            <u>1</u>952
                                                               779.
Afghanistan Asia
                                       28.8 8425333
                                                                                   > select(gapminder, year, lifeExp) # this...
# A tibble: 1.704 x 2
Afghanistan Asia
                             1957
                                       30.3
                                             9240934
                                                               821.
                                                                                      year lifeExp
Afghanistan Asia
                             <u>1</u>962
                                       32.0 10<u>267</u>083
                                                               853.
                             <u>1</u>967
Afghanistan Asia
                                       34.0 11537966
                                                               836.
                                                                                      <u>1</u>952
<u>1</u>957
Afghanistan Asia
                             <u>1</u>972
                                       36.1 13079460
                                                               740.
                                       38.4 14<u>880</u>372
Afghanistan Asia
                             <u>1</u>977
                                                                                      1962
                                                                                               32.0
head(gapminder) #
                      ...is the same as this!
A tibble: 6 x 6
                                                                                      1977
1982
1987
                                                                                               38.4
country
               continent year lifeExp
                                                   pop gdpPercap
                                                                                               39.9
40.8
               <fct>
                                      <db1>
                                                              <db1>
                                                               779.
Afghanistan Asia
                            1952
                                       28.8 8<u>425</u>333
                                                                                               41.7
Afghanistan Asia
                             <u>1</u>957
                                       30.3 9240934
                                                               821.
                                                                                    ... with 1,694 more rows gapminder %5% select(year, lifeExp) # ...is the same as this! A tibble: 1,704 \times 2
Afghanistan Asia
                             1962
                                       32.0 10267083
                                                               853.
Afghanistan Asia
                            <u>1</u>967
                                       34.0 11537966
                                                               836.
Afghanistan Asia
                             <u>1</u>972
                                       36.1 13079460
                                                               740.
                                                                                      year lifeExp
Afghanistan Asia
                            <u>1</u>977
                                       38.4 14<u>880</u>372
                                                               786.
                                                                                      1957
                                                                                               30.3
                                                                                      1962
1967
gapminder %>% head(3) # can pass arguments! this...
                                                                                               32.0
                                                                                      1972
1977
                                                                                               36.1
               continent year lifeExp
country
                                      <db1>
                                                               <db1>
                                                                                      1982
                                                                                               39.9
Afghanistan Asia
                             <u>1</u>952
                                       28.8
                                               8<u>425</u>333
                                                               779.
                                                                                      1987
                                                                                               40.8
Afghanistan Asia
                             1957
                                       30.3
                                               9240934
                                                                821.
Afghanistan Asia
                             <u>1</u>962
                                       32.0 10267083
                                                               853.
head(gapminder, 3) # ...is the same as this! A tibble: 3 \times 6
country
               continent year lifeExp
                                                    pop gdpPercap
               <fct>
                            <int>
                                      < dh7 >
                                                  <int>
                                                               <db1>
Afghanistan Asia
                             <u>1</u>952
                                                               779.
                                       28.8 8425333
Afghanistan Asia
                             <u>1</u>957
                                       30.3 9240934
                                                                821.
                                       32.0 10267083
                             <u>1</u>962
Afghanistan Asia
                                                               853.
```

"%>%" is so much better when typing out because it does not have to use the dplyr package.

Compared to the base R, "%>%" allows coders to fix the code much easier because they can tell exactly the issue of the code and sometimes do not need to type in one consecutive line.

```
gapminder %>%
  filter(country == "Cambodia") %>%
  select(-continent, -lifeExp) # same as
A tibble: 12 x 4
                          pop gdpPercap
 country year
  <fct>
             <int>
                        <int>
                                   <db1>
 Cambodia <u>1</u>952 4<u>693</u>836
                                      368.
 Cambodia <u>1</u>957
                                      434.
                     5<u>322</u>536
 Cambodia
              <u>1</u>962
                     6083619
                                      497.
 Cambodia <u>1</u>967
                      6960067
                                      523.
 Cambodia <u>1</u>972 7<u>450</u>606
                                      422.
Cambodia <u>1</u>977 6<u>978</u>607
Cambodia <u>1</u>982 7<u>272</u>485
                                       525.
                                      624.
 Cambodia <u>1</u>987 8<u>371</u>791
                                      684.
```

Third, the "mutate" adds a new column to the data frame that you are working. For example, f we decide to filter out to only any inputs of Cambodia, then select all the variable except "continent", "lifeExp" and mutate/ create a new column gdp = pop * gdpPercap, we will have something like this

```
gapminder %>%
    filter(country == "Cambodia") %>%
    select(-continent, -lifeExp) %>%
mutate(gdp = pop * gdpPercap)
 A tibble: 12 x 5
  country year
                                   pop gdpPercap
                  <int>
                                                <db1>
   <fct>
                                                  368. <u>1</u>729<u>534</u>398.
434. <u>2</u>310<u>184</u>671.
497. <u>3023033</u>308.
L Cambodia <u>1</u>952 4<u>693</u>836
                   <u>1</u>957
  Cambodia
                             5<u>322</u>536
                   <u>1</u>962
3 Cambodia
                             6083619
                   <u>1</u>967
                             6<u>960</u>067
                                                  523. <u>3</u>643<u>123</u>977.
422. <u>3</u>141<u>354</u>496.
4 Cambodia
5 Cambodia
                   <u>1</u>972 7<u>450</u>606
                                                  525. <u>3</u>663<u>574</u>552.
624. <u>4541488550</u>.
5 Cambodia <u>1</u>977 6<u>978</u>607
7 Cambodia 1982 7272485
```

Forth, the "summarise" or "summarize" adds new column when grouping to together. For In the first example, we see that we calculate the "mean_gdp" for Cambodia because we use filter option but for the second, we calculate for all of them using the additional "group by" function

```
gapminder %>%
    filter(country == "Cambodia") %>%
    select(-continent, -lifeExp) %>%
mutate(gdp = pop * gdpPercap) %>%
group_by(country) %>%
    summarize(mean_gdp = mean(gdp)) %>%
    ungroup() # if you use group_by, also use
later
A tibble: 1 x 2
 country
                mean_qdp
  Cambodia 6596612377.
  ## summarize for all countries (replaces ou
 gapminder %>%
  select(-continent, -lifeExp) %>%
    mutate(gdp = pop * gdpPercap) %>%
group_by(country) %>%
    summarize(mean_gdp = mean(gdp)) %>%
    ungroup() # if you use group_by, also use
A tibble: 142 x 2
                        mean_gdp
   country
1 Afghanistan <u>12</u>709<u>647</u>583.
 2 Albania
                     9094669267.
 3 Algeria
                    <u>96</u>735<u>171</u>261.
                    <u>25</u>532<u>681</u>843.
 4 Angola
5 Argentina <u>266</u>754<u>123</u>835.
6 Australia <u>320</u>253<u>755</u>823.
                 <u>158</u>579<u>002</u>935.
  Austria
                    <u>7</u>694<u>793</u>798.
8 Bahrain
9 Bangladesh <u>80</u>648<u>494</u>456.
O Belgium <u>197</u>371<u>599</u>665.
```

2/ tidyr

While dplyr is certainly helpful, researchers also employ tidyr as well to create tidy dataset. Tidy data is a normalized data that is stored in a standard way and can be used to analyze or visualize insights.

First, using "gather", we can group the different columns together into one. "gather" is often used with "separate", which turns a single character column into multiple columns. Thus, there are four ways to write the code as followed

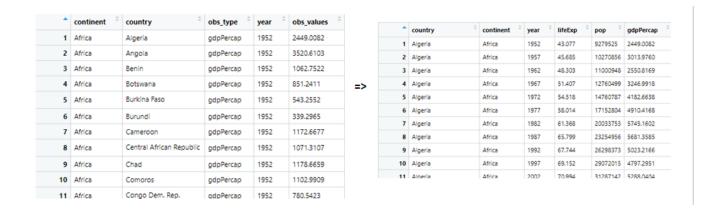
And the result turns the unnormalized dataset in the left to the noramlized one in the right, which is much cleaner for process

| | | | | | | | | • | continent | country | obstype_year | obs_values |
|----|-----------|--------------------------|----------------|----------------|----------------|------|----|----|-----------|--------------------------|----------------|------------|
| | | | | | | | | 1 | Africa | Algeria | gdpPercap_1952 | 2449.0082 |
| • | continent | country | gdpPercap_1952 | gdpPercap_1957 | gdpPercap_1962 | gdp | | 2 | Africa | Angola | gdpPercap_1952 | 3520.6103 |
| 1 | Africa | Algeria | 2449.0082 | 3013.9760 | 2550.8169 | 32 1 | | 3 | Africa | Benin | gdpPercap_1952 | 1062.7522 |
| 2 | Africa | Angola | 3520.6103 | 3827.9405 | 4269.2767 | 55 | | 4 | Africa | Botswana | gdpPercap_1952 | 851.2411 |
| 3 | Africa | Benin | 1062.7522 | 959.6011 | 949.4991 | 10 | | 5 | Africa | Burkina Faso | gdpPercap_1952 | 543.2552 |
| 4 | Africa | Botswana | 851.2411 | 918.2325 | 983.6540 | 12 | | | Africa | Burundi | gdpPercap_1952 | 339.2965 |
| 5 | Africa | Burkina Faso | 543.2552 | 617.1835 | 722.5120 | 79 | | | | | | |
| 6 | Africa | Burundi | 339.2965 | 379.5646 | 355.2032 | 41 | | 7 | Africa | Cameroon | gdpPercap_1952 | 1172.6677 |
| 7 | Africa | Cameroon | 1172.6677 | 1313.0481 | 1399.6074 | 15 | | 8 | Africa | Central African Republic | gdpPercap_1952 | 1071.3107 |
| 8 | Africa | Central African Republic | 1071.3107 | 1190.8443 | 1193.0688 | 11 | | 9 | Africa | Chad | gdpPercap_1952 | 1178.6659 |
| 9 | Africa | Chad | 1178.6659 | 1308.4956 | 1389.8176 | 11 | => | 10 | Africa | Comoros | gdpPercap_1952 | 1102.9909 |
| 10 | Africa | Comoros | 1102.9909 | 1211.1485 | 1406.6483 | 18 | | 11 | Africa | Congo Dem. Rep. | gdpPercap_1952 | 780.5423 |
| 11 | Africa | Congo Dem, Rep. | 780.5423 | 905.8602 | 896.3146 | 86 | | 12 | Africa | Congo Rep. | gdpPercap_1952 | 2125.6214 |
| | Africa | Congo Rep. | 2125.6214 | 2315.0566 | 2464.7832 | 26 | | 13 | Africa | Cote d'Ivoire | gdpPercap_1952 | 1388.5947 |

Second, "spread" is the opposite of "gather", spreading the data from one column into multiple ones. Here are three ways to do it

```
## spread() from normal to wide
gap_normal <- gap_long %>%
    spread(obs_type, obs_values) %>%
    select(country, continent, year, lifeExp, pop, gdpPercap)
# or
gap_normal <- gap_long %>%
    spread(obs_type, obs_values)
gap_normal <- gap_normal[,names(gapminder)]</pre>
```

And the result will turns the normalized dataset in the left to the one in the right



Third, "all.equal" is used to compare R objects x and y testing "near equality". In the case that they are significantly different, it will make an effort to compare them and produce a report of differences. After comparing the "gap normal" and "gapreminder" we have this

```
all.equal(gap_normal,gapminder)
[1] "Attributes: < Component "class": Lengths (1, 3) differ (string e on first 1) >"
[2] "Attributes: < Component "class": 1 string mismatch >"
[3] "Component "country": Modes: character, numeric"
[4] "Component "country": Attributes: < target is NULL, current is 1 >"
[5] "Component "country": target is character, current is factor"
[6] "Component "continent": Modes: character, numeric"
[7] "Component "continent": Attributes: < target is NULL, current is 2 ""
[8] "Component "continent": target is character, current is factor"
[9] "Component "year": Modes: character, numeric"
[9] "Component "year": target is character, current is numeric"
[10] "Component "gear": target is character, current is numeric"
[11] "Component "lifeExp": Mean relative difference: 0.203822"
[12] "Component "gdpPercap": Mean relative difference: 1.634504"
[13] "Component "gdpPercap": Mean relative difference: 1.162302"
```

Forth, "unite" function allows us to paste together multiple columns into one.

Final project's proposal

In America, there are 11,000 new cases of invasive cervical cancers being found annually. Despite that fact that the number of new cases are on the downward trend over the past decade, it still kills about 4,000 women in America and 300,000 women globally.(Fontham et al., 2020) The sooner one discovers it, the more chances they have to survive. Thus, in order to effectively eliminate such diseases, we need to create a Machine Learning algorithm that can detect the cancer as soon as possible

For the final project, I am planning to analyze the dataset: Cervical cancer (Risk Factors) from UCI. The dataset was collected at 'Hospital Universitario de Caracas' in Caracas, Venezuela. Cervical cancer has 858 inputs with variables ranging from demographic information, habits to historic medical records. Some variables have missing data as patients refused to answer due to privacy concerns. I am planning to do a mix of EDA, Clustering and Classification https://archive.ics.uci.edu/ml/datasets/Cervical+cancer+%28Risk+Factors%29

References

- A Grammar of Data Manipulation dplyr. (2020). Retrieved February 7, 2021, from https://dplyr.tidyverse.org/
- Fontham, E. T. H., Wolf, A. M. D., Church, T. R., Etzioni, R., Flowers, C. R., Herzig, A., ... Smith, R. A. (2020). Cervical cancer screening for individuals at average risk: 2020 guideline update from the American Cancer Society. *CA: A Cancer Journal for Clinicians*, 70(5), 321–346. https://doi.org/10.3322/caac.21628
- Tidy Messy Data tidyr. (2020). Retrieved February 7, 2021, from https://tidyr.tidyverse.org/