
Lecture 10: Tidy Data Principles

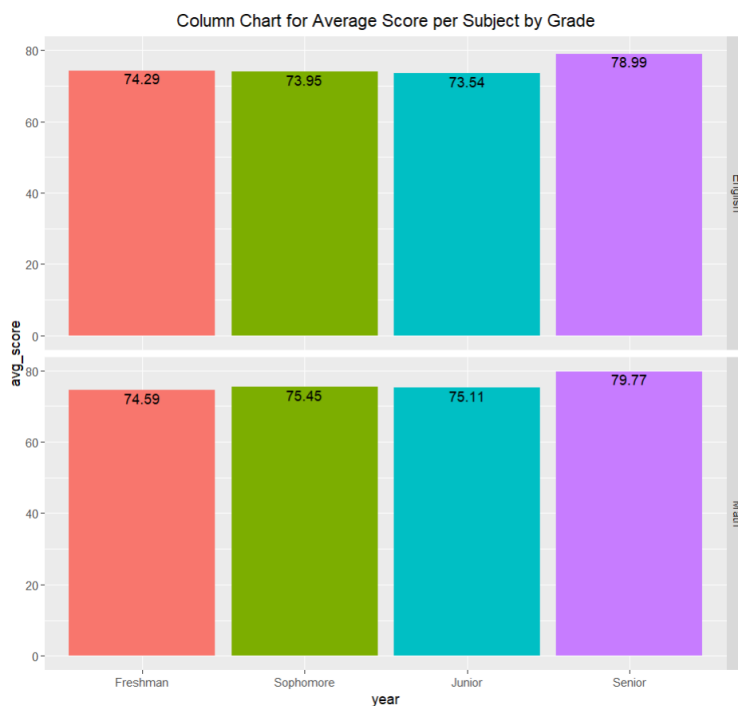
Problem Set for Lecture 10

These questions will assist in bolstering your understanding of the material in Lecture 10. Emphasis should be placed on having the correct code/output as well as communication. The deliverable should be a knitted RMarkdown document to pdf without any code running off the page. You will be able to present these in office hours as an oral assessment interview along with the problem set for Lecture 8 and Lecture 9.

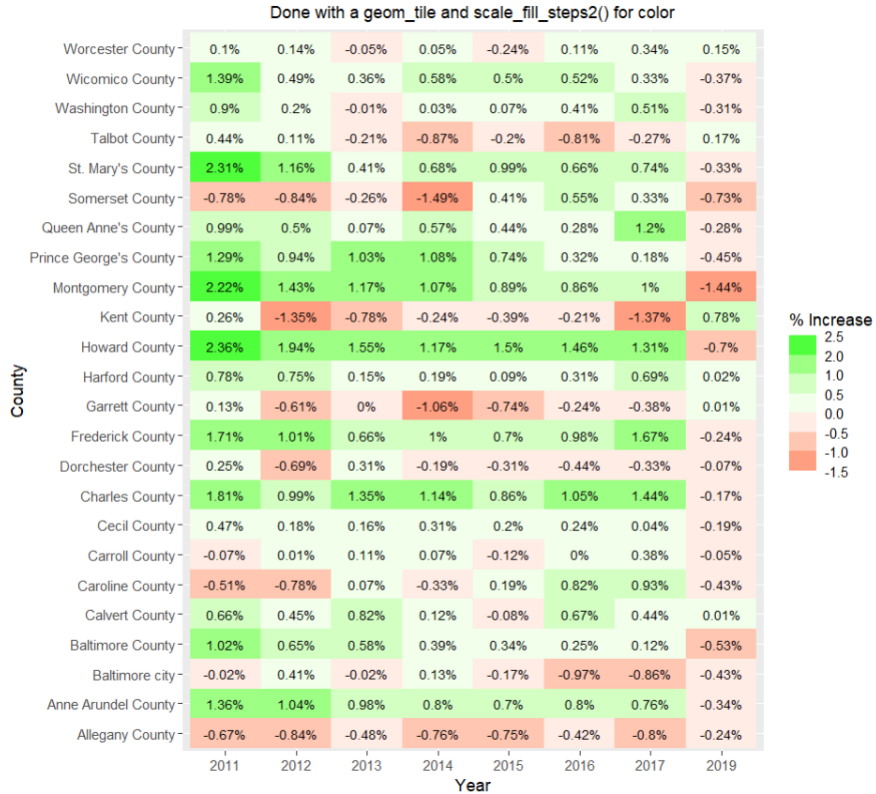
***** First, make sure you have read the Lecture Material's writeup *****

1. Using the `course_scores` dataset in the `MSMU` library, tidy the data into a form that allows you to analyze student performance across subject and grade level. Then answer the following questions.
 - (a) Compute the average numeric score for each subject across all grade levels.
 - (b) Compute the average numeric score for each grade level across all subjects.
 - (c) Compute the average numeric score for each subject-grade combination.
 - (d) Determine the distribution of letter grades for each subject.
 - (e) Compute a GPA for each student using the scale F=0, D=1, C=2, B=3, A=4 across all of their letter-grade observations.
 - (f) If a student needs at least a 2.0 GPA to graduate, how many students would graduate? (51)
 - (g) Identify the top 5 students with the largest absolute difference between their average English score and average Math score. (Student 39, 14, 44, 18, 55)
 - (h) Create the visualization (found below) that compares numeric scores across grade levels and subjects. Write 2–3 sentences describing what the graph suggests.
2. Using the `county_complete` dataset (only looking at the `state`, `name`, and any column containing "pop") in the `usdata` library, tidy the population data so that year is treated as a variable rather than being embedded in the column names. Then answer the following questions.
 - (a) Determine the total U.S. county-level population in each available year.
 - (b) For each county, determine the absolute population change from 2000 to 2019 and identify the 10 counties with the largest increase. (Maricopa County)
 - (c) For each county, determine the percent population change from 2000 to 2017 and identify the 10 counties with the largest percent increase. (Pinal County– 141%)
 - (d) How many counties had population decline from 2010 to 2019?

-
- (e) For each state, determine how many of its counties increased in population from 2010 to 2019 and how many decreased.
- (f) Create the visualization (found below) using your tidied population data to show how the counties in Maryland have changed over time. Write 2–3 sentences describing what the graph suggests.
3. Using the `billboard` dataset in the `tidyr` library, tidy the data so that week number is stored as a variable and chart position is stored as a value. Also note that a low `position` is better (considered higher on the chart). Then answer the following questions.
- (a) For each song, determine the best chart position it ever reached.
- (b) Which artist had the most songs appear in the dataset? (Jay-Z with 5)
- (c) Determine how many weeks each song remained on the chart. (ie was not NA)
- (d) Among songs that reached the top 10 at some point, which song stayed on the chart the longest?
- (e) For the artists with multiple songs on the list, determine which one has on average the best chart position across all of their songs. (Christina Aguilera and Destiny's Child with average top spot of 1.67)
- (f) Determine which song had the largest improvement from its first recorded chart position to its best chart position. (He Loves U Not)
- (g) Create the visualization (found below) using your tidied `billboard` dataset that shows how chart position changes from week to week for the 5 songs that stayed on the chart the longest. Write 2–3 sentences describing what your graph suggests.
-



Percent Increase of Maryland Counties Over Time



Plot of Position for Songs on the chart the longest

