

Exam 02 (Open Notes)

DSST289: Introduction to Data Science

Erik Fredner

2024-10-23

Table of contents

Deadline	1
Honor Pledge	1
Setup	2
Instructions	2
Data: storms	3
Features	3
Metadata	4
Questions	5
Max wind speed over hurricane lifetime	5
Average speed by hurricane category	5
Days of the year with midnight hurricanes	5
Last letter of the year	6
Average max storm wind speed by storm name gender	6
Storms by status	7
Max wind speed by first letter per year	7
Trend in Average Wind Speed Over Time by Storm Name Gender	7

Deadline

Monday, October 28 *before* the start of class.

Honor Pledge

“I pledge that I will neither give nor receive unauthorized assistance during the completion of this work.”

For this exam, you may use class notes, notebooks, and slides. Any other resource (e.g., non-class websites, ChatGPT, etc.) is unauthorized.

Signature (type your full name):

UR email:

Section start time:

Setup

1. Navigate to Blackboard > Course Documents > Exams > Exam 02
2. Download this exam from Blackboard: `exam02_open.qmd`
3. Download the data from Blackboard:
 - `storms.csv`
 - `storm_gender.csv`
 - `storm_codes.csv`
4. Move the exam to the `nb` folder in your DSST289 folder, just as we do when working on new notebooks in class: (`...DSST289/nb/exam02_open.qmd`)
5. Move the data to `...DSST289/data/`.

Instructions

1. Although there are multiple ways of producing the results requested in each question, I expect to see you use patterns and techniques that we have discussed.
2. If you are unable to complete a question, explain your attempt to maximize partial credit.
3. If you encounter R or RStudio errors that you cannot resolve on your own, contact me ASAP. I can help you with configuration issues, but will not help you answer questions. If you run into any issues with your personal computer, use the computers in the library to complete the exam.
4. When you have finished the exam, **render** your `.qmd` file to `.pdf`. If rendering fails, upload the `.qmd` file.
5. Go to Blackboard > Assignments > Exam 02 (open notes). Upload your **rendered** document there.

Data: storms

The data for this exam consists of information about tropical storms in the Atlantic Ocean between 1950 and 2020.

```
library(tidyverse)

storms <- read_csv("../data/storms.csv")

storms |>
  slice_sample(n = 5)
```

```
# A tibble: 5 x 10
   year name   letter   doy  hour   lat   lon status category  wind
  <dbl> <chr>   <chr>   <dbl> <dbl> <dbl> <dbl> <chr>      <dbl> <dbl>
1  2001 Karen    K     287    12  39.3 -63.9 TS         0     60
2  1988 Gilbert  G     259    12  21.9 -91.7 HU         2     85
3  1999 Floyd   F     262     6  48.5 -52.5 EX         0     35
4  2008 Marco   M     280    12  18.9 -93.7 TS         0     40
5  1998 Jeanne  J     265    18  13.1 -25.2 HU         1     65
```

`storms` contains one row for each time a particular storm was measured. Storms are generally measured once every six hours.

Features

Variable	Description
<code>year</code>	The year in which the storm was recorded
<code>name</code>	A common name for the storm. Names can be reused for different storms in different years.
<code>letter</code>	The first letter of the name; storms are (usually) named in alphabetical order
<code>doy</code>	The day of the year (1-365) of the record
<code>hour</code>	The hour of the day (0-23) of the record in Eastern time
<code>lat</code>	Latitude of the record in degrees
<code>lon</code>	Longitude of the record in degrees
<code>status</code>	A two-digit status code of the storm system; see <code>storm_codes.csv</code> for full names
<code>category</code>	For hurricanes (<code>status == "HU"</code>), a number giving the category of the storm from 0-5
<code>wind</code>	The observed sustained wind speed in miles per hour

Metadata

In addition to the main `storms` table, there are two metadata tables. `storm_gender` provides an automatically determined estimate of whether storm's name is male or female. Its `prob` column gives a confidence score for the accuracy of the `gender` determination. A higher score indicates a higher confidence.

```
storm_gender <- read_csv("../data/storm_gender.csv")

storm_gender |>
  filter(prob < 1) |>
  arrange(desc(prob)) |>
  slice_head(n = 3)
```

```
# A tibble: 3 x 3
  name  gender  prob
<chr> <chr>  <dbl>
1 Anna  female  0.999
2 Grace female  0.999
3 Julia female  0.999
```

```
storm_gender |>
  arrange(desc(prob)) |>
  slice_tail(n = 3)
```

```
# A tibble: 3 x 3
  name  gender  prob
<chr> <chr>  <dbl>
1 Nana  female  0.688
2 Charley female  0.642
3 Joan  female  0.510
```

There is a column in `storms` called `status` that describes the type of storm with a two letter code. `storm_codes` provides a full name for each of these codes:

```
storm_codes <- read_csv("../data/storm_codes.csv")
storm_codes
```

```
# A tibble: 9 x 2
  status status_name
```

	<chr>	<chr>
1	TD	tropical depression
2	TS	tropical storm
3	HU	hurricane
4	EX	extratropical cyclone
5	SD	subtropical depression
6	SS	subtropical storm
7	LO	low
8	WV	tropical wave
9	DB	disturbance

Questions

Max wind speed over hurricane lifetime

Output a table with one row for each storm in the data set that provides the maximum wind speed the storm achieved over its lifetime.

Average speed by hurricane category

Hurricanes get assigned one of six different categories based on their sustained wind speed. When a hurricane is covered on the news, you may hear it described as a “Category 3” storm, for example.

Create a new data set that has one row for each hurricane category in each year that shows the average wind speed of hurricanes in that category during that year.

Tip

Not every storm in `storms` is a hurricane.

Using the table you just created, create a line plot with a points layer showing the average wind speed over time by hurricane category. Color the points and lines by hurricane category using a colorblind-friendly scale. Label the axes and legend.

Days of the year with midnight hurricanes

Create a new table where the unit of observation is the day of the year. Count the total number of hurricanes observed at midnight on each day of the year.

Tip

There are days of the year without hurricanes observed at midnight. You do **not** need rows for those days.

Using the table you just created, make a bar plot that shows the number of hurricanes on each day of the year, with the day of the year on the x-axis and the number of hurricanes observed on that day on the y-axis.

Add a layer on top of that bar plot that highlights the days of the year with the *median* number of hurricanes observed at midnight. Color the bars for the median days "#440154" and the non-median days "#fde725". Label the axes and title the plot.

Last letter of the year

Storms are named in alphabetical order, with the first storm of the year starting with the letter A, the second with the letter B, and so on.

Output a table with two columns: **letter** and **n**. **n** should indicate the number of years in which each letter was the *last* letter used to name a storm in that year. For example, in 1972 and 1983, the last storm of the year started with "D."

Note

The data has been filtered to exclude some storms, such as those with Greek letters, so do not expect these results to exactly match other sources.

Average max storm wind speed by storm name gender

Create a table with two rows showing the average maximum wind speeds of storms with male or female names.

Note

By "average maximum," I mean that you should first compute each storm's maximum wind speed, *then* take the average of these maximum values.

Not all storm names appear in **storm_gender**. Write code that returns an alphabetical list of the unique names that appear in **storms** but do not have a **storm_gender**.

Storms by status

`storms` contains codes describing the status of the storm at the point of observation. Produce a table containing the number of distinct storms observed with each status. This table should have two columns: One containing the full name of the storm status (*not* the abbreviated code), and the other containing the count of observed storms, sorted by the most frequent status name.

Tip

The same storm can have multiple different statuses across different observations.

Max wind speed by first letter per year

Use the `storms` data set to create a table, the *first four* rows and columns of which look like the following:

year	A	B	C
1950	110	90	95
1951	80	50	115
1952	85	95	105
1953	60	80	140
1954	95	50	100

Your table must contain one column for all of the years and one column for each letter in the `storms` data set.

The values in each cell other than `year` should represent the max wind speed attained by a storm in that year with a name starting with the corresponding letter. For example, the 1950 storm named Baker had a maximum wind speed of 90.

Trend in Average Wind Speed Over Time by Storm Name Gender

Determine if there is a trend in the average wind speed of storms over time by the gender of the storm's name.

First, calculate the average wind speed for all storms in each year by storm name gender. Then, create a scatter plot of the average wind speed per year by gender, and add linear trend lines *within* each group.