



Descriptive Statistics

Objective: Use descriptive statistics to gather information about a data set.

In this lab, you will complete a series of exercises to calculate various summary statistics of the `dsfda.ht_daily_metrics` data set.

```
1 %run "../Includes/Classroom-Setup"
```

Part 1: Measures of Central Tendency

Exercise 1: Mean

Use the `AVG` SQL function to calculate the average resting heartrate, as an integer, for each lifestyle.

 **Hint:** If needed, check out the [Spark SQL documentation on avg](#).

```
1 %sql
2 -- ANSWER
3
4 SELECT CAST(AVG(resting_heart_rate) AS int) as avg_resting_heart_rate, lifestyle
5 FROM dsfa.ht_daily_metrics
6 GROUP BY lifestyle
```

Exercise 2: Median

Use the `percentile` SQL function to calculate the median resting heartrate, as an integer, for the "Athlete" lifestyle.

Note that this code block takes longer to run; this is because calculating the median requires shuffling across worker nodes.


 Hint: If needed, check out the [Spark SQL documentation on percentile](#).

```
1 %sql
2 -- ANSWER
3
4 SELECT CAST(percentile(resting_heartrate, 0.5) AS int) AS median_resting_heartrate
5 FROM dsfda.ht_daily_metrics
6 WHERE lifestyle = 'Athlete'
```

Exercise 3: Mode

Calculate the mode, as an integer, of the resting heartrate for each lifestyle.

 Hint: Cast the `heartrate` to the nearest integer first. You can use the [Spark SQL documentation](#) on `cast` for reference.


 **Hint:** Try to find the `mode` Spark SQL function in the documentation on your own.


```
1 %sql
2 -- ANSWER
3
4 SELECT COUNT(*) AS count, CAST(resting_heartrate AS int) AS resting_heartrate
5 FROM dsfa.ht_daily_metrics
6 GROUP BY resting_heartrate
7 ORDER BY count DESC
```

Part 2: Measures of Dispersion

Exercise 4: Standard Deviation

Calculate the standard deviation of the resting heartrate for each lifestyle.

 **Hint:** Cast the heartrate to the nearest integer first.

 **Hint:** Use the Spark SQL documentation to determine which function computes the standard deviation.


Cmd 13

```
1 %sql
2 -- ANSWER
3
4 SELECT CAST(STD(resting_hearttrate) AS int) as std_resting_hearttrate, lifestyle
5 FROM dsfda.ht_daily_metrics
6 GROUP BY lifestyle
```

Cmd 14

Exercise 5: Interquartile Range

Calculate the interquartile range of resting heartrate, as an integer, for each lifestyle.

 **Hint:** Recall that the interquartile range is the difference between the 75th percentile and the 25th percentile.

 **Hint:** Refer to previous exercises in this lab to determine which Spark SQL function can be used to compute the value at a given percentile.

Cmd 15

```
1 %sql
2 -- ANSWER
3
4 SELECT CAST(percentile(resting_hearttrate, 0.75) - percentile(resting_hearttrate, 0.25) AS int) AS iqr_resting_hearttrate
5 FROM dsfda.ht_daily_metrics
6 GROUP BY lifestyle
```

Cmd 16

Congratulations! You've completed the Descriptive Statistics Lab.

If you have any trouble with any of this course's labs, be sure to check out the solutions to the labs in the Solutions folder.

Cmd 17

© 2021 Databricks, Inc. All rights reserved.

Apache, Apache Spark, Spark and the Spark logo are trademarks of the [Apache Software Foundation](#).

[Privacy Policy](#) | [Terms of Use](#) | [Support](#)

Shift+Enter to run