

FIT 5202 Assignment 1 Feedback Sheet						
Student Name: ISOBEL ROWE						
Marked By: David _C.						
Part A : Working with RDDs and DataFrames						
Tasks		Criteria	Yes	Partial	No	Comments
1.1 Data Preparation and Loading	1	No of processors and title of application	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Created SparkSession and SparkConf	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	- No SparkConf implemented - Additional configuration unnecessary (e.g. mongodb connectors)
	2	Correctly imported the RDDs for Units	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	3	Correctly imported the RDDs for Crashes	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	4	Headers removed, count and first 10 rows displayed for both RDDs	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	- Rows displayed would be better to show in tabular format
1.2 Data Partitioning in RDD	1	Correct number of partitions displayed	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Default partition strategy answered correctly	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	- Correct partition strategy, but not explaining the reason of 5 partitions
	2	a. Key Value Pair RDDs created correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		b. Hash function correctly defined	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		b. Partitioning implemented correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		c. Number of records correctly displayed in each partition	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		c. Skewness in partition discussed correctly	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	- No attempt
1.3 Query/Analysis	1	Valid drivers filtered out	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	- Filtering criteria is not proper (e.g. pedestrian not filtered out), suggest to examine the data to identify the wrong data value
		Average Age calculated correctly	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	- Incorrect values due to missing filter
	2	Valid vehicle filtered out	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Oldest and newest vehicles calculated correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	- There are multiple vehicles with max and min year. Only showed 1 for each of them
2.1 Data Preparation and Loading	1	Data loaded into dataframes correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	

<b>2.1 Data Preparation and Loading</b>	<b>2</b>	Schema correctly displayed	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
<b>2.2 Query/Analysis</b>	<b>1</b>	Filters implemented correctly and data displayed	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	<b>2</b>	10 crash events with highest casualties correctly displayed	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	- Incorrect results because 'Total Cas' was sorted alphabetically rather than numerically. Column should have been converted to int or float
	<b>3</b>	Total fatalities for each crash type displayed	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Unlicensed driver filter implemented	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		join, group by and aggregation correctly implemented	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Results correctly displayed	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
<b>2.3 Severity Analysis</b>	<b>1</b>	Group by and count correctly implemented	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Most common severity level answered correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
	<b>2</b>	a. Positive on drugs only calculated correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		b. Positive on alc. only calculated correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		c. Positive on both calculated correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		d. Negative for both calculated correctly	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Brief explanation of the observation	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	- No attempt in notebook
<b>2.4 RDDs vs DataFrame vs SparkSQL</b>	<b>1</b>	Correct implementation for RDD	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	- No attempt
		Correct implementation for DataFrame	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	- No attempt
		Correct implementation for Spark SQL	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	- No attempt
	<b>2</b>	Correct implementation for RDD	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	- No attempt
		Correct implementation for DataFrame	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Correct implementation for Spark SQL	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
		Discussion of the performance differences	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	- No attempt in notebook

<b>Qualitative Aspect</b>	Organization of tasks in jupyter notebook Adherence to python standards Use of appropriate comments, output readability	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	- To increase the readability in the notebook, you should only show the relevant columns, all columns would make it messy and decrease the readability
<b>Part B : Pre-recorded Video Presentation</b>					
1. RDD Partitioning	Partitioning strategy, data distribution, data skewness, approaches to manage skew	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	- How to manage skewness was not fully explained
2. Crash Severity Analysis	Correctness of the observations based on the bar graph visualization	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
3. RDDs vs DataFrame vs Spark SQL	Performance findings , explanation for df/sql being faster than RDD	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	- DataFrame and SQL queries run faster thanks to optimized execution plans and custom memory management. Suggest doing more reading on those and taking a look at Spark UI to see what happens
<b>Qualitative Aspect</b>	Overall quality of presentation delivery (video/audio clarity) content (use of graphs/quality of slides)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	- Suggest adding interaction with the slide, so that it could be easier for the audience to follow your talk - More graphs for slides would be better to be added
<b>Final Grade</b>		Late Submission		2	<b>C</b>