

Question 1		
/12	(a)	<b>1 point</b> per type of bias named <b>2 points</b> per correct definition
/4	(b)	<b>1 point</b> per correct pairing
/4	(c)	<b>2 points</b> for a reasonable solution <b>2 points</b> for naming a reasonable stage
Question 2		
/6	(a)	<b>3 points</b> per valid benefit [productivity/iteration/data in memory/good for EDA/etc]
/2	(b)	<b>2 points</b> for mentioning a SQL, queries, tables, etc.
/4	(ci)	<b>4 points</b> for a mentioning in some way replication of data
/4	(cii)	<b>4 points</b> for mentioning storing lineage graphs, or recomputation, or something similar
/4	(d)	<b>2 points</b> for stating that transformations are lazy, not evaluated immediately, etc <b>2 points</b> for stating that actions are evaluated immediately, trigger computation, etc
/5	(e)	<b>1 point</b> for naming a valid transformation [map, filter, union, etc] <b>2 points</b> for a DAG which shows at least 2 RDDs <b>2 points</b> for having multiple partitions per RDD, with one child per parent
/5	(f)	<b>1 point</b> for naming a valid transformation [reduceByKey, groupByKey, etc] <b>2 points</b> for a DAG which shows at least 2 RDDs <b>2 points</b> for having multiple partitions per RDD, with multiple children per parent
Question 3		
/10	(a)	<b>1 point</b> per RDD in one iteration [start, unlabeled, unlabeled, middle] <b>1 point</b> per correct transformation label in one iteration [flatMap, map, reduceByKey] <b>2 points</b> for having the 2nd iteration represented <b>1 point</b> for having the fullOuterJoin between middle and result
/10	(b)	<b>4 points</b> for correctly stating 3 stages ( <b>2 points</b> for stating 2 stages) <b>3 points</b> per boundary drawn through the reduceByKey transformations
/12	(c)	<b>4 points</b> for identifying .reduceByKey on line 6 <b>4 points</b> for identifying .flatMap (line 4) or .map (line 5) <b>4 points</b> for identifying .collect on line 12
/4	(d)	<b>4 points</b> for stating that it runs 1 job
/4	(e)	<b>4 points</b> for stating that it implements Word Count

		Question 4
/6	(a)	<b>1 point</b> per classifier mentioned [k-NN, Naive Bayes, Logistic Regression] <b>1 point</b> per correct categorization [structural, statistical, statistical]
/4	(b)	<b>1 point</b> for stating that words are the features in spam classification <b>1 point</b> for stating that urls visited are the features in ad click probability <b>2 points</b> for stating that k-NN works better for low dimensionality (or Curse of Dim)
/2	(c)	<b>2 points</b> for mentioning that it works well with large amounts of data
/3	(d)	<b>3 points</b> for mentioning we run on each class, and pick the most likely
		Question 5
/10	(a)	<b>10 points</b> for getting the correct answer (or close with just math errors) <b>5 points</b> for having some correct steps/formulas
/10	(b)	<b>10 points</b> for getting the correct answer (or close with just math errors) <b>5 points</b> for having some correct steps/formulas
/6	(c)	<b>6 points</b> for getting the correct answer (or close with just math errors) <b>3 points</b> for having some correct steps/formulas
/4	(d)	<b>2 points</b> for mentioning Laplace Smoothing <b>2 points</b> for mentioning getting more data (or other reasonable answer)
/5	(e)	<b>5 points</b> for a reasonable explanation that points to a specific type of bias, or issue <b>2 points</b> for just mentioning bias but not anything specific
		Question 6
/2	(a)	<b>0.5 points</b> (rounded down) per correct answer [volume, velocity, variety, veracity]
/2	(b)	<b>1 point</b> per correct answer [random process generating data, random sampling]
/2	(c)	<b>2 points</b> for stating that odds of machine going down approaches 1.0 as number of machines increases
/2	(d)	<b>2 points</b> for any feasible answer
/2	(e)	<b>1 points</b> per correct answer [sample, distribute computation, Spark, MR]