



Used Car Prediction Analysis

Michael Miranda
Roger Ramirez
Andrew Pang

CIS 5560 Intro to Big Data Science | Professor Jongwook Woo



Github for Code Files

<https://github.com/Mgratricz/CIS-5560-Used-Cars-Prediction-Analysis>



Project Summary

In recent years, rising vehicle prices, limited availability of new cars, and ongoing supply chain issues have made buying a used car a more common and complex.

Our goal is to create models, using different types of algorithms learned in class, that can estimate the price of a used car based on its individual attributes and what primary attributes greatly affect a car with an accident.



About the Dataset



Data Source: US Used Cars Dataset

- <https://www.kaggle.com/datasets/ananyamittal/us-used-cars-dataset>
- Data Size: ~2.5GB

Dataset Preview

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
1	vin	back_legroo	bed	bed_height	bed_length	body_type	cabin	city	city_fuel_ec	combine_fuel	daysonmarket	dealer_zip	description	engine_cylitr	engine_displ	engine_type	exterior_col	fleet
2	KL7CJKS84LB0	35.7 in				SUV / Crossover		Raynham	26	2767	All specials inclu	183	2767	All specials inclu	1400	I4	Summit White	
3	1C4PJMDX0LD4	40.3 in				SUV / Crossover		Rochester Hills	19	130	48307	110	48307	110	3200	V6	Diamond Black Crystal Pearlcoat	
4	3N1AB7AP5KLE	37.4 in				Sedan		Tulsa	29	44	74133	44	74133	44	1800	I4	Gray	TRUE
5	3N1AB8BV2LY2	37.4 in				Sedan		Atlanta	29	41	30360	41	30360	41	2000	I4	Gun Metallic	
6	1GNKRHKD0HJ	36.8 in				SUV / Crossover		Beavercreek Twp	15	48	45434	48	45434	48	3600	V6	Mosaic Black Metallic	FALSE
7	3CZRU6H37LM	39.3 in				SUV / Crossover		Lancaster	27	71	1523	71	1523	71	1800	I4	MODERN STE	
8	5TFUY5F11X9	34.7 in	--		78.7 in	Pickup Truck		Langhorne	13	17	19047	17	19047	17	5700	V8	Super White	
9	1J4FA9S3P3C	35 in				SUV / Crossover		Austin	15	4	78757	4	78757	4	4000	I6	Red	FALSE
10	1GCRYDED1LZ	35.2 in	--		79.4 in	Pickup Truck		Big Rapids	16	9	49307	9	49307	9	5300	V8	Northsky Blue Metallic	
11	2C4RDGEG9HF	36.5 in				Minivan		Asheville	17	33	28805	33	28805	33	3600	V6 Flex Fuel Vel	Billet Clear Coat	TRUE
12	JA4MW51SX5K	35.2 in				SUV / Crossover		Cumming	13	552	30028	552	30028	552	3800	V6		FALSE
13	1FTEW1E53KF	43.6 in	--		67.1 in	Pickup Truck		Pennsville	16	580	8070	580	8070	580	5000	V8 Flex Fuel Vel	Oxford White/Met	FALSE
14	2HGFA1F96AH	34.6 in				Sedan		Maryville	25	110	37801	110	37801	110	1800	I4	Gray	FALSE
15	5FNFY5H50MB	38.4 in				SUV / Crossover		Jefferson City	20	15	65109	15	65109	15	3500	V6	Black	
16	3GKALPEXLL	39.7 in				SUV / Crossover		Plymouth	22	200	48170	200	48170	200	2000	I4	Summit White	
17	JTMZFREV2ED	37.2 in				SUV / Crossover		Philadelphia	24	5	19111	5	19111	5	2500	I4	Shoreline Blue F	FALSE
18	1GKKVREDXBJ	36.9 in				SUV / Crossover		Sacramento	16	377	95822	377	95822	377	2600	V6	White	TRUE
19	1FMCU9G61LU	40.7 in				SUV / Crossover		Ukiah	26									
20	3MVDMBEM5LJ	36.3 in				SUV / Crossover		Asheville	24									
21	1G1ZB5T4LF1	38.1 in				Sedan		Andover	29									
22	1FMCU0G62LU	40.7 in				SUV / Crossover		San Bernardino	27									
23	1GYKPCRS5LZ	39.1 in				SUV / Crossover		Las Vegas	18									
24	3GCUYDED9LC	43.4 in	--		69.9 in	Pickup Truck		Clarksburg	16									
25	2G1115SL7E92	39.8 in				Sedan		Los Angeles	21									
26	2GNAXKEV1L6	39.9 in				SUV / Crossover		Brunswick	26									
27	19XFB2F59CE0	36.2 in				Sedan		Dallas	28									
28	3GTP1NEC7HG	40.9 in	--		69.3 in	Pickup Truck		Sheboygan										
29	1FTER1EH9LLA	30.4 in	--		72.8 in	Pickup Truck		Seaford	21									
30	1VWCN7A30EC	39.1 in				Sedan		Watertown	30									
31	2B3KA43G08H3	40.2 in				Sedan		Modesto										
32	3GTP9EEL4LG	43.4 in	--		69.9 in	Pickup Truck		Austin	15									
33	19UUB2F42KAC	34.5 in				Sedan		Gilbert	20									
34	2GCUYDE9H6	40.9 in				Pickup Truck		Austin	16									

%pyspark

Random Forest Regression

4. Assemble + scale

```
from pyspark.ml.feature import VectorAssembler, StandardScaler
```

```
feature_columns = [
    "city_fuel_economy", "highway_fuel_economy", "daysonmarket",
    "engine_displacement", "horsepower", "mileage", "seller_rating",
    "year", "make_indexed", "model_indexed",
    "torque", "engine_cylinders", "power", "wheelbase", "width"
]
```

```
assembler = VectorAssembler(
    inputCols=feature_columns,
    outputCol="assembled_features"
```

```
)
```

The Goal

To create accurate models that can be used to predict the *price* and *accident* variables using Hadoop Distributed File System and PySpark.



Experimental Specifications



Hadoop Version	3.1.2
HDFS Configuration	390.17 gb Capacity
Total Nodes	3
Spark Version	3.0.2

Workflow Overview

01

Pre-Processing



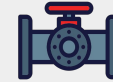
02

Splitting the Data



03

Assembler & Pipeline



04

Model



05

Train Validation & Cross Validation



06

Evaluation



01

Pre-Processing

```
df.printSchema()

root
 |-- vin: string (nullable = true)
 |-- back_legroom: string (nullable = true)
 |-- bed: string (nullable = true)
 |-- bed_height: string (nullable = true)
 |-- bed_length: string (nullable = true)
 |-- body_type: string (nullable = true)
 |-- cabin: string (nullable = true)
 |-- city: string (nullable = true)
 |-- city_fuel_economy: double (nullable = true)
 |-- combine_fuel_economy: string (nullable = true)
 |-- daysonmarket: integer (nullable = true)
 |-- dealer_zip: string (nullable = true)
 |-- description: string (nullable = true)
 |-- engine_cylinders: string (nullable = true)
 |-- engine_displacement: double (nullable = true)
 |-- engine_type: string (nullable = true)
```

```
for name, dtype in [
    ("city_fuel_economy", "double"),
    ("highway_fuel_economy", "double"),
    ("daysonmarket", "int"),
    ("engine_displacement", "double"),
    ("horsepower", "double"),
    ("mileage", "double"),
    ("seller_rating", "double"),
    ("year", "int"),
    ("price", "double"),
    ("torque", "double"),
    ("engine_cylinders", "double"),
    ("power", "double"),
    ("wheelbase", "double"),
    ("width", "double"),
]:
    data = data.withColumn(name, col(name).cast(dtype))
```

```
%pyspark
# 3. Index categorical features
from pyspark.ml.feature import StringIndexer

make_indexer = StringIndexer(inputCol="make_name", outputCol="make_indexed")
model_indexer = StringIndexer(inputCol="model_name", outputCol="model_indexed")

data = make_indexer.fit(data).transform(data)
data = model_indexer.fit(data).transform(data)
```

Took 3 sec. Last updated by anonymous at May 03 2025, 9:11:17 PM.

- Essential data cleaning processes, for example, converting data types, indexing categorical variables and handling missing values.

01

Feature Importance

```
%pyspark
# 9. Feature importance
importances = rf_model.featureImportances.toArray()
for feature, score in zip(feature_columns, importances):
    print(f"{feature}: {score:.4f}")
```

```
city_fuel_economy: 0.0658
highway_fuel_economy: 0.0716
daysonmarket: 0.0476
engine_displacement: 0.0757
horsepower: 0.2187
mileage: 0.1915
seller_rating: 0.0271
year: 0.1550
make_indexed: 0.0846
model_indexed: 0.0624
torque: 0.0000
engine_cylinders: 0.0000
power: 0.0000
wheelbase: 0.0000
width: 0.0000
```

Feature Importance from Random Forest

	feature	coefficient
0	city_fuel_economy	-14833.229917
4	engine_type_index	-14686.785567
7	maximum_seating_index	-11350.793752
6	fuel_type_index	-11280.454991
1	highway_fuel_economy	-9386.707043
3	body_type_index	-8176.955298
10	year_index	6893.039711
5	make_name_index	-3601.329178
8	model_name_index	-1646.974202
9	wheel_system_index	-483.822212
2	mileage	-468.488020

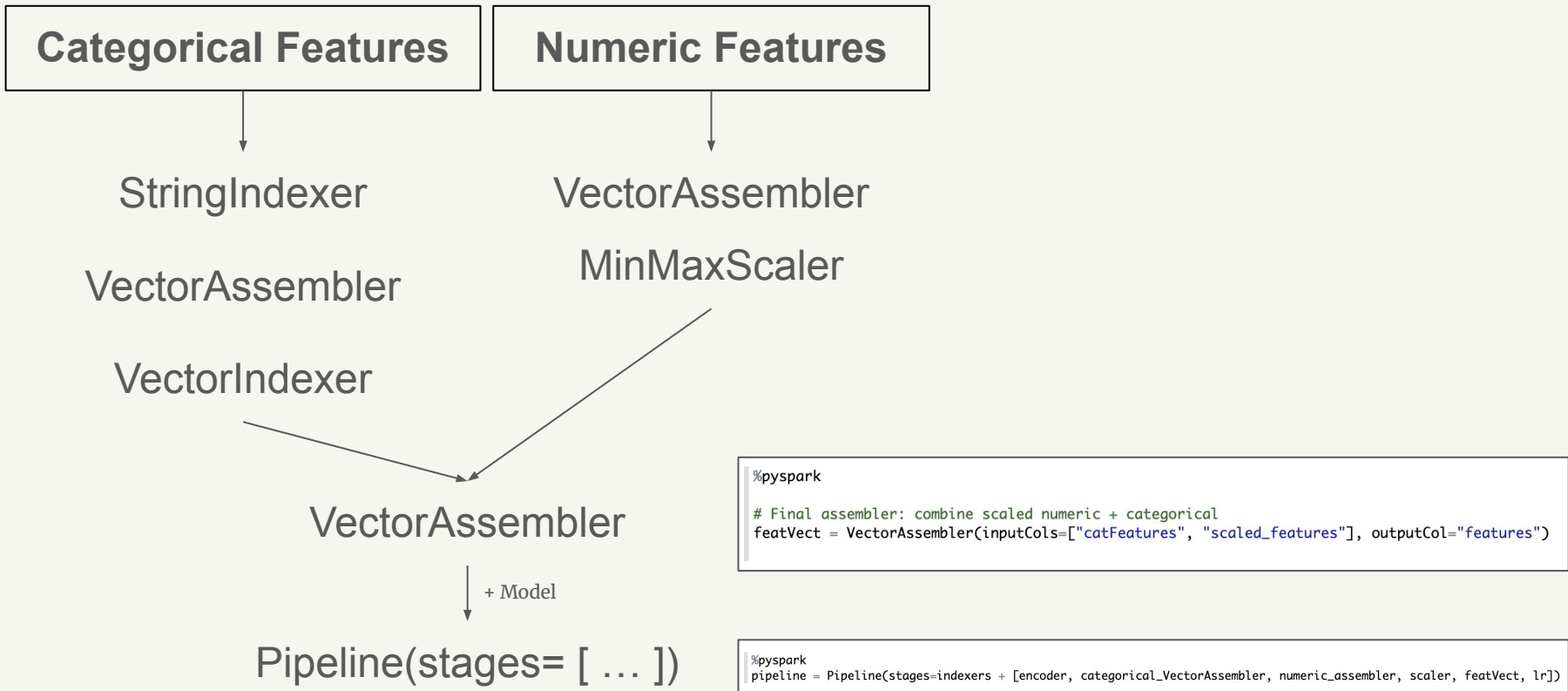
Coefficients from Linear Regression

02

Splitting the Data

```
%pyspark

splits = data.randomSplit([0.7, 0.3])
train = splits[0]
test = splits[1]
train_rows = train.count()
test_rows = test.count()
print("Training Rows:", train_rows, "\nTesting Rows:", test_rows)
```



#1 Linear regression

```
from pyspark.ml.regression import LinearRegression
```

#2 Random forest regression

```
from pyspark.ml.regression import RandomForestRegressor
```

#3 Random forest classifier

```
from pyspark.ml.classification import RandomForestClassifier
```

#4 Binomial logistic regression

```
from pyspark.ml.classification import LogisticRegression
```

price

has_accidents

6A. Cross-Validation

```
rf = RandomForestRegressor(featuresCol="features", labelCol="price")
paramGrid = ParamGridBuilder() \
    .addGrid(rf.numTrees, [50]) \
    .addGrid(rf.maxDepth, [5, 10]) \
    .build()

evaluator = RegressionEvaluator(labelCol="price", predictionCol="prediction", metricName="rmse")
crossval = CrossValidator(estimator=rf,
                          estimatorParamMaps=paramGrid,
                          evaluator=evaluator,
                          numFolds=3)
```

```
%pyspark

tvsv = TrainValidationSplit(estimator=pipeline,
                           estimatorParamMaps=paramGrid,
                           evaluator=RegressionEvaluator(labelCol="price", predictionCol="prediction", metricName="rmse"))

start = time.time()
tvsvModel = tvsv.fit(train)
end = time.time()

tvsv_time

tvsv_predictions = tvsvModel.transform(test)

print(f"tvsvModel fit time: {end - start:.2f} seconds")

# EVAL

evaluator_rmse = RegressionEvaluator(labelCol="price", predictionCol="prediction", metricName="rmse")
evaluator_r2 = RegressionEvaluator(labelCol="price", predictionCol="prediction", metricName="r2")

tvsv_rmse = evaluator_rmse.evaluate(tvsv_predictions)
tvsv_r2 = evaluator_r2.evaluate(tvsv_predictions)

print("TVS RMSE =", tvsv_rmse)
print("TVS R2 =", tvsv_r2)
```

	Algorithm	Predictor Variables	RMSE	R2	Computing Time (Training)
Regression	#1 Linear Regression	price	CV: 9537.86 TVS: 9546.95	CV: 0.7755 TVS: 0.7751	CV: 731.42 seconds TVS: 348.63 seconds
Regression	#2 Random Forest Regression	price	CV: 5978.83 TVS: 5978.84	CV: 0.8824 TVS: 0.8824	CV: 1131.87 seconds TVS: 763.25 seconds
		Predictor Variables	Accuracy & AUC		Computing Time (Training)
Classification	#3 Random Forest Classifier	has_accidents	AUC: 0.6600 Accuracy: 0.8842 Precision: 0.7818 Recall: 0.8842 F1: 0.8299		CV: 1084.35 seconds TVS: 630.73 seconds
Classification	#4 Binomial Logistic Regression	has_accidents	AUC: 0.6830 Accuracy: 0.8829 Precision: 0.8186 Recall: 0.8829 F1: 0.8281		CV: 400.45 seconds TVS: 251.54 seconds

Thank you!

<https://github.com/Mgratricz/CIS-5560-Used-Cars-Prediction-Analysis>





Github for Code Files

<https://github.com/Mgratricz/CIS-5560-Used-Cars-Prediction-Analysis>



```
%pyspark

# 4. Assemble + scale
from pyspark.ml.feature import VectorAssembler, StandardScaler

feature_columns = [
    "city_fuel_economy", "highway_fuel_economy", "daysonmarket",
    "engine_displacement", "horsepower", "mileage", "seller_rating",
    "year", "make_indexed", "model_indexed",
    "torque", "engine_cylinders", "power", "wheelbase", "width"
]

assembler = VectorAssembler(
    inputCols=feature_columns,
    outputCol="assembled_features"
)
assembled_data = assembler.transform(data)

scaler = StandardScaler(
    inputCol="assembled_features",
    outputCol="features",
    withMean=True,
    withStd=True
)
scaler_model = scaler.fit(assembled_data)
final_data = scaler_model.transform(assembled_data)

final_data.select("features").show(3, truncate=False)
```

- Random Forest Regression (Price)
- Selecting the feature columns plus feature handling

```
+-----+
|features|
+-----+
|[-0.5668006073478629,-0.4167300982404232,0.46520517013174334,0.387011|
|[0.9316074980846417,1.1481263077387083,-0.29629995817986815,-0.849857|
|[0.6319258769981408,0.2092124641512294,-0.05722276673319944,-0.849857|
+-----+
only showing top 3 rows
```

Linear Regression (Regression)

```

25/05/05 08:49:37 INFO DAGScheduler: Job 971 is finished. Cancelling potential speculative or zombie tasks for this job
25/05/05 08:49:37 INFO VarnScheduler: Killing all running tasks in stage 171: Stage finished
25/05/05 08:49:37 INFO DAGScheduler: Job 972 finished: treeAggregate at Statistics.scala#85, took 10.45138 s
CV RMSE = 9148.763583364
CV R^2 = 0.7455660923866951
CV test expParam 0.1
Best elasticNetParam: 0.0
feature coefficient
4 engine_type_index -16486.679791
2 mileage -14588.299585
0 city_fuel_economy -13463.853028
6 fuel_type_index -13841.823882
7 max_min_cooling_index -9217.893999
1 highway_fuel_economy -8659.655222
3 body_type_index -7485.853971
5 year_index -3486.143465
8 model_name_index -3469.757615
9 wheel_system_index -1946.331446
6 make_name_index -1148.343558

--- Summary of All Results ---
Training Rows: 525,123
Testing Rows: 224,857

--- Summary of Cross-Validation Results ---
Cross-Validated RMSE: 9148.7743
Cross-Validated R2: 0.7455
CV Time: 796.74 seconds

--- Summary of Train-Validation-Split Results ---
TrainValidationSplit RMSE: 9180.9726
TrainValidationSplit R2: 0.7448
VFS Time: 309.4 seconds
25/05/05 08:49:37 INFO SparkContext: Invoking stop() from shutdown hook
25/05/05 08:49:37 INFO AbstractConnector: Stopped SparkUIWebUIConnector at http://10.4.1.18:8080/
25/05/05 08:49:37 INFO SparkUI: Stopped Spark web UI at http://10.4.1.18:8080/
25/05/05 08:49:37 INFO VarnClientSchedulerBackend: Interrupting monitor thread

```

Random Forest (Regression)

```

25/05/05 00:28:28 INFO DAGScheduler: Job 149 finished: treeAggregate at Statistics.scala:58, took 4.129032 s
Final RMSE: 6040.05
Final R2: 0.8800
city_fuel_economy: 0.0423
highway_fuel_economy: 0.0490
daysonmarket: 0.0384
engine_displacement: 0.0932
horsepower: 0.2167
mileage: 0.2171
seller_rating: 0.0406
year: 0.1081
make_indexed: 0.0792
model_indexed: 0.0555
torque: 0.0000
engine_cylinders: 0.0000
power: 0.0000
wheelbase: 0.0000
width: 0.0000
Cross-validated RMSE: 5978.83
Cross-validated R2: 0.8824
CV Time: 339.80 seconds
trainValidationSplit RMSE: 5978.67
trainValidationSplit R2: 0.8824
FIS Time: 61.21 seconds
Final RMSE: 6040.05
Final R2: 0.8800
25/05/05 00:28:28 INFO SparkContext: Invoking stop() from shutdown hook
25/05/05 00:28:28 INFO AbstractConnector: Stopped Spark#45393[HTTP/1.1,http://1.1.1.1][0.0.0.0:4041]
25/05/05 00:28:28 INFO SparkUI: Stopped Spark web UI at http://bigdata-uio0.s303291929060.trainingcn.oraclecn.com:4041
25/05/05 00:28:28 INFO yarnClientschedulereBackend: Interrupting monitor thread
25/05/05 00:28:28 INFO yarnClientschedulereBackend: Shutting down all executors
25/05/05 00:28:28 INFO yarnSchedulerBackendYarnDriverEndpoint: Asking each executor to stop shutdown
25/05/05 00:28:28 INFO yarnClientschedulereBackend: YARN client scheduler backend stopped
25/05/05 00:28:28 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
25/05/05 00:28:28 INFO MemoryStore: MemoryStore cleared
25/05/05 00:28:28 INFO BlockManager: BlockManager stopped
25/05/05 00:28:28 INFO BlockManagerMaster: BlockManagerMaster stopped
25/05/05 00:28:28 INFO OutputCommitCoordinatorOutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
25/05/05 00:28:28 INFO SparkContext: Successfully stopped SparkContext
25/05/05 00:28:28 INFO ShutdownHookManager: Shutdown hook called

```

Logistic Regression (Classification)

```

26/05/05 01:24:24 INFO DAGScheduler: Job 679 finished: collectAdeMap at MulticlassMetrics.scala:61, took 12.460891 s
Train Validation Split AUC (Logistic Regression): 0.8808
Train Validation Split Accuracy (Logistic Regression): 0.8836
Train Validation Split Precision (Logistic Regression): 0.8836
Train Validation Split Recall (Logistic Regression): 0.8838
Train Validation Split F1 (Logistic Regression): 0.8808

--- Summary of All Results ---
Training Rows: 74,400
Testing Rows: 32,197

--- Summary of Cross-Validation Results ---
CV AUC: 0.8808
CV Accuracy: 0.8838
CV Precision: 0.8836
CV Recall: 0.8838
CV F1: 0.8807
CV Time: 400.45 seconds

--- Summary of Train-Validation-Split Results ---
TVS AUC: 0.8808
TVS Accuracy: 0.8838
TVS Precision: 0.8836
TVS Recall: 0.8838
TVS F1: 0.8801
TVS Time: 291.84 seconds

26/05/05 01:24:24 INFO SparkContext: Invoking stop() from shutdown hook
26/05/05 01:24:24 INFO AbstractConnector: Stopped SparkRSocketClient[JVMT/2.1, (http://131.0.0.0:4042)
26/05/05 01:24:24 INFO SparkUI: Stopped Spark web UI at http://bigdatau01.smb0291299606.trainingynco.oracleevcn.com:4042
26/05/05 01:24:24 INFO YarnClientSchedulerBackend: Interrupting monitor thread
26/05/05 01:24:24 INFO YarnClientSchedulerBackend: Shutting down all executors
26/05/05 01:24:24 INFO YarnClientSchedulerBackend: Yarn client scheduler backend stopped. Asking each executor to shut down
26/05/05 01:24:24 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
26/05/05 01:24:24 INFO MemoryStore: MemoryStore cleared
26/05/05 01:24:24 INFO BlockManager: BlockManager stopped
26/05/05 01:24:24 INFO BlockManagerMaster: BlockManagerMaster stopped

```

[illegible]