NETAJI SUBHAS UNIVERSITY OF TECHNOLOGY



COMPUTER HARDWARE SOFTWARE WORKSHOP COCSC19

SUBMITTED BY:

KHUSHI KHARKE- 2021UCS1617

YASH CHAUHAN- 2021UCS1628

ADITYA YADAV- 2021UCS1640

SPARK TASK(UNIT-4)

Task 1: Explore RDD in spark

RDD stands for Resilient Distributed Dataset and it's one of the fundamental data structures in Apache Spark. RDDs are immutable, distributed collections of objects, which can be processed in parallel across a cluster. They provide fault tolerance, allowing computations to be rerun on failure by using lineage information to rebuild lost data.

Here are some key characteristics and features of RDDs:

1. Immutability:

RDDs are immutable, meaning once created, their content cannot be changed. This
immutability ensures consistency and simplifies parallel processing, as RDDs can be
shared across multiple operations without worrying about concurrent modifications.

2. Distributed Nature:

 RDDs are distributed across the nodes of a cluster, allowing for parallel processing of large datasets. This distribution enables Spark to leverage the computing power of multiple machines, making it suitable for handling big data workloads.

3. Resilience to Failures:

 RDDs are resilient to failures through lineage information. Spark keeps track of the sequence of transformations applied to each RDD, allowing it to reconstruct lost partitions by reapplying those transformations. This fault tolerance ensures reliability in distributed computing environments where failures are common.

4. Lazy Evaluation:

Spark employs lazy evaluation with RDDs, meaning transformations are not
executed immediately when called. Instead, Spark builds up a Directed Acyclic Graph
(DAG) representing the computation plan. The actual computation occurs only when
an action is triggered. Lazy evaluation optimizes performance by allowing Spark to
optimize and reorder operations before execution.

5. **In-Memory Computation**:

 RDDs can be cached or persisted in memory, allowing for faster access during subsequent computations. This in-memory processing capability is crucial for iterative algorithms and interactive data analysis, as it reduces disk I/O overhead and improves overall performance.

6. Fault Tolerance Mechanism:

• RDDs achieve fault tolerance through lineage information and resilient distributed storage. When a partition of an RDD is lost due to a node failure, Spark can reconstruct it by reapplying the transformations that led to its creation. This mechanism ensures data integrity and reliability in distributed computations.

7. Wide Range of Operations:

RDDs support a wide range of transformations and actions for data manipulation
and analysis. Transformations include map, filter, flatMap, reduceByKey, join, etc.,
while actions include collect, count, reduce, foreach, saveAsTextFile, etc. These
operations enable complex data processing workflows and analytics tasks.

8. Persistence Options:

 RDDs can be persisted in memory, on disk, or both, using various storage levels (e.g., MEMORY_ONLY, MEMORY_AND_DISK, DISK_ONLY, etc.). Persistence allows RDDs to be reused across multiple computations, reducing redundant computation and improving performance.

9. Multi-Language Support:

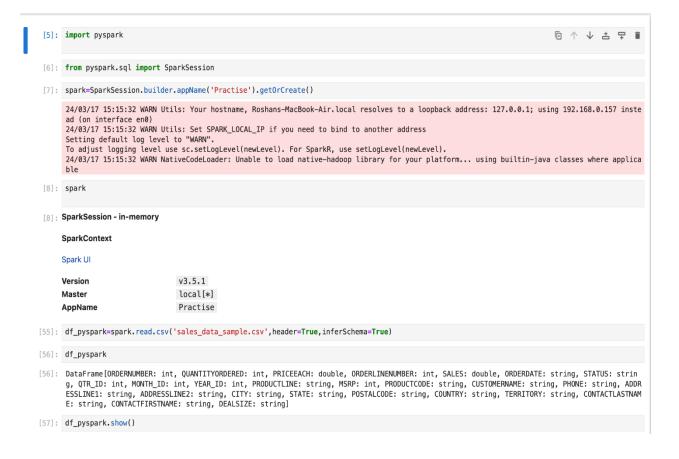
 RDDs are accessible from multiple programming languages, including Scala, Java, and Python, making Spark a versatile framework for developers with different language preferences.

Overall, RDDs serve as the foundational building blocks of distributed data processing in Apache Spark, providing resilience, parallelism, and flexibility for handling diverse big data workloads. While DataFrame and Dataset APIs offer higher-level abstractions and optimizations, RDDs remain integral for low-level control and specialized use cases in Spark applications.

Task 2: In PySpark, create a program that reads a CSV file containing sales data, performs data cleaning by handling missing values and removing duplicates, calculates the total sales amount for each product, and finally, outputs the results to a new CSV file. Ensure to use transformations and actions in your PySpark script.

Steps we will be doing are:

- 1. Read the CSV file containing sales data.
- 2. Perform data cleaning: Handle missing values. Remove duplicate entries.
- 3. Calculate the total sales amount for each product.
- 4. Output the results to a new CSV file.



ORDERNUMBER QUANTITYORI CUSTOMERNAME NAME CONTACTFIRSTNAME I	PHONE	ADDRESSLINE1 AD		CITY	STATE PO	1_ID YEAR_ID PRODUCTLINE MSF DSTALCODE COUNTRY TERRITOF	RY CONTACTL
+	+-		+	+	+		
10107 Land of Toys Inc. Kwai !		2 2871.0 97 Long Airport	2/24/2003 NULL	0:00 Shipped NYC	1 NY	2 2003 Motorcycles 9 10022 USA M	95 S10_10 NA
10121 Reims Collectables riot Paul	34 81.35 26.47.1555 Small	5 2765.9 59 rue de l'Abbaye	5/7/2003 NULL	0:00 Shipped Reims	2 NULL	5 2003 Motorcycles 9 51100 France EME	
10134 Lyon Souveniers -	41 94.74 +33 1 46 62 7555 2	2 3884.34 7 rue du Colonel	7/1/2003 NULL	0:00 Shipped Paris	3 NULL	7 2003 Motorcycles 9 75508 France EME	
10145 Toys4GrownUps.com		6 3746.7 78934 Hillside Dr.	8/25/2003 NULL	0:00 Shipped Pasadena	3 CA	8 2003 Motorcycles 9 90003 USA M	95 S10_1 NA
oung Julie 10159 Corporate Gift Id	Medium 49 100.0 6505551386	14 5205.27 7734 Strong St.		0:00 Shipped San Francisco	4 CA		95 S10_1 NA
rown Julie 10168 Technics Stores Inc.	Medium 36 96.66 6505556809			0:00 Shipped Burlingame	4 CA	10 2003 Motorcycles 9 94217 USA N	95 S10_1 NA
rano Juri 10180 Daedalus Designs	Medium 29 86.13			0:00 Shipped	4 NULL	11 2003 Motorcycles 9 59000 France EME	95 S10_1
ance Martine 10188	Small 48 100.0	1 5512.32	11/18/2003	0:00 Shipped	4	11 2003 Motorcycles 9	95 S10_1
Herkku Gifts ztan Veysel 10201	+47 2267 3215 D Medium 22 98.57	rammen 121, PR 7 2 2168.54		Bergen 0:00 Shipped	NULL	N 5804 Norway EME 12 2003 Motorcycles 9	
Mini Wheels Co. rphy Julie 10211	6505555787 59 Small 41 100.0	557 North Pendal	NULL	San Francisco 0:00 Shipped	CA	NULL USA N	NA
Auto Canal Petit rier Dominique 10223	(1) 47.55.6555	25, rue Lauriston	NULL		NULL	75016 France EME	EA
Australian Collec uson Peter	03 9520 4555 Medium	636 St Kilda Road	Level 3	Melbourne	/ictoria	3004 Australia APA	AC F
10237 Vitachrome Inc. rick Michael	23 100.0 2125551500 Small	7 2333.12 2678 Kingston Rd.		0:00 Shipped NYC	2 NY	4 2004 Motorcycles 9 10022 USA N	95 S10_: NA
10251 Tekni Collectable rown William	28 100.0 2015559350 Medium	2 3188.64 7476 Moss Rd.	5/18/2004 NULL	0:00 Shipped Newark	2 NJ		95 S10_ NA
10263 Gift Depot Inc.	34 100.0 2035552570	2 3676.76 25593 South Bay Ln.		0:00 Shipped Bridgewater	2 CT		95 S10_: NA
10275 La Rochelle Gifts		1 4177.35 7, rue des Cinqu	7/23/2004 NULL	0:00 Shipped Nantes	3 NULL	7 2004 Motorcycles 9 44000 France EME	
rune Janine 10285 Marta's Replicas Co.	36 100.0	6 4099.68 39323 Spinnaker Dr.		0:00 Shipped Cambridge	3 MA	8 2004 Motorcycles 9 51247 USA N	95 S10_
ndez Marta 10299 Toys of Finland, Co.	Medium 23 100.0 90-224 8555		9/30/2004 NULL	0:00 Shipped Helsinki	3 NULL	9 2004 Motorcycles 9 21240 Finland EME	
unen Matti 10309	Small 41 100.0	5 4394.38	10/15/2004	0:00 Shipped	4 NULL	10 2004 Motorcycles 9	95 S10_:
Baane Mini Imports fsen Jonas 10318	Medium 46 94.74		11/2/2004	0:00 Shipped	4	11 2004 Motorcycles 9	95 S10_:
Diecast Classics u Kyung Me 10329	42 100.0	1 4396.14	11/15/2004	0:00 Shipped	PA 4	70267 USA M	NA 95 S10_:
Land of Toys Inc. Kwai Me	2125557818 8	97 Long Airport	NULL	NYC			NA į

[58]: df_pyspark.head(3)

[58]: [Row(ORDERNUMBER=10107, QUANTITYORDERED=30, PRICEEACH=95.7, ORDERLINENUMBER=2, SALES=2871.0, ORDERDATE='2/24/2003 0:00', STATUS='Shipped', QTR_ID=1, MONTH_ID=2, YEAR_ID=2003, PRODUCTLINE='Motorcycles', MSRP=95, PRODUCTCODE='S10_1678', CUSTOMERNAME='Land of Toys Inc.', PHONE='2 125557818', ADDRESSLINE1='897 Long Airport Avenue', ADDRESSLINE2=Mone, CITY='NYC', STATE='NY', POSTALCODE='10022', COUNTRY='USA', TERRITOR Y='NA', CONTACTLASTNAME='Yu', CONTACTFIRSTNAME='Kwai', DEALSIZE='Small'),

Row(ORDERNUMBER=10121, QUANTITYORDERED=34, PRICEEACH=81.35, ORDERLINENUMBER=5, SALES=2765.9, ORDERDATE='5/7/2003 0:00', STATUS='Shipped', OTR_ID=2, MONTH_ID=5, YEAR_ID=2003, PRODUCTLINE='Motorcycles', MSRP=95, PRODUCTCODE='S10_1678', CUSTOMERNAME='Reims Collectables', PHONE ='26.47.1555', ADDRESSLINE1="59 rue de l'Abbaye", ADDRESSLINE2=None, CITY='Reims', STATE=None, POSTALCODE='51100', COUNTRY='France', TERRI TORY='EMEA', CONTACTLASTNAME='Henriot', CONTACTFIRSTNAME='Paul', DEALSIZE='Small'),

ROW(ORDERNUMBER=10134, QUANTITYORDERED=41, PRICEEACH=94.74, ORBERLINENUMBER=2, SALES=3884.34, ORDERDATE='7/1/2003 0:00', STATUS='Shipped', QTR_ID=3, MONTH_ID=7, YEAR_ID=2003, PRODUCTLINE='Motorcycles', MSRP=95, PRODUCTCODE='S10_1678', CUSTOMERNAME='Lyon Souveniers', PHONE ='433 1 46 62 7555', ADDRESSLINE1='27 rue du Colonel Pierre Avia', ADDRESSLINE2=Mone, CITY='Paris', STATE=None, POSTALCODE='75508', COUNTR Y='France', TERRITORY='EMEA', CONTACTLASTNAME='Da Cunha', CONTACTFIRSTNAME='Daniel', DEALSIZE='Medium')]

[59]: df_pyspark=df_pyspark.drop('ORDERDATE')

[60]: df_pyspark

[60]: DataFrame[ORDERNUMBER: int, QUANTITYORDERED: int, PRICEEACH: double, ORDERLINENUMBER: int, SALES: double, STATUS: string, QTR_ID: int, MON TH_ID: int, YEAR_ID: int, PRODUCTLINE: string, MSRP: int, PRODUCTCODE: string, CUSTOMERNAME: string, PHONE: string, ADDRESSLINE1: string, ADDRESSLINE2: string, CITY: string, STATE: string, POSTALCODE: string, COUNTRY: string, TERRITORY: string, CONTACTLASTNAME: string, CONTACT TFIRSTNAME: string, DEALSIZE: string]

[61]: df_pyspark.show()

		+		+	+	+	+	++-	
	+								-
ERNAME	UMBER QUANTITYORDE PHONE DEALSIZE	RED PRICEEACH ORDERL ADDRESSLINE1							PRODUCTCODE CONTACTLASTNAME CONT
+	+	+		+	+	+	+	++-	
+	+							+-	
s Inc.		30 95.7 397 Long Airport		.0 Shipped NYC	1 NY	2 2003 10022	Motorcycles USA	95 NA	S10_1678 Land Yu
 tables	10121 26.47.1555	34 81.35 59 rue de l'Abbaye		.9 Shipped Reims	2 NULL		Motorcycles France		S10_1678 Reims Henriot
 eniers	Small 10134 +33 1 46 62 7555 2	41 94.74 27 rue du Colonel		34 Shipped Paris	3 NULL		Motorcycles France		S10_1678 Lyc Da Cunha
Daniel ps.com	Medium 10145 6265557265	45 83.26 78934 Hillside Dr.	6 3746. NULL	.7 Shipped Pasadena	3 CA	8 2003 90003	Motorcycles USA	95 NA	S10_1678 Toys4 Young
Julie 	Medium 10159 6505551386	49 100.0 7734 Strong St.	14 5205.2	27 Shipped Francisco	4 CAI		Motorcycles USA	95 NAI	S10_1678 Corporat
Julie	Medium 10168	36 96.66 9408 Furth Circle	1 3479.7	76 Shipped Burlingame	4 CA		Motorcycles USA		S10_1678 Technics
Juri 	Medium 10180	29 86.13 184, chausse de T	9 2497.7	77 Shipped Lille	4 NULL	11 2003	Motorcycles France	95	S10_1678 Daedalus Rance
artine	Small 10188	48 100.0	1 5512.3	32 Shipped	4	11 2003	Motorcycles	95	S10_1678
eysel	Medium 10201	22 98.57	2 2168.5	Bergen 54 Shipped	NULL	12 2003	Norway Motorcycles	95	0eztan S10_1678 Mir
Julie	6505555787 5 Small 10211	5557 North Pendal 41 100.0		n Francisco 44 Shipped	(CA)		USA Motorcycles		Murphy S10_1678 Auto
Petit nique		25, rue Lauriston 37 100.0	NULL	Paris 56 Shipped	NULL	75016	France Motorcycles	EMEA	Perrier S10_1678 Australi
lec Peter	03 9520 4555 Medium	636 St Kilda Road	Level 3	Melbourne V		3004 Au	stralia Motorcycles	APAC	Ferguson S10 1678 Vit
e Inc. ichael	2125551500 Small	2678 Kingston Rd.	Suite 101	NYC	NY	10022	USA	NA	Frick
ble illiam	10251 2015559350 Medium	28 100.0 7476 Moss Rd.	NULL	54 Shipped Newark	2 NJ	94019	Motorcycles USA	NA	S10_1678 Tekni Co Brown
 t Inc.	10263	34 100.0 25593 South Bay Ln.		76 Shipped Bridgewater	2 CT	6 2004 97562	Motorcycles USA	95 NA	S10_1678 Gif King
 Gifts	10275	45 92.83 7, rue des Cinqu	1 4177.3 NULL	35 Shipped Nantes	3 NULL		Motorcycles France		S10_1678 La Ro Labrune
as Co.	10285	36 100.0 39323 Spinnaker Dr.		68 Shipped Cambridge	3 MA		Motorcycles USA	95 NA	S10_1678 Marta's Hernandez
 d, Co.	10299 90-224 8555	23 100.0 Keskuskatu 45		39 Shipped Helsinki	3 NULL		Motorcycles Finland		S10_1678 Toys of Karttunen
 mports	Small 10309 07-98 9555 E	41 100.0 Erling Skakkes ga	5 4394.3 NULL	38 Shipped Stavern	4 NULL		Motorcycles Norway	95 EMEA	S10_1678 Baane Bergulfsen
Jonas cs	Medium 10318 2155551555	46 94.74 7586 Pompton St.	1 4358.6	04 Shipped Allentown	4 PA	11 2004 70267	Motorcycles USA	95 NA	S10_1678 Diecast Yu
1		42 100.0 397 Long Airport	1 4396.1 NULL	14 Shipped NYC	4 NY	11 2004 10022	Motorcycles USA	95 NA	S10_1678 Land Yu
	Medium		·	+	·	· +	+	++-	·

ORDERLINENUMBER SALES STATUS QTR_ID MONTH_ID YEAR_ID PRODUCTLINE MSRP PF LINE1 ADDRESSLINE2 CITY STATE POSTALCODE COUNTRY TERRITORY CONT	
1 3965.66 Shipped 1 2 2004 Motorcycles 95 Road Level 3 Melbourne Victoria 3004 Australia APAC	S10_1678 Australian
Road Level 3 Melbourne Victoria 3004 Australia APAC	Ferguson
7 2333.12 Shipped 2 4 2004 Motorcycles 95 n Rd. Suite 101 NYC NY 10022 USA NA	S10_1678 Vitac Frick
13 1451.0 Shipped 4 12 2004 Motorcycles 95 ui Level 6 Chatswood NSW 2067 Australia APAC	S10_1678 Souveniers Huxley
1 4860.24 Shipped 4 10 2003 Classic Cars 214 n St. Suite 750 NYC NY 10022 USA NA	S10_1949 Classic Le Hernandez
9 4905.39 Shipped 3 7 2004 Classic Cars 214 ui Level 6 Chatswood NSW 2067 Australia APAC	S10_1949 Souveniers Huxley
1 3944.7 Shipped 4 11 2004 Classic Cars 214 Road Level 3 Melbourne Victoria 3004 Australia APAC	S10_1949 Australian Ferguson
4 2416.56 Shipped 1 3 2005 Classic Cars 214 treet Level 15 North Sydney NSW 2060 Australia APAC	S10_1949 Anna's Dec O'Hara
212702 00101/	640 204614
3 2793.86 Shipped 2 4 2003 Motorcycles 118 Road Level 3 Melbourne Victoria 3004 Australia APAC	S10_2016 Australia Ferguson
4 5422.39 Shipped 1 2 2004 Motorcycles 118 Road Level 3 Melbourne Victoria 3004 Australia APAC	S10_2016 Australia Ferguson
8 1329.9 Shipped 4 12 2004 Motorcycles 118	
ui Level 6 Chatswood NSW 2067 Australia APAC 2 9264.86 Shipped 2 4 2003 Motorcycles 193	Huxley
Road Level 3 Melbourne Victoria 3004 Australia APAC	Ferguson
3 9774.03 Shipped 1 2 2004 Motorcycles 193 Road Level 3 Melbourne Victoria 3004 Australia APAC	S10_4698 Australia Ferguson
9 7023.9 Shipped 2 4 2004 Motorcycles 193	S10_4698 Vita
9 7023.9 Shipped 2 4 2004 Motorcycles 193 n Rd. Suite 101 NYC NY 10022 USA NA	Frick
8 1201.25 Shipped 4 11 2004 Classic Cars 136 n St. Suite 750 NYC NY 10022 USA NA	S10_4757 Classic Lo Hernandez
2 4302.08 Shipped 3 7 2004 Classic Cars 147	
ui Level 6 Chatswood NSW 2067 Australia APAC 2 4428.0 Shipped 4 11 2004 Classic Cars 147	Huxley S10_4962 Australia
Road Level 3 Melbourne Victoria 3004 Australia APAC	Ferguson
4 2297.05 Shipped 1 1 2005 Classic Cars 147 treet Level 15 North Sydney NSW 2060 Australia APAC	S10_4962 Anna's De O'Hara
7 1735.3 Shipped 1 3 2005 Classic Cars 147 treet Level 15 North Sydney NSW 2060 Australia APAC	
2 5019.9 Shipped 4 11 2003 Classic Cars 194	
treet Level 15 North Sydney NSW 2060 Australia APAC	0'Hara
2 11279.2 Shipped 2 6 2003 Classic Cars 207 ircle Suite 400 NYC NY 10022 USA NA	

+	+-		+		
++	DED DDT CEE 4 CII DDD ED T	UENUMBER L. CALEGE STATUS OF	D TD MONIT	THE TRIVEAR ARE DROPHICT THE INCORDER	DODUCTOODEL
	RED PRICEEACH ORDERLI ADDRESSLINE1 A			TH_ID YEAR_ID PRODUCTLINE MSRP P DSTALCODE COUNTRY TERRITORY CON	
NAME DEALSIZE	ADDRESSEINET A	DDRESSLINEZ CITT	SIAIEJP	DSTALCODE COUNTRY TERRITORY CON	TACTEAS INAME CONTACT
	+		+	+	
+	+-		+		
+					
10223	37 100.0			2 2004 Motorcycles 95	
	636 St Kilda Road	Level 3 Melbourne Vi	ctoria	3004 Australia APAC	Ferguson
eter Medium 10237	23 100.0	7 2333.12 Shipped	21	4 2004 Motorcycles 95	S10_1678 Vita
	2679 Kingston Pd	Suite 101 NYC	2 NY	10022 USA NA	Frick
hael Small	2070 Kingston Ka.	Suite 101 Nic	NI I	10022 03A 14A	TTTCK
10361	20 72.55	13 1451.0 Shipped	41	12 2004 Motorcycles 95	S10_1678 Souvenie
		Level 6 Chatswood	NSW	2067 Australia APAC	Huxley
rian Small					
10163	21 100.0	1 4860.24 Shipped	4	10 2003 Classic Cars 214	S10_1949 Classic
s Inc. 2125558493	5905 Pompton St.	Suite 750 NYC	NY	10022 USA NA	Hernandez
aria Medium	211 100 01	014005 20156	2.1	71 2004161	C10 10401C
10270 Th 1+61 2 0405 85551M	21 100.0	9 4905.39 Shipped Level 6 Chatswood	3 NSW	7 2004 Classic Cars 214 2067 Australia APAC	S10_1949 Souvenie Huxley
rian Medium	onition money builting	Level of Charswood	WCN	2007 MUSCI acta APAC	nux cey [
10347	30 100.0	1 3944.7 Shipped	41	11 2004 Classic Cars 214	S10_1949 Australia
lec 03 9520 4555		Level 3 Melbourne Vi		3004 Australia APAC	Ferguson
eter Medium	·				- '
10391	24 100.0	4 2416.56 Shipped	1	3 2005 Classic Cars 214	S10_1949 Anna's D
ion 02 9936 8555	201 Miller Street	Level 15 North Sydney	NSW	2060 Australia APAC	0'Hara
nna Small	201 06 241	212702 001014	2.1	41 20021 Massacratical 2221	C10 2016
10120	29 96.34	3 2793.86 Shipped		4 2003 Motorcycles 118	S10_2016 Australia
lec 03 9520 4555 eter Small	030 St VIlda KO90	Level 3 Melbourne Vi	ccoria	3004 Australia APAC	Ferguson
10223	47 100.0	4 5422.39 Shipped	11	2 2004 Motorcycles 118	S10_2016 Australia
lec 03 9520 4555		Level 3 Melbourne Vi		3004 Australia APAC	Ferguson
eter Medium					- '
10361	26 51.15	8 1329.9 Shipped	4	12 2004 Motorcycles 118	
Th +61 2 9495 8555 M	onitor Money Bui	Level 6 Chatswood	NSW	2067 Australia APAC	Huxley
rian Small 10120	46 100.0	210264 0616hi	21	41 20021 Motorcycles 1021	C10 4600 Auc+==1:-
10120 lec 03 9520 4555		2 9264.86 Shipped Level 3 Melbourne Vi	2 ctorial	4 2003 Motorcycles 193 3004 Australia APAC	S10_4698 Australia
eter Large	220 21 VIII VOGU	rever 21 Herbourne VI	c coi ta j	2004 Mastrattal MEMC	i ei gustii į
10223	49 100.0	3 9774.03 Shipped	1	2 2004 Motorcycles 193	S10_4698 Australia
lec 03 9520 4555		Level 3 Melbourne Vi		3004 Australia APAC	Ferguson
eter Large					- '
10237	39 100.0	9 7023.9 Shipped	2	4 2004 Motorcycles 193	S10_4698 Vita
	2678 Kingston Rd.	Suite 101 NYC	NY	10022 USA NA	Frick
hael Large	251 40 051	011201 25165	41	111 2004161	C10 47571C1
10337 s Inc. 2125558493	25 48.05 5005 Pompton St	8 1201.25 Shipped Suite 750 NYC	4 NYI	11 2004 Classic Cars 136 10022 USA NA	S10_4757 Classic I Hernandez
s inc. 2125558493 aria Small	Jago rompton St.	Surice /Sul NIC	INT	10022 USA NA	nernanuez
10270	32 100.0	2 4302.08 Shipped	3	7 2004 Classic Cars 147	S10_4962 Souvenie
		Level 6 Chatswood	NSW	2067 Australia APAC	Huxley
rian Medium	,		- 1		
10347	27 100.0	2 4428.0 Shipped	4	11 2004 Classic Cars 147	
lec 03 9520 4555	636 St Kilda Road	Level 3 Melbourne Vi	ctoria	3004 Australia APAC	Ferguson
eter Medium	251 65 551	410007		41 2005163	
10370	35 65.63	4 2297.05 Shipped	1 NSWI	1 2005 Classic Cars 147	S10_4962 Anna's D
ion 02 9936 8555 nna Small	ZWI MILLER STREET	Level 15 North Sydney	M2M	2060 Australia APAC	0'Hara
10391	37 46.9	7 1735.3 Shipped	1	3 2005 Classic Cars 147	S10_4962 Anna's D
ion 02 9936 8555		Level 15 North Sydney	NSW	2060 Australia APAC	0'Hara
nna Small					
10169	30 100.0	2 5019.9 Shipped	4	11 2003 Classic Cars 194	
ion 02 9936 8555	201 Miller Street	Level 15 North Sydney	NSW	2060 Australia APAC	0'Hara
nna Medium					
10127	46 100.0	2 11279.2 Shipped	2	6 2003 Classic Cars 207	
ne Inc 2125557413 eff Large	4092 Furth Circle	Suite 400 NYC	NY	10022 USA NA	Young

```
imputer = Imputer(
                             inputCols=["{}_imputed",'PRICEEACH', 'SALES'],
outputCols=["{}_imputed".format(c) for c in ['QUANTITYORDERED','PRICEEACH', 'SALES']]
                             ).setStrategy("median")
[67]: imputer.fit(df_pyspark).transform(df_pyspark).show()
                  |ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER| SALES| STATUS|QTR_ID|MONTH_ID|YEAR_ID| PRODUCTLINE|MSRP|PRODUCTCODE| CUSTO MERNAME| PHONE| ADDRESSLINE1|ADDRESSLINE2| CITY| STATE|POSTALCODE| COUNTRY|TERRITORY|CONTACTLASTNAME|CONTACTFIRS TNAME|DEALSIZE|QUANTITYORDERED_imputed|PRICEEACH_imputed|SALES_imputed|
                 | 10223| 37| 100.0|
|llec...| 03 9520 4555| 636 St Kilda Road|
|Peter| Medium| 37|
| 10237| 23| 100.0|
|me Inc.| 2125551500| 2678 Kingston Rd.|
                                                                                                                                 ----+----
                                                                                                                                                                                                                                                                   2| 2004| Motorcycles| 95| S10_1678|Australian Co
3004|Australia| APAC| Ferguson|
                                                                                                                                                               Level 3| Melbourne|Victoria|
                                                                                                                                                                 100.01
                                                                                                                                                                                                   3965.661
                                                                                                                                                                                                                                                                  4| 2004| Motorcycles| 95|
10022| USA| NA|
                                                                                                                                                         7|2333.12|Shipped|
Suite 101| NYC|
100.0| 2333.12|
                                                                                                                                                                                                                                                                                                                                                        S10_1678|
                    | 1045.
me Inc.| 2125551500| - 23|
chael| Small| 23| 72.55|
20| 72.55|
                 me Inc.|
chael| Small|
| 10361| 20| 72.55|
d Th...|+61 2 9495 8555|Monitor Money Bui...|
drian| Small| 20|
| 10163| 21| 100.0|
ds Inc.| 2125558493| 5905 Pompton St.|
211
100.0|
                                                                                                                                                              13| 1451.0|Shipped|
Level 6| Chatswood|
                                                                                                                                                                                                                                                               12| 2004| Motorcycles| 95|
                                                                                                                                                                                                                                                                                                                                                        S10_1678|Souveniers An
                                                                                                                                                                                                                                      NSWI
                                                                                                                                                                                                                                                                    2067|Australia| APAC|
                                                                                                                                                         Level 6| 72.55| 1451.0| 1|4860.24|Shipped| Suite 750| NYC| 100.0| 4860.24|
                                                                                                                                                                                                                                                                                                                                                                    Huxley|
                                                                                                                                                                                                                                                                            2003|Classic Cars| 214|
                                                                                                                                                                                                                                                                                                                                                        S10_1949|Classic Legen
                 10022| USA|
                                                                                                                                                                                                                                                                                                                                                            Hernandez|
                  100.0| 4860.24| 9|4905.39|Shipped| 3|
Level 6| Chatswood| NSW| 100.0| 4905.39| 1| 3944.7|Shipped| 4|
Level 3| Melbourne|Victoria| 100.01 3944.1
                                                                                                                                                                                                                                                                           2004|Classic Cars| 214|
                                                                                                                                                                                                                                                                                                                                                        S10_1949|Souveniers An
                                                                                                                                                                                                                                                                     2067|Australia| APAC|
                                                                                                                                                                                                                                                                                                                                                                    Huxley|
                                                                                                                                                                                                                                                               11| 2004|Classic Cars| 214|
3004|Australia| APAC|
                                                                                                                                                                                                                                                                                                                                                        S10_1949|Australian Co
                 | Net 
                                                                                                                                                                                                                                                                                                                                                               Ferguson
                                                                                                                                                                 100.01
                                                                                                                                                                                                      3944.7
                                                                                                                                                            4|2416.56|Shipped|
Level 15|North Sydney|
100.0| 2416.56|
                                                                                                                                                                                                                                                                           2005|Classic Cars| 214|
                                                                                                                                                                                                                                                                                                                                                        S10_1949|Anna's Decora
                                                                                                                                                               100.0| 2416.56|
3|2793.86|Shipped|
                                                                                                                                                                                                                                                                            2003| Motorcycles| 118|
                                                                                                                                                                                                                                                                                                                                                        S10 2016|Australian Co
                 | 101/v| | 102 | 103 9520 4555| 636 St Kilda Roau| | 29| | 10223| | 47| 100.0| | 102.0| | 03 9520 4555| 636 St Kilda Road| | 47| | 10361| 26| 51.15| | 10361| 26| 51.15|
                                                                                                                                                              Level 3| Melbourne|Victoria|
96.34| 2793.86|
4|5422.39|Shipped| 1|
                                                                                                                                                                                                                                                                    3004|Australia| APAC|
                                                                                                                                                                                                                                                                  2| 2004| Motorcycles| 118|
3004|Australia| APAC|
                                                                                                                                                                                                                                                                                                                                                        S10_2016|Australian Co
                                                                                                                                                              Level 3| Melbourne|Victoria|
                                                                                                                                                                                                                                                                                                                                                               Ferguson
               | Telecond 
                                                                                                                                                                 100.01
                                                                                                                                                                                                   5422.391
                                                                                                                                                                                                                                                                     !| 2004| Motorcycles| 118|
2067|Australia| APAC|
                                                                                                                                                                8| 1329.9|Shipped|
Level 6| Chatswood|
                                                                                                                                                                                                                                                                                                                                                        S10_2016|Souveniers An
                                                                                                                                                                                                                                                                                                                                                                    Huxley|
                                                                                                                                                                 51.15| 1329.9|
2|9264.86|Shipped|
                                                                                                                                                                                                                                                                            2003| Motorcycles| 193|
                                                                                                                                                                                                                                                                                                                                                        S10_4698|Australian Co
                                                                                                                                                              Level 3| Melbourne|Victoria|
100.0| 9264.86|
3|9774.03|Shipped| 1|
                                                                                                                                                                                                                                                                    3004|Australia| APAC|
                                                                                                                                                                                                                                                                                                                                                              Ferguson|
                                                                                                                                                                                                                                                                 2| 2004| Motorcycles| 193
3004|Australia| APAC|
                                                                                                                                                                                                                                                                           2004| Motorcycles| 193|
                                                                                                                                                                                                                                                                                                                                                         S10_4698|Australian Co
                                                                                                                                                              Level 3| Melbourne|Victoria|
100.0| 9774.03|
                                                                                                                                                                                                                                                                                                                                                              Ferguson|
                                                                                                                                                                                                                                                                                                                                                                     4698| Vitachro
Frick|
                                                                                                                                                         100.0| 9//4.05|
9| 7023.9|Shipped|
Suite 101| NYC|
100.0| 7023.9|
                                                                                                                                                                                                                                                                  4| 2004| Motorcycles| 193|
10022| USA| NA|
                                                                                                                                                                                                                                                                                                                                                        S10_4698|
                 ΝΥΙ
                                                                                                                                                                  100.0| 7023.9|
8|1201.25|Shipped|
                                                                                                                                                                                                                                    4 |
NY |
                                                                                                                                                                                                                                                                           2004|Classic Cars| 136|
                                                                                                                                                                                                                                                                                                                                                         S10_4757|Classic Legen
                 ds Inc.| 2125558495,
Maria| Small| 25,
Maria| 32| 100.0|
                                                                                                                                                         Suite 750| Nici
                                                                                                                                                                                                                                                                                                                                 NA |
                                                                                                                                                                                                                                                                  10022| USA|
                                                                                                                                                                                                                                                                                                                                                            Hernandez|
                 Maria| Small|
| 10270|
| Th...|+61 2 9495 8555|Monitor Money Bui...|
                                                                                                                                                                                                                                                                     7| 2004|Classic Cars| 147|
2067|Australia| APAC|
                                                                                                                                                               2|4302.08|Shipped|
Level 6| Chatswood|
100.0| 4302.08|
                                                                                                                                                                                                                                                                                                                                                        S10_4962|Souveniers An
                                                                                                                                                                                                                                                                                                                                                                    Huxley|
                                       | 4302.08|
| 2| 4428.0|Shipped| 4|
| Level 3| Melbourne|Victoria|
                                                                                                                                                                                                                                                                           2004|Classic Cars| 147|
                 | 10347| 27| 100.0| |
|lec...| 03 9520 4555| 636 St Kilda Road|
|Peter| Medium| 27|
| 10370| 35| 65.63|
|tion...| 02 9936 8555| 201 Miller Street|
| Anna| Small| 37| 46.9|
| 10169| 30| 100.0|
| 10169| 30| 100.0|
| 101731| 46| 100.0|
| 101731| 46| 100.0|
                                     103471
                                                                                                                                                                                                                                                                                                                                                        S10 4962|Australian Co
                                                                                                                                                                                                                                                                     3004|Australia| APAC|
                                                                                                                                                                                                                                                                                                                                                               Ferguson|
                                                                                                                                                            100.0| 4428.0|
4|2297.05|Shipped|
Level 15|North Sydney|
                                                                                                                                                                                                                                                                             2005|Classic Cars| 147|
                                                                                                                                                                                                                                                                                                                                                         S10_4962|Anna's Decora
                                                                                                                                                                                                                                                                     2060|Australia| APAC|
                                                                                                                                                                                                                                       NSW|
                                                                                                                                                                                                                                                                                                                                                                    0'Hara|
                                                                                                                                                           65.63| 2297.05|
7| 1735.3|Shipped|
Level 15|North Sydney|
46.9| 1735.3|
2| 5019.9|Shipped|
                                                                                                                                                                                                                                                                           2005|Classic Cars| 147|
                                                                                                                                                                                                                                                                                                                                                        S10_4962|Anna's Decora
                                                                                                                                                                                                                                                                     2060|Australia| APAC|
                                                                                                                                                                                                                                       NSW|
                                                                                                                                                                                                                                                                            2003|Classic Cars| 194|
                                                                                                                                                                                                                                                                                                                                                        S12_1099|Anna's Decora
                                                                                                                                                            Level 15|North Sydney|
                                                                                                                                                                                                                                        NSWI
                                                                                                                                                                                                                                                                     2060|Australia| APAC|
                                                                                                                                                                                                                                                                                                                                                                    0'Haral
                                                                                                                                                                100.0| 5
                                                                                                                                                                                                   5019.91
```

[66]: from pyspark.ml.feature import Imputer

						II+	T V T T
			-+	+	+	+	
+ DRDERNUMBER QUANTITYORDERED F	RICEEACHIORDERLINENU	IBERI SALESI STATUSIOT	R IDIMON	NTH IDIY	'EAR ID! PRODUCTLINE!	MSRP I PR	ODUCTCODE I
STOMERNAME PHONE ACTFIRSTNAME DEALSIZE	ADDRESSLINE1 A				ALCODE COUNTRY TER		
	+	++	+	+-	+	+	
·+++	t-		-+	+	+	+	+-
10223 37 Collec 03 9520 4555	100.0 636 St Kilda Road	1 3965.66 Shipped Level 3 Melbourr	1 e Victo	2 ria	2004 Motorcycles 3004 Australia	95 APAC	S10_1678 Austral Ferguson
eter Medium 10237 23 nrome Inc. 2125551500	100.0 2678 Kingston Rd.	7 2333.12 Shipped Suite 101 NY	2	4 NY	2004 Motorcycles 10022 USA	95 NA	S10_1678 Vi
ichael Śmall 10361 20	72.55	13 1451.0 Shipped	4	12	2004 Motorcycles	95	S10_1678 Souveni
And Th +61 2 9495 8555 Mc Irian Small 10163 21	nitor Money Bui	Level 6 Chatswood 1 4860.24 Shipped	d N 4	NSW 10	2067 Australia 2003 Classic Cars	APAC	Huxley S10_1949 Classic
ends Inc. 2125558493 ria Medium	5905 Pompton St.	Suite 750 NY		NY	10022 USA	NA	Hernandez
10270 21 And Th +61 2 9495 8555 Mc drian Medium	100.0 onitor Money Bui	9 4905.39 Shipped Level 6 Chatswoo	3 d N	7 NSW	2004 Classic Cars 2067 Australia	214 APAC	S10_1949 Souveni Huxley
10347 30 Collec 03 9520 4555	100.0 636 St Kilda Road	1 3944.7 Shipped Level 3 Melbourn	4 e Victo	11 ria	2004 Classic Cars 3004 Australia	214 APAC	S10_1949 Austral Ferguson
eter Medium 10391 24 pration 02 9936 8555	100.0 201 Miller Street	4 2416.56 Shipped Level 15 North Sydne	1 y N	3 NSW	2005 Classic Cars 2060 Australia	214 APAC	S10_1949 Anna's O'Hara
na Small 10120 29 Collec 03 9520 4555	96.34 636 St Kilda Road	3 2793.86 Shipped Level 3 Melbourn	2 elVictor	4 rial	2003 Motorcycles 3004 Australia	118 APAC	S10_2016 Austra
ter Small 10223 47	100.0	4 5422.39 Shipped	1	2	2004 Motorcycles	118	S10_2016 Austra
Collec 03 9520 4555 eter Medium 10361 26	636 St Kilda Road 51.15	Level 3 Melbourr 8 1329.9 Shipped	4	12	3004 Australia 2004 Motorcycles	APAC	Ferguson S10_2016 Souveni
And Th +61 2 9495 8555 Mo		Level 6 Chatswoo		NSW	2067 Australia	APAC	Huxley
10120 46 Collec 03 9520 4555 ter Large	100.0 636 St Kilda Road	2 9264.86 Shipped Level 3 Melbourn	2 e Victo	4 ria	2003 Motorcycles 3004 Australia	APAC	S10_4698 Austra Ferguson
10223 49 Collec 03 9520 4555 ter Large	100.0 636 St Kilda Road	3 9774.03 Shipped Level 3 Melbourr	1 e Victo	2 ria	2004 Motorcycles 3004 Australia	193 APAC	S10_4698 Austra Ferguson
10237 39 rome Inc. 2125551500	100.0 2678 Kingston Rd.	9 7023.9 Shipped Suite 101 NY	2 C	4 NY	2004 Motorcycles 10022 USA	193 NA	S10_4698 Vi Frick
ichael Large 10337 25 gends Inc. 2125558493	48.05 5905 Pompton St.	8 1201.25 Shipped Suite 750 NY	4 C	11 NY	2004 Classic Cars 10022 USA	136 NA	S10_4757 Classic
ria Small 10270 32 And Th +61 2 9495 8555 Mo	100.0	2 4302.08 Shipped Level 6 Chatswoo	3	7 NSW	2004 Classic Cars 2067 Australia	147 APAC	S10_4962 Souveni
rian Medium							
10347 27 Collec 03 9520 4555 ter Medium	100.0 636 St Kilda Road	2 4428.0 Shipped Level 3 Melbourr	4 e Victo	11 ria	2004 Classic Cars 3004 Australia	147 APAC	S10_4962 Austra Ferguson
10370 35 ration 02 9936 8555	65.63 201 Miller Street	4 2297.05 Shipped Level 15 North Sydne	1 y N	1 NSW	2005 Classic Cars 2060 Australia	147 APAC	S10_4962 Anna's O'Hara
na Small 10391 37 ration 02 9936 8555	46.9 201 Miller Street	7 1735.3 Shipped Level 15 North Sydne	1 y N	3 NSW	2005 Classic Cars 2060 Australia	147 APAC	S10_4962 Anna's 0'Hara
na Small 10169 30 ration 02 9936 8555	100.0 201 Miller Street	2 5019.9 Shipped Level 15 North Sydne	4	11 NSW	2003 Classic Cars 2060 Australia	194 APAC	S12_1099 Anna's 0'Hara
na Medium 10127 46	100.0	2 11279.2 Shipped	2	45W 6	2003 Classic Cars		
chine Inc 2125557413 eff Large	4092 Furth Circle	Suite 400 NY	Cl	NY	10022 USA	NA	Young

only showing top 20 rows

| ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENMBER| SALES| STATUS|QTR_ID||MONTH_ID|YEAR_ID| PRODUCTLINE|MSRP|PRODUCTCODE| C USTOMERNAME| PHONE| ADDRESSLINE1|ADDRESSLINE2| CITY| STATE|POSTALCODE| COUNTRY|TERRITORY|CONTACTLASTNAME|CONTACTFIRSTNAME|DEALSIZE|QUANTITYORDERED_imputed|PRICEEACH_imputed|SALES_imputed|

· · · · · · · · · · · · · · · · · · ·	+	
10223 37	100.0	1 3965.66 Shipped 1 2 2004 Motorcycles 95 S10_1678 Australia
n Collec 03 9520 4555	636 St Kilda Road	Level 3 Melbourne Victoria 3004 Australia APAC Ferguson
Peter Medium	37	100.0 3965.66
10237 23		7 2333.12 Shipped 2 4 2004 Motorcycles 95 S10_1678 Vita
chrome Inc. 2125551500	2678 Kingston Rd.	
Michael Small	23	100.0 2333.12
10361 20	72.55	13 1451.0 Shipped 4 12 2004 Motorcycles 95 S10_1678 Souvenier
s And Th +61 2 9495 8555 M		
Adrian Small	20	72.55 1451.0
10163 21	100.0	1 4860.24 Shipped 4 10 2003 Classic Cars 214 S10_1949 Classic L
egends Inc. 2125558493	5905 Pompton St.	
Maria Medium	21	100.0 4860.24
10270 21	100.0	9 4905.39 Shipped 3 7 2004 Classic Cars 214 S10_1949 Souvenier
s And Th +61 2 9495 8555 M		
Adrian Medium	21	100.0 4905.39
10347 30	100.0	1 3944.7 Shipped 4 11 2004 Classic Cars 214 S10_1949 Australia
n Collec 03 9520 4555	636 St Kilda Road	
Peter Medium	30	100.0 3944.7
10391 24	100.0	4 2416.56 Shipped 1 3 2005 Classic Cars 214 S10_1949 Anna's De
coration 02 9936 8555	201 Miller Street	
Anna Small	24	100.0 2416.56
10120 29	96.34	3 2793.86 Shipped 2 4 2003 Motorcycles 118 S10_2016 Australia
n Collec 03 9520 4555	636 St Kilda Road	
Peter Small	29	96.34 2793.86
10223 47	100.0	4 5422.39 Shipped 1 2 2004 Motorcycles 118 S10_2016 Australia
n Collec 03 9520 4555	636 St Kilda Road	
Peter Medium	47	100.0 5422.39
10361 26 s And Th +61 2 9495 8555 M	51.15	8 1329.9 Shipped 4 12 2004 Motorcycles 118 S10_2016 Souvenier Level 6 Chatswood NSW 2067 Australia APAC Huxley
	, ,	1 1 1 1 1 21
Adrian Small 10120 46	26 100.0	51.15 1329.9 2 9264.86 Shipped 2 4 2003 Motorcycles 193 S10_4698 Australia
n Collec 03 9520 4555	636 St Kilda Road 46	Level 3 Melbourne Victoria 3004 Australia APAC Ferguson 100.0 9264.86
Peter Large 10223 49	100.0	
n Collec 03 9520 4555	636 St Kilda Road	
Peter Large	49	100.0 9774.03
10237 39	100.0	9 7023.9 Shipped 2 4 2004 Motorcycles 193 S10_4698 Vita
chrome Inc. 2125551500	2678 Kingston Rd.	
Michael Large	39	100.0 7023.9
10337 25	48.05	8 1201.25 Shipped 4 11 2004 Classic Cars 136 S10_4757 Classic L
egends Inc. 2125558493	5905 Pompton St.	
Maria Small	25	48.05 1201.25
10270 32	100.0	2 4302.08 Shipped 3 7 2004 Classic Cars 147 S10_4962 Souvenier
s And Th +61 2 9495 8555 M		
Adrian Medium	321	100.0 4302.08
10347 27	100.0	2 4428.0 Shipped 4 11 2004 Classic Cars 147 S10 4962 Australia
n Collec 03 9520 4555	636 St Kilda Road	
Peter Medium	27	100.0 4428.0
10370 35	65.63	4 2297.05 Shipped 1 1 2005 Classic Cars 147 S10 4962 Anna's De
coration 02 9936 8555	201 Miller Street	
Anna Small	35	65.63 2297.05
10391 37	46.9	7 1735.3 Shipped 1 3 2005 Classic Cars 147 S10_4962 Anna's De
coration 02 9936 8555	201 Miller Street	
Annal Small	371	46.9 1735.3
10169 30	100.0	2 5019.9 Shipped 4 11 2003 Classic Cars 194 S12_1099 Anna's De
coration 02 9936 8555	201 Miller Street	
Anna Medium	30	100.0 5019.9
10127 46	100.0	2 11279.2 Shipped 2 6 2003 Classic Cars 207 S12_1108 Muscle
Machine Inc 2125557413	4092 Furth Circle	
Jeff Large	46	100.0 11279.2
+		

only showing ton 20 rows

[70]: imputed_df.head()

[70]: Row(ORDERNUMBER=10223, QUANTITYORDERED=37, PRICEEACH=100.0, ORDERLINENUMBER=1, SALES=3965.66, STATUS='Shipped', QTR_ID=1, MONTH_ID=2, YEAR_ID=2004, PRODUCTLINE='Motorcycles', MSRP=95, PRODUCTCODE='S10_1678', CUSTOMERNAME='Australian Collectors, Co.', PHONE='03 9520 45 55', ADDRESSLINE1='636 5t Kilda Road', ADDRESSLINE2='Level 3', CITY='Melbourne', STATE='Victoria', POSTALCODE='3004', COUNTRY='Australia', TERRITORY='APAC', CONTACTLASTNAME='Ferguson', CONTACTFIRSTNAME='Peter', DEALSIZE='Medium', QUANTITYORDERED_imputed=37, PRICEEACH_imputed=100.0, SALES_imputed=3965.66)

[72]: imputed_df=imputed_df.drop('QUANTITYORDERED','PRICEEACH', 'SALES')

[73]: # Drop specified columns
 columns_to_drop = ['QUANTITYORDERED','PRICEEACH', 'SALES']
 imputed_df = imputed_df.drop(*columns_to_drop)

Show the DataFrame after dropping columns
 imputed_df.show()

DRESSLINE1 A	DRESSLINE2	R STATUS QTR_: CITY d SALES_imputed	STATE PO					DDUCTCODE CU:	STOMERNAME FIRSTNAME DEA	PHONE LSIZE QUANTIT	
		++-			+		+-		+-		
10223 : Kilda Road		l Shipped Melbourne Vi	+ 1 ictoria			Motorcycles ustralia	95 APAC	S10_1678 Australian Ferguson	Collec Peter	03 9520 4555 Medium	63
7 10237 Kingston Rd.		3965.66 7 Shipped NYC	2 NY		2004 0022	Motorcycles USA	95 NA	S10_1678 Vitach	nrome Inc. Michaell	2125551500 Small	26
10361	100.0	2333.12 3 Shipped	4	12	2004	Motorcycles	95	S10_1678 Souveniers		,	Monit
Money Bui 0 10163	72.55	Chatswood 1451.0 1 Shipped	NSW			Australia Classic Cars	APAC	Huxley S10 1949 Classic Leo	Adrian	Small 2125558493	
5 Pompton St. L	Suite 750 100.0	NYC 4860.24	. NY	1	0022	USA	ŅΑĮ	Hernandez	Maria	Medium	
10270 Money Bui	Level 6	9 Shipped Chatswood 4905.39	3 NSW			Classic Cars Australia	214 APAC	S10_1949 Souveniers Huxley		. 2 9495 8555 N Medium	Monit
10347 Kilda Road	Level 3	1 Shipped Melbourne V:				Classic Cars ustralia	214 APAC	S10_1949 Australian Ferguson		03 9520 4555 Medium	63
0 10391 iller Street		3944.7 4 Shipped North Sydney	1 NSW			Classic Cars ustralia	214 APAC	S10_1949 Anna's Deco O'Hara	oration Anna	02 9936 8555 Small	2
1 10120 Kilda Road	100.0	2416.56 B Shipped Melbourne V	2	4	2003	Motorcycles ustralia		S10_2016 Australian Ferguson		03 9520 4555 Small	6
10223	96.34	2793.86 4 Shipped	1	2	2004	Motorcycles	118	S10_2016 Australian	Collec	03 9520 4555	6
Kilda Road 	100.0	Melbourne V: 5422.39 B Shipped				ustralia Motorcycles	APAC	Ferguson S10_2016 Souveniers	Peter		Moni
Money Bui	Level 6 51.15	Chatswood 1329.9	NSW		2067	Australia	APAC	Huxley	Adrian	Small	
10120 Kilda Road	Level 3	2 Shipped Melbourne V: 9264.86	2 ictoria			Motorcycles ustralia	193 APAC	S10_4698 Australian Ferguson	Peter	03 9520 4555 Large	6
10223 Kilda Road	Level 3	B Shipped Melbourne V: 9774.03	1 ictoria			Motorcycles ustralia	193 APAC	S10_4698 Australian Ferguson	Collec Peter	03 9520 4555 Large	6
10237 Kingston Rd.	Suite 101	Shipped NYC	2 NY		2004 0022	Motorcycles USA	193 NA	S10_4698 Vitach Frick	nrome Inc. Michael	2125551500 Large	2
10337 Pompton St.		7023.9 3 Shipped NYC	4 NY		2004 0022	Classic Cars USA	136 NA	S10_4757 Classic Leg	gends Inc. Maria	2125558493 Small	
10270		1201.25 2 Shipped Chatswood	3 NSW			Classic Cars Australia	147 APAC	S10_4962 Souveniers Huxley	And Th +61		Moni
Money Bui 10347	100.0	4302.08 2 Shipped	4	11	2004	Classic Cars	147	S10_4962 Australian	Collec	03 9520 4555	6
Kilda Road 10370	100.0	Melbourne V: 4428.0 4 Shipped	ictoria 1			ustralia Classic Cars	APAC	Ferguson S10_4962 Anna's Deco	Peter oration	Medium 02 9936 8555	2
ller Street	Level 15 65.63	North Sydney 2297.05	NSW	. 2	060 A	ustralia	APAC	0'Hara	Anna	Small	
10391 ller Street 		7 Shipped North Sydney 1735.3	1 NSW			Classic Cars ustralia	147 APAC	S10_4962 Anna's Deco O'Hara	oration Anna	02 9936 8555 Small	2
10169 ller Street	Level 15	2 Shipped North Sydney 5019.9	4 NSW			Classic Cars ustralia	194 APAC	S12_1099 Anna's Deco O'Hara	oration Anna	02 9936 8555 Medium	2
' 10127 Furth Circle	Suite 400	2 Shipped NYC	2 NY			Classic Cars USA	207 NA	S12_1108 Muscle Ma Young	achine Inc Jeff	2125557413 Large	4
il .	100.0	11279.2									