[0]:	<pre>import pandas as pd import numpy as np import seaborn as sns #df=spark.read.format("csv").option("header", "true").load("/FileStore/tables/Sales/Car_Sales.csv")</pre>
	#dr=spark.read.tormat("csv").option("neader", "true").load("/Filestore/tables/Sales/Car_Sales.csv") df=spark.read.csv("/FileStore/tables/Sales/Car_Sales.csv", header=True, inferSchema=True) df.show(5) +++ ORDERNUMBER QUANTITYORDERED PRICEEACH ORDERLINENUMBER SALES ORDERDATE STATUS QTR_ID MONTH_ID YEAR_ID PRODUCTLINE MSRP PRODUCTCODE
	CUSTOMERNAME PHONE ADDRESSLINE1 ADDRESSLINE2 CITY STATE POSTALCODE COUNTRY TERRITORY CONTACTLASTNAME CONTACTFIRSTNAME DEALSIZE +++
	ai Small
	10145 45 83.26 6 3746.7 8/25/2003 0:00 Shipped 3 8 2003 Motorcycles 95 S10_1678 Toys 4GrownUps.com 6265557265 78934 Hillside Dr. null Pasadena CA 90003 USA NA Young Julie Medium 10159 49 100.0 14 5205.27 10/10/2003 0:00 Shipped 4 10 2003 Motorcycles 95 S10_1678 Corporate Gift Id 6505551386 7734 Strong St. null San Francisco CA null USA NA Brown Julie Medium +
	df.printSchema()
	root ORDERNUMBER: integer (nullable = true) QUANTITYORDERED: integer (nullable = true) PRICEEACH: double (nullable = true) ORDERLINENUMBER: integer (nullable = true) SALES: double (nullable = true) ORDERDATE: string (nullable = true)
	STATUS: string (nullable = true) QTR_ID: integer (nullable = true) MONTH_ID: integer (nullable = true) YEAR_ID: integer (nullable = true) PRODUCTLINE: string (nullable = true) MSRP: integer (nullable = true) PRODUCTCODE: string (nullable = true)
	CUSTOMERNAME: string (nullable = true) PHONE: string (nullable = true) ADDRESSLINE1: string (nullable = true) ADDRESSLINE2: string (nullable = true) CITY: string (nullable = true) STATE: string (nullable = true)
	POSTALCODE: string (nullable = true) COUNTRY: string (nullable = true) TERRITORY: string (nullable = true) CONTACTLASTNAME: string (nullable = true) CONTACTFIRSTNAME: string (nullable = true) DEALSIZE: string (nullable = true)
[0]:	<pre>df.select(["SALES","CITY"]).show(5) ++ SALES CITY ++ 2871.0 NYC </pre>
	2765.9 Reims 3884.34 Paris 3746.7 Pasadena 5205.27 San Francisco ++ only showing top 5 rows
[0]:	type(df) Out[88]: pyspark.sql.dataframe.DataFrame df.columns
	Out[89]: ['ORDERNUMBER', 'QUANTITYORDERED', 'PRICEEACH', 'ORDERLINENUMBER', 'SALES', 'ORDERDATE', 'STATUS',
	'QTR_ID', 'MONTH_ID', 'YEAR_ID', 'PRODUCTLINE', 'MSRP', 'PRODUCTCODE', 'CUSTOMERNAME',
	'PHONE', 'ADDRESSLINE1', 'ADDRESSLINE2', 'CITY', 'STATE', 'POSTALCODE', 'COUNTRY',
[0]:	'TERRITORY', 'CONTACTLASTNAME', 'CONTACTFIRSTNAME', 'DEALSIZE'] df_sel=df.select(['PRICEEACH','MSRP','SALES'])
[0]:	df_sel.describe().show() ++
	mean 83.65854410201929 100.71555083244775 3553.88907190932 stddev 20.174276527840536 40.18791167720266 1841.8651057401842 min 26.88 33 482.13 max 100.0 214 14082.8 ++
[0]: [0]:	<pre>df=df.na.drop(how='any', subset=['SALES']) from pyspark.sql.functions import isnan, when, count, col df.select([count(when(isnan(c) col(c).isNull(), c)).alias(c) for c in df.columns]).show()</pre>
	++++++
[0]:	0 0 2521 0 1486 76 0 0 0 0 0 0 0 0 0 0 0 0 0
	<pre>zip(["q1", "q3"], df_sel.approxQuantile(c, [0.25, 0.75], 0))) for c in df_sel.columns } for c in bounds: iqr = bounds[c]['q3'] - bounds[c]['q1'] bounds[c]['lower'] = bounds[c]['q1'] - (iqr * 1.5) bounds[c]['upper'] = bounds[c]['q3'] + (iqr * 1.5) print(bounds)</pre>
[0]:	{'PRICEEACH': {'q1': 68.8, 'q3': 100.0, 'lower': 21.99999999999993, 'upper': 146.8}, 'MSRP': {'q1': 68.0, 'q3': 124.0, 'lower': -16.0, 'upper': 208.0}, 'SALES': {'q1': 2203.11, 'q3': 4508.0, 'lower': -1254.225, 'upper': 7965.335}} df.describe().toPandas()
	summary ORDERNUMBER QUANTITYORDERED PRICEEACH ORDERLINENUMBER SALES ORDERDATE STATUS QTR_ID MONTH 0 count 2823
[0]:	3 min 10100 6 26.88 1 482.13 1/10/2003 Cancelled 1 4 max 10425 97 100.0 18 14082.8 9/9/2004 0:00 Shipped 4 df.toPandas().isnull().sum() Out[96]: ORDERNUMBER 0
	QUANTITYORDERED 0 PRICEEACH 0 ORDERLINENUMBER 0 SALES 0 ORDERDATE 0 STATUS 0
	QTR_ID 0 MONTH_ID 0 YEAR_ID 0 PRODUCTLINE 0 MSRP 0 PRODUCTCODE 0 CUSTOMERNAME 0 PHONE 0
	ADDRESSLINE1 0 ADDRESSLINE2 2521 CITY 0 STATE 1486 POSTALCODE 76 COUNTRY 0 TERRITORY 0
	CONTACTLASTNAME 0 CONTACTFIRSTNAME 0 DEALSIZE 0 dtype: int64 df.toPandas().shape
[0]:	Out[97]: (2823, 25) sns.heatmap(df.toPandas().isnull(),yticklabels=False,cbar=False,cmap='viridis') Out[98]:
	ORDERNUMBER - PRICEEACH - ORDERLINES - ORDERDATE - SALES - ORDERDATE - STATUS - TATUS - TATUS - TO TR ID - YEAR ID - YEAR ID - PRODUCTCODE - CUSTOMERNAME - PHONE - ADDRESSLINE - ADDRESSLINE - COUNTRY - TERRITORY - TERRITOR
[0]:	### ### ### ### ### ##################
[0]:	<pre>df=df.na.drop(how='any', subset=['POSTALCODE']) df.toPandas().isnull().sum()</pre>
	Out[101]: ORDERNUMBER 0 QUANTITYORDERED 0 PRICEEACH 0 ORDERLINENUMBER 0 SALES 0 ORDERDATE 0 STATUS 0
	QTR_ID 0 MONTH_ID 0 YEAR_ID 0 PRODUCTLINE 0 MSRP 0 PRODUCTCODE 0 CUSTOMERNAME 0
	PHONE 0 ADDRESSLINE1 0 CITY 0 POSTALCODE 0 COUNTRY 0 TERRITORY 0 CONTACTLASTNAME 0
[0]:	CONTACTFIRSTNAME 0 DEALSIZE 0 dtype: int64 sns.heatmap(df.toPandas().isnull(),yticklabels=False,cbar=False,cmap='viridis') Out[102]:
	ORDERNUMBER - DERLINENUMBER - SALES - ORDERDATE - STATUS - GTR_ID - MONTH_ID - YEAR_ID - PRODUCTLINE - MOSP - PRODUCTCODE - COLVITRY - TERRITORY - TERRITORY - TERRITORY - TERRITORY - DEALSIZE -
[0]:	<pre> do d</pre>
	Out[103]: PRICEEACH
[0]:	PRICEEACH ++ ++
	Out[105]: PRICEEACH 100.0 Name: 0.95, dtype: float64 df.select('PRICEEACH').where(df.PRICEEACH>100).show()
	PRICEEACH ++ ++ Q1 = np.percentile(df.select('PRICEEACH').toPandas(), 25, interpolation = 'midpoint')
	<pre>Q2 = np.percentile(df.select('PRICEEACH').toPandas(), 50, interpolation = 'midpoint') Q3 = np.percentile(df.select('PRICEEACH').toPandas(), 75, interpolation = 'midpoint') IQR=Q3-Q1 low_lim = Q1 - 1.5 * IQR up_lim = Q3 + 1.5 * IQR df.where(df.PRICEEACH < low_lim).toPandas().shape</pre> Out[107]: (0, 23)
[0]:	<pre>df.where(df.PRICEEACH > up_lim).toPandas().shape Out[108]: (0, 23) Q4 = np.percentile(df.select('SALES').toPandas(), 25, interpolation = 'midpoint') Q5 = np.percentile(df.select('SALES').toPandas(), 50, interpolation = 'midpoint')</pre>
	Q6 = np.percentile(df.select('SALES').toPandas(), 75, interpolation = 'midpoint') IQR1=Q6-Q4 low_lim1 = Q4 - 1.5 * IQR1 up_lim1 = Q6 + 1.5 * IQR1 df.where(df.SALES > up_lim1).toPandas().shape Out[109]: (79, 23)
[0]:	<pre>sns.boxplot(df.select('SALES').toPandas()) Out[110]:</pre>
[0]:	0 2000 4000 6000 8000 10000 12000 14000 <matplotlib.axessubplots.axessubplot 0x7f71a7ded430="" at=""> df.where(df.SALES<low_lim1).topandas().shape< td=""></low_lim1).topandas().shape<></matplotlib.axessubplots.axessubplot>
[0]:	<pre>Out[111]: (0, 23)</pre> $Q7 = np.percentile(df.select('MSRP').toPandas(), 25, interpolation = 'midpoint')$ $Q8 = np.percentile(df.select('MSRP').toPandas(), 50, interpolation = 'midpoint')$ $Q9 = np.percentile(df.select('MSRP').toPandas(), 75, interpolation = 'midpoint')$ $IQR2=Q9-Q7$ $low_lim2 = Q7 - 1.5 * IQR2$ $up_lim2 = Q9 + 1.5 * IQR2$ <pre> Uf 100 100</pre>
[0]:	<pre>up_lim2 = Q9 + 1.5 * IQR2 df.where(df.MSRP > up_lim2).toPandas().shape Out[112]: (26, 23) sns.boxplot(df.select('MSRP').toPandas()) Out[113]:</pre>
	<pre>matplotlib.axessubplots.AxesSubplot at 0x7f71a7eff790></pre>
[0]:	<pre><matplotlib.axessubplots.axessubplot 0x7f71a7eff790="" at=""> df.where(df.MSRP<low_lim2).topandas().shape (0,="" (df.count(),="" 23)="" len(df.columns))<="" out[114]:="" pre=""></low_lim2).topandas().shape></matplotlib.axessubplots.axessubplot></pre>
[0]:	<pre>(df.count(), len(df.columns)) Out[115]: (2747, 23) df2 = df.where(df.SALES<up_lim1) (2668,="" (df2.count(),len(df2.columns))="" 23)<="" out[116]:="" pre=""></up_lim1)></pre>
[0]:	<pre>df2 = df2.where(df.MSRP<up_lim2) (2649,="" (df2.count(),len(df2.columns))="" 23)<="" out[117]:="" pre=""></up_lim2)></pre>
	Q10 = np.percentile(df2.select('SALES').toPandas(), 25, interpolation = 'midpoint') Q11 = np.percentile(df2.select('SALES').toPandas(), 50, interpolation = 'midpoint') Q12 = np.percentile(df2.select('SALES').toPandas(), 75, interpolation = 'midpoint') IQR3=Q12-Q10 low_lim3 = Q10 - 1.5 * IQR3 up_lim3 = Q12 + 1.5 * IQR3 df2.where(df2.SALES > up_lim3).toPandas().shape Out[118]: (17, 23)
[0]:	<pre>Out[118]: (17, 23) sns.boxplot(df2.select('SALES').toPandas()) Out[119]:</pre>
[0]:	1000 2000 3000 4000 5000 6000 7000 8000 <matplotlib.axessubplots.axessubplot 0x7f71a7e3d730="" at=""> df2 = df2.where(df2.SALES<up_lim3) (df2.count().len(df2.columns))<="" td=""></up_lim3)></matplotlib.axessubplots.axessubplot>
[0]:	<pre>df2 = df2.where(df2.SALES<up_11m3) (2632,="" (df2.count(),len(df2.columns))="" 23)="" out[120]:="" out[122]:<="" pre="" sns.boxplot(df2.select('sales').topandas())=""></up_11m3)></pre>
	1000 2000 3000 4000 5000 6000 7000 <matplotlib.axessubplots.axessubplot 0x7f71a7a43880="" at=""></matplotlib.axessubplots.axessubplot>
[0]:	Q13 = np.percentile(df2.select('SALES').toPandas(), 25, interpolation = 'midpoint') Q14 = np.percentile(df2.select('SALES').toPandas(), 50, interpolation = 'midpoint') Q15 = np.percentile(df2.select('SALES').toPandas(), 75, interpolation = 'midpoint') IQR4=Q15-Q13 low_lim4 = Q13 - 1.5 * IQR4 up_lim4 = Q15 + 1.5 * IQR4
[0]:	<pre>up_lim4 = Q15 + 1.5 * IQR4 df2.where(df2.SALES > up_lim4).toPandas().shape Out[124]: (7, 23) df2 = df2.where(df2.SALES<up_lim4) (df2.count(),len(df2.columns))<="" pre=""></up_lim4)></pre>
[0]:	Out[125]: (2625, 23) sns.boxplot(df2.select('SALES').toPandas()) Out[126]:
	1000 2000 3000 4000 5000 6000 7000 <matplotlib.axessubplots.axessubplot 0x7f71a7ee8880="" at=""></matplotlib.axessubplots.axessubplot>