Ejercicio B Last Checkpoint: hace una hora (unsaved changes) Current Kernel Logo Logout  Menu  Python 3	
Trusted  • File  • New NotebookToggle Dropdown  • Python 3	
<ul> <li>Open</li> <li>Make a Copy</li> <li>Save as</li> <li>Rename</li> <li>Save and Checkpoint</li> </ul>	
<ul> <li>Revert to CheckpointToggle Dropdown domingo, 16 de mayo de 2021 20:56</li> <li>Print Preview</li> <li>Download asToggle Dropdown asciidoc (asciidoc)</li> </ul>	
<ul> <li>html (.html)</li> <li>latex (.tex)</li> <li>markdown (.md)</li> <li>notebook (.ipynb)</li> <li>pdf (.pdf)</li> <li>rst (.rst)</li> <li>Python (.py)</li> </ul>	
<ul> <li>slides (.slides.html)</li> <li>Deploy as</li> <li>Trusted Notebook</li> <li>Close and Halt</li> <li>Edit</li> </ul>	
<ul> <li>Cut Cells</li> <li>Copy Cells</li> <li>Paste Cells Above</li> <li>Paste Cells Below</li> <li>Paste Cells &amp; Replace</li> <li>Delete Cells</li> </ul>	
<ul> <li>Undo Delete Cells</li> <li>Split Cell</li> <li>Merge Cell Above</li> <li>Merge Cell Below</li> </ul>	
<ul> <li>Move Cell Up</li> <li>Move Cell Down</li> <li>Edit Notebook Metadata</li> <li>Find and Replace</li> </ul>	
<ul> <li>Cut Cell Attachments</li> <li>Copy Cell Attachments</li> <li>Paste Cell Attachments</li> <li>Insert Image</li> </ul> <ul> <li>View</li> </ul>	
<ul> <li>Toggle Header</li> <li>Toggle Toolbar</li> <li>Toggle Line Numbers</li> <li>Cell Toolbar</li> <li>None</li> <li>Edit Metadata</li> </ul>	
<ul> <li>Raw Cell Format</li> <li>Slideshow</li> <li>Attachments</li> <li>Tags</li> <li>Insert</li> <li>Insert Cell Above</li> <li>Insert Cell Below</li> </ul>	
• Cell  • Cell  • Run Cells  • Run Cells and Select Below  • Run Cells and Insert Below  • Run Cells and Insert Below  • Run All  • Run All	
<ul> <li>Run All Below</li> <li>Cell Type</li> <li>Code</li> <li>Markdown</li> <li>Raw NBConvert</li> </ul>	
<ul> <li>Current Outputs</li> <li>■ Toggle</li> <li>■ Toggle Scrolling</li> <li>■ Clear</li> <li>Output</li> <li>■ Toggle</li> </ul>	
<ul> <li>Toggle Scrolling</li> <li>Clear</li> <li>Kernel</li> <li>Interrupt</li> <li>Restart</li> <li>Restart &amp; Clear Output</li> </ul>	
<ul> <li>Restart &amp; Run All</li> <li>Reconnect</li> <li>Shutdown</li> <li>Change kernel</li> <li>Python 3</li> </ul>	
<ul> <li>Widgets <ul> <li>Save Notebook Widget State</li> <li>Clear Notebook Widget State</li> <li>Download Widget State</li> <li>Embed Widgets</li> </ul> </li> <li>Help <ul> <li>User Interface Tour</li> </ul> </li> </ul>	
<ul> <li>Keyboard Shortcuts</li> <li>Edit Keyboard Shortcuts</li> <li>Notebook Help</li> <li>Markdown</li> </ul>	
<ul> <li>Python Reference</li> <li>IPython Reference</li> <li>NumPy Reference</li> <li>SciPy Reference</li> <li>Matplotlib Reference</li> <li>SymPy Reference</li> <li>SymPy Reference</li> </ul>	
o pandas Reference o About	
Run Code V In [26]:	
import pyspark from pyspark.sql import SparkSession	
<pre>from pyspark.sql.functions import rank, col from pyspark.sql.window import Window spark = SparkSession \     .builder \</pre>	
.appName("Python Spark SQL Opi") \ .config("spark.some.config.option", "some-value") \ .getOrCreate() In [2]:	
<pre>df = spark.read.csv("all_data.csv", header=True) In [3]:  df.show()</pre>	d  longitud +
CUADERNO FORMA IT 96 HOJAS PASTA DU    ESTRELLA  MATERIAL ESCOLAR  UTILES ESCOLARES  25.9 2011-05-18 00:00: ABASTECEDORA LUMEN  PAPELERIAS ABASTECEDORA LUME CANNES No. 6 ESQ DISTRITO FEDERAL TLALPAN 19.29699	9   -99.125417 9   -99.125417 9   -99.125417 9   -99.125417 9   -99.125417 9   -99.125417
CRAYONES CAJA 24 CERAS. TA   PAPER MATE. CARMEN MATERIAL ESCOLAR UTILES ESCOLARES 23.2 2011-05-18 00:00: ABASTECEDORA LUMEN PAPELERIAS ABASTECEDORA LUME CANNES NO. 6 ESQ DISTRITO FEDERAL TLALPAN 19.29699  PAN BLANCO BOLILLO PIEZA S/M PAN BASICOS 1.2 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICANA AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  HARINA HOT CAKES CAJA 800 GR. PRONTO TRADICIONALES GALLETAS PASTAS Y BASICOS 21.63 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICANA AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  PASTA PARA SOPA PAQUETE 200 GR. S GAMESA GALLETAS PASTAS Y BASICOS 3.45 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  GALLETAS DULCES PAQUETE 280 GR. P MARINELA GALLETAS PASTAS Y BASICOS 1.3 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  GALLETAS DULCES CAJA 752 GR. SAND NABISCO GALLETAS PASTAS Y BASICOS 41.97 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  MARINELA GALLETAS DULCES CAJA 752 GR. SAND NABISCO GALLETAS PASTAS Y BASICOS 41.97 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  MARINELA GRACA TIENDA DE AUTOSER COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  MARINELA GRACA TIENDA DE AUTOSER COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  MARINELA GRACA TIENDA DE AUTOSER COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  MARINELA GRACA TIENDA DE AUTOSER COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICANA AV. LAGO DE GUADA MÉXICO ATIZAPAN NA  MARINELA GRACA TIENDA DE AUTOSER	9   -99.125417 A   NA A   NA A   NA A   NA A   NA
SHAMPOO BOTELLA 400 ML P  HEAD & SHOULDERS ARTS. PARA EL CUI  BASICOS  49.9 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA  MÉXICO ATIZAPAN  NA 1598 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA  MÉXICO ATIZAPAN  NA 1598 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA  MÉXICO ATIZAPAN  NA 1598 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA  MÉXICO ATIZAPAN  NA 1598 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA  MÉXICO ATIZAPAN  NA 151NTE PARA EL CAB  CAJA  REVITALIQUE. 4 ARTS. PARA EL CUI  BASICOS  74 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA  MÉXICO ATIZAPAN  NA 151NTE PARA EL CAB  CAJA  WELLA KOLESTON. 40 ARTS. PARA EL CUI  BASICOS  58.9 2011-01-10 00:00: COMERCIAL MEXICANA TIENDA DE AUTOSER COMERCIAL MEXICAN AV. LAGO DE GUADA  MÉXICO ATIZAPAN  NA 151NTE PARA EL CAB  MÉXICO ATIZAPAN	AN   A AN   A AN   A AN   A
<pre>In [4]: x  df.printSchema() root   producto: string (nullable = true)</pre>	
presentacion: string (nullable = true)   marca: string (nullable = true)   categoria: string (nullable = true)   fechaRegistro: string (nullable = true)   fechaRegistro: string (nullable = true)   cadenaComercial: string (nullable = true)	
giro: string (nullable = true)   nombreComercial: string (nullable = true)   direccion: string (nullable = true)   estado: string (nullable = true)   municipio: string (nullable = true)   latitud: string (nullable = true)   longitud: string (nullable = true)	
<pre>In [5]: # Pregunta 1_a: encontrando numero de registros en la base de datos df.count() Out[5]:</pre>	
Out[5]: 62530715 In [7]: x	
<pre>df.select("categoria").distinct().show() ++  </pre>	
categoria   DETERGENTES Y PRO   CARNE Y VISCERAS   PRODUCTOS DE TEMP   GALLETAS PASTAS Y    HORTALIZAS FRESCAS    null	
DERIVADOS DE LECHE   TORTILLAS Y DERIV    GRASAS ANIMALES C    APARATOS ELECTRON    LEGUMBRES SECAS    CAFE    MUEBLES DE COCINA    CARNES FRIAS SECA	
CHOCOLATES Y GOLO  +	
<pre>#Pregunta 1_b: encontrando el numero de categorias  df.select("categoria").distinct().count() Out[8]: 42</pre>	
<pre>In [10]: #Pregunta 1_c: encontando el numero de cadenas comerciales monitoreadas df.select("cadenaComercial").distinct().count() Out[10]:</pre>	
706 In [18]: x	
<pre>conteo_productos_estado = df.groupBy("estado","producto").count() In [19]:  x  conteo_productos_estado.show() ++   estado  producto count  ++   producto count </pre>	
MÉXICO TINTE PARA EL CAB 44007    MÉXICO  TELEVISORES 29702    MÉXICO  ACELGA  7691    MÉXICO  QUESO. COTIJA  4414    DISTRITO FEDERAL  AZUCAR 18078    MÉXICO  DESENFRIOL-ITO  642    JALISCO  ARROZ 11735	
0AXACA PEDIALYTE. ELECTR  302    TLAXCALA  AGUA SIN GAS 14505   VERACRUZ DE IGNAC  TOMATE  652    MICHOACÁN DE OCAMPO  PAN DE CAJA 13003    YUCATÁN  FLAGENASE 400  313    MICHOACÁN DE OCAMPO  ECTIVA  39    YUCATÁN  SALSA CATSUP  6549    YUCATÁN  CLAVULIN  183	
YUCATÁN  CAPOTENA 271   JALISCO  FLAGENASE 400  699    HIDALGO  VERMOX  121   OAXACA  MAIZ POZOLERO  1387    OAXACA  AJO  783  +	
<pre>In [29]: #Pregunta 1_e: Encontrando los 2 productos mas moniutoreados en cada entidad window = Window.partitionBy(conteo_productos_estado['count'].desc())</pre>	
<pre>conteo_productos_estado.select('*', rank().over(window).alias('rank')).filter(col('rank') &lt;= 2).show()  +++++   estado  producto count rank  +++</pre>	
QUINTANA R00  REFRESC0 34367  2    NUEVO LEÓN  DETERGENTE P/ROPA 50307  1    NUEVO LEÓN  REFRESC0 49592  2    SINALOA  REFRESC0 33115  1    SINALOA  DETERGENTE P/ROPA 27177  2    TABASCO  REFRESCO 28754  1    TABASCO  DETERGENTE P/ROPA 26431  2	
BAJA CALIFORNIA  REFRESCO 37243  1    BAJA CALIFORNIA  DETERGENTE P/ROPA 23395  2    TLAXCALA  REFRESCO 43904  1    TLAXCALA  DETERGENTE P/ROPA 41398  2   COAHUILA DE ZARAGOZA  FUD 28613  1   COAHUILA DE ZARAGOZA  REFRESCO 26889  2    null LECHE ULTRAPASTEU  804  1    null  REFRESCO  553  2	
ESQ. SUR 125"  PAN BLANCO BOLILLO  130  1    ESQ. SUR 125"  TORTILLA DE MAIZ  2  2    COL. EDUARDO GUERRA  REFRESCO  275  1    COL. EDUARDO GUERRA  JABON DE TOCADOR  270  2  +	
<pre>In [47]: x df.createOrReplaceTempView("all_data") In [67]:</pre>	
sqlDF = spark.sql("SELECT cadenaComercial, producto, COUNT(*) AS numberProduct FROM all_data GROUP BY cadenaComercial, producto ORDER BY numberProduct desc") sqlDF.show()  +++++   cadenaComercial  producto numberProduct  +	
TORTILLERIAS TRAD  TORTILLA DE MAIZ  206950    WAL-MART  REFRESCO  182066    BODEGA AURRERA  REFRESCO  173538    BODEGA AURRERA  FUD  136876    WAL-MART  DETERGENTE P/ROPA  134237    WAL-MART  FUD  129023    SORIANA  REFRESCO  128758	
SORIANA  FUD   120610    WAL-MART LECHE ULTRAPASTEU  118766    SORIANA  DETERGENTE P/ROPA  116610    BODEGA AURRERA LECHE ULTRAPASTEU  115742    WAL-MART  JABON DE TOCADOR  107971    WAL-MART  YOGHURT  104072    WAL-MART  CERVEZA  102961    WAL-MART  DESODORANTE  102042	
WAL-MART  DESODORANTE  102042    WAL-MART  SHAMPOO  101301    WAL-MART  CHILES EN LATA  100236   MEGA COMERCIAL ME  REFRESCO  98946    BODEGA AURRERA  SHAMPOO  98942    WAL-MART  MAYONESA  97702  +	
In [77]:  #Pregunta 1_f: encontrando la cadena comercial con mas variedad de productos sqlDF.createOrReplaceTempView("sqlDF")	
sqlDF.createOrReplaceTempView("sqlDF")  cadna_conMasProductos = spark.sql('SELECT cadenaComercial, COUNT(DISTINCT(producto)) AS numberProducts FROM sqlDF GROUP BY cadenaComercial ORDER BY numberProducts desc')  cadna_conMasProductos.show()  +	
SORIANA  1059    WAL-MART  1051   MEGA COMERCIAL ME  1049    COMERCIAL MEXICANA  1036    CHEDRAUI  1026    MERCADO SORIANA  1024	
BODEGA AURRERA   1012    HIPERMERCADO SORIANA   1006     H.E.B.   1001     SORIANA PLUS   999     SORIANA SUPER   996     BODEGA COMERCIAL   979     I.S.S.S.T.E.   937     SUPERAMA   936	
S MART  851   SUPERMERCADOS SAN  849    SUMESA  848    CITY MARKET  844   FARMACIA GUADALAJARA  819    CASA LEY  808	
only showing top 20 rows	