

Customer Value

CLV stands for "Customer Lifetime Value". calculation using Spark/PySpark and FRM (Frequency, Recency, and Monetary Value)

```
In [20]: # Athena execution
spark.version
```

Calculation started (calculation_id=52c36bad-027a-5510-b92f-1f8c3fbe6e4a) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| |elapsed time = 00:00s

Calculation completed.

'3.2.1-amzn-0'

```
In [6]: # ## EMR and Athena for Spark Job already have spark session set-up
# ## EXECUTE ONLY IN LOCAL DEVELOPMENT
# import findspark
# findspark.init()

# import pandas as pd
# from pyspark.sql import SparkSession

# ## default Spark appName - se preferir
# spark = SparkSession.builder.appName('Spark3-quick-demo-app').master('local[*]').getOrCreate()
# sc = spark.sparkContext
# spark
```

Aux functions

```
In [21]: ## Aux function

def fshape(dataframe1):
    print('Shape : ', dataframe1.count(), len(dataframe1.columns))

def fhead(dataframe1, num_records=3):
    ## Show all columns - pandas dataframe
    # import pandas as pd
    # pd.options.display.max_columns = None
```

```

    return dataframe1.limit(num_records).toPandas()

def fsummary(dataframe1):
    return dataframe1.summary().toPandas()

```

Calculation started (calculation_id=b4c36bad-040f-541d-1407-f910c7254bef) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

Quick info related to the dataset

Original dataset - converted to Parquet (typical file format stored in S3)

- <https://archive.ics.uci.edu/ml/datasets/online+retail>

```

In [22]: s3_bucket = 's3://S3_BUCKET_NAME/'
parquet_file_name = 'S3_NAME_DIR/data_input/OnlineRetail__AWS.parquet'

s3_filename = s3_bucket + parquet_file_name

```

Calculation started (calculation_id=32c36bad-0567-6ef1-b954-53f4c44642f9) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

```

In [7]: ##### SETUP ACCESS - IAM ROLE
# # - sample AWSAthenaSparkExecutionRole-xpto DEV environment
# - Give the production access according with company security policies

# {
#     "Effect": "Allow",
#     "Action": [
#         "s3:*"
#     ],
#     "Resource": "*"
# },

# {
#     "Sid": "VisualEditor11",
#     "Effect": "Allow",
#     "Action": "glue:*",

```

```
#           "Resource": "*"
# },
```

```
In [23]: ## read local file
sdf = spark.read.parquet(s3_filename)
# sdf.printSchema()

fshape(sdf)
fhead(sdf)
```

Calculation started (calculation_id=94c36bad-0697-9c3b-2c93-a31836df5dbd) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

Shape : 541909 8

	InvoiceNo	StockCode	...	CustomerID	Country
0	536365	85123A	...	17850.0	United Kingdom
1	536365	71053	...	17850.0	United Kingdom
2	536365	84406B	...	17850.0	United Kingdom

[3 rows x 8 columns]

Create dataset with customer purchase history and apply CLV formula

- customer_id
- invoice_date
- revenue : monetary value

```
In [24]: sdf.createOrReplaceTempView('TB_SALES_SDF')
spark.sql('select max(TO_DATE(InvoiceDate)) as current_date_for_FRMV_CLV, current_date as not_today from TB_SALES_SDF').show()
```

Calculation started (calculation_id=6ac36bad-1c43-7f23-e50f-a121f82c4e3b) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

```
+-----+-----+
|current_date_for_FRMV_CLV| not_today|
+-----+-----+
|          2011-12-09|2023-03-12|
+-----+-----+
```

Information to understand the formula

The formula to calculates: **Customer Lifetime Value (CLV)** using the **FRM (Frequency, Recency, Monetary Value)** approach with a **discount rate of 10%** .

- monetary_value: the total monetary value spent by the customer.
- frequency: the frequency of customer purchases, i.e., how many times they made a purchase.
- recency_dt: the recency of the customer's purchases, i.e., how many days ago they made their last purchase.
- 365: the number of days in a year.
- 0.1: the discount rate used to calculate the present value of future cash flows.

The formula itself consists of three parts:

- (monetary_value / frequency): this part calculates the average value of each purchase made by the customer.
- (1 - ((recency + 1) / 365)): this part calculates the probability of the customer returning to make a purchase based on the time since their last purchase. The longer the time since the last purchase, the lower the probability of the customer returning to make a purchase.
- / (1 + discount): this part applies the discount rate to calculate the present value of future cash flows.

```
In [25]: ## formula to calculate CLV
def fnc_customer_clv_udf(monetary_value_f, frequency_f, recency_f, discount_f=0.1):
    return round ( ( (monetary_value_f / frequency_f) * (1 - ((recency_f + 1) / 365)) / (1 + discount_f) ) , 2)

## Register the formula to be used by Spark-SQL
from pyspark.sql.types import FloatType

spark.udf.register('fnc_customer_clv_udf', fnc_customer_clv_udf, FloatType())

print("Catalog Entry:")
[print(r) for r in spark.catalog.listFunctions() if "fnc_customer_clv_udf" in r.name]
```

Calculation started (calculation_id=bcc36bad-3209-2f90-8673-2ed6a6c99f91) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

Catalog Entry:

```
Function(name='fnc_customer_clv_udf', description=None, className='org.apache.spark.sql.UDFRegistration$$Lambda$4575/40860001',
isTemporary=True)
[None]
```

```
In [26]: ## Apply some filters and create the main customer purchase history as an example
sql_query_clv = """
WITH TB_SALES_V AS
(
    SELECT CustomerID as customer_id
        , COUNT(DISTINCT (InvoiceDate)) as frequency
        , DATEDIFF( current_date , MAX (InvoiceDate) ) as recency_now
        , ROUND(SUM(Quantity * UnitPrice), 2) as monetary_value
        , ROUND(avg(Quantity * UnitPrice), 2) as avg_revenue
        , MIN(InvoiceDate) as dt_first_Invoice
        , MAX(InvoiceDate) as dt_last_Invoice
        -- , ROUND(AVG(Quantity), 2) as avg_items
        -- , ROUND(SUM(Quantity), 2) as total_items
    FROM TB_SALES_SDF
    WHERE 1 = 1
        AND InvoiceDate IS NOT NULL
        AND Quantity > 0
        AND UnitPrice > 0
    GROUP BY customer_id
)
SELECT tb3.*
    , ROUND ( ( (monetary_value / frequency) * (1 - ((recency_dt + 1) / 365)) / (1 + 0.1) ) , 2) AS CLV_SQL -- discount of 0.1
    , fnc_customer_clv_udf(monetary_value,frequency,recency_dt) AS CLV_UDF
FROM (
    SELECT tb1.*
        , CAST( DATEDIFF(tb2.dt_current_date , tb1.dt_last_Invoice ) as float) as recency_dt
    FROM TB_SALES_V as tb1
    CROSS JOIN (SELECT MAX(dt_last_Invoice) AS dt_current_date FROM TB_SALES_V) tb2
) tb3
WHERE 1 = 1
    AND monetary_value > 0
    AND frequency > 0
    AND customer_id IS NOT NULL
ORDER BY monetary_value DESC
"""
```

```
sdf_clv = spark.sql(sql_query_clv)
sdf_clv.printSchema()
```

Calculation started (calculation_id=70c36bad-3fba-3876-0405-adf94009a3a6) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| |elapsed time = 00:00s

Calculation completed.

root

```
|-- customer_id: double (nullable = true)
|-- frequency: long (nullable = false)
|-- recency_now: integer (nullable = true)
|-- monetary_value: double (nullable = true)
|-- avg_revenue: double (nullable = true)
|-- dt_first_Invoice: timestamp (nullable = true)
|-- dt_last_Invoice: timestamp (nullable = true)
|-- recency_dt: float (nullable = true)
|-- CLV_SQL: double (nullable = true)
|-- CLV_UDF: float (nullable = true)
```

In [27]: `print('clv_sql and clv_udf provide the same information - just show how to implement it using 2 solutions... SQL and UDF')`
`fhead(sdf_clv)`

Calculation started (calculation_id=82c36bad-451b-27e6-dedb-6bc9a7c3787b) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| |elapsed time = 00:00s

Calculation completed.

clv_sql and clv_udf provide the same information - just show how to implement it using 2 solutions... SQL and UDF

	customer_id	frequency	recency_now	...	recency_dt	CLV_SQL	CLV_UDF
0	14646.0	51	4112	...	1.0	3555.12	3555.120117
1	16446.0	2	4111	...	0.0	76368.60	76368.601562
2	17450.0	27	4121	...	10.0	3961.80	3961.800049

[3 rows x 10 columns]

Machine Learning - Customer segmentation and plot

- Predictive Power (KI) = 0.741 and Prediction Confidence (KR) = 0.917

In [28]: `sdf_clv.createOrReplaceTempView('TB_CLV_SDF')`

Calculation started (calculation_id=56c36bad-5f27-88a5-f8f2-31933f1a67bb) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...
Progress: 0%| | elapsed time = 00:00s
Calculation completed.

```
In [29]: def ml_sql_prediction():
        text_sql_ml2 = """
        SELECT
            TB_CLV_SDF.*,
            ( CASE
                WHEN ( ((abs(`frequency` - 7.0e0) <= 10e-9) OR ( `frequency` >= 8.0e0 AND `frequency` <= 1.3e1 ) ) AND ((abs(month(`dt_
                WHEN ( ( ( `recency_now` >= 4.109e3 AND `recency_now` <= 4.113e3 ) ) AND ((abs(`frequency` - 6.0e0) <= 10e-9) OR (abs(`
                WHEN ( ((abs(year(`dt_first_Invoice`) - 2.01e3) <= 10e-9) OR ( (`dt_first_Invoice` IS NULL ) ) ) AND ((abs(`frequency` -
                WHEN ( ((abs(`frequency` - 7.0e0) <= 10e-9) OR ( `frequency` >= 8.0e0 AND `frequency` <= 1.3e1 ) ) AND ((abs(year(`dt_
                WHEN ( ( ( `frequency` > 1.0e1 AND `frequency` <= 1.14e2 ) ) ) THEN 10
                WHEN ( ( ( `recency_now` >= 4.109e3 AND `recency_now` <= 4.113e3 ) ) AND ((abs(`frequency` - 1.0e0) <= 10e-9) OR (abs(`
                WHEN ( ((abs(month(`dt_last_Invoice`) - 1.0e0) <= 10e-9) OR (abs(month(`dt_last_Invoice`) - 2.0e0) <= 10e-9) OR (abs(mon
                WHEN ( ((abs(day(`dt_last_Invoice`) - 6.0e0) <= 10e-9) OR (abs(day(`dt_last_Invoice`) - 7.0e0) <= 10e-9) OR (abs(day(`dt_
                WHEN ( ( ( (datediff(concat(year(`dt_last_Invoice`),'-',month(`dt_last_Invoice`),'-',day(`dt_last_Invoice`)),concat(year(
                WHEN ( ( ( `recency_now` >= 4.112e3 AND `recency_now` <= 4.152e3 ) OR ( `recency_now` > 4.434e3 AND `recency_now` <=
                ELSE 11
            END ) AS kc_monetary_value
        FROM TB_CLV_SDF
        """
        return text_sql_ml2
```

Calculation started (calculation_id=e2c36bad-6478-03b8-4206-a91940d6c77c) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...
Progress: 0%| | elapsed time = 00:00s
Calculation completed.

```
In [30]: ml_spark = ml_sql_prediction()

sdf_ml = spark.sql(ml_spark)

sdf_ml.printSchema()
# fhead(sdf_ml)
sdf_ml.show(3, vertical=True)
```

Calculation started (calculation_id=6cc36bad-65e0-be8e-043b-cb178704f702) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...
Progress: 0%| | elapsed time = 00:00s

Calculation completed.

root

```
|-- customer_id: double (nullable = true)
|-- frequency: long (nullable = false)
|-- recency_now: integer (nullable = true)
|-- monetary_value: double (nullable = true)
|-- avg_revenue: double (nullable = true)
|-- dt_first_Invoice: timestamp (nullable = true)
|-- dt_last_Invoice: timestamp (nullable = true)
|-- recency_dt: float (nullable = true)
|-- CLV_SQL: double (nullable = true)
|-- CLV_UDF: float (nullable = true)
|-- kc_monetary_value: integer (nullable = false)
```

-RECORD 0-----

customer_id	14646.0
frequency	51
recency_now	4112
monetary_value	200541.0
avg_revenue	137.36
dt_first_Invoice	2010-12-20 10:09:00
dt_last_Invoice	2011-12-08 00:12:00
recency_dt	1.0
CLV_SQL	3555.12
CLV_UDF	3555.12
kc_monetary_value	10

-RECORD 1-----

customer_id	16446.0
frequency	2
recency_now	4111
monetary_value	168472.49
avg_revenue	56157.5
dt_first_Invoice	2011-05-18 09:52:00
dt_last_Invoice	2011-12-09 09:15:00
recency_dt	0.0
CLV_SQL	76368.6
CLV_UDF	76368.6
kc_monetary_value	5

-RECORD 2-----

customer_id	17450.0
frequency	27
recency_now	4121
monetary_value	121321.71
avg_revenue	588.94


```

dt_first_Invoice | 2010-12-07 09:23:00
dt_last_Invoice  | 2011-11-29 09:56:00
recency_dt       | 10.0
CLV_SQL          | 3961.8
CLV_UDF          | 3961.8
kc_monetary_value | 10
only showing top 3 rows

```

In [31]: `fhead(sdf_clv,num_records=4)`

Calculation started (calculation_id=fec36bad-7ff1-bfc0-2e02-af610ef7f64c) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

	customer_id	frequency	recency_now	...	recency_dt	CLV_SQL	CLV_UDF
0	14646.0	51	4112	...	1.0	3555.12	3555.120117
1	16446.0	2	4111	...	0.0	76368.60	76368.601562
2	17450.0	27	4121	...	10.0	3961.80	3961.800049
3	18102.0	30	4111	...	0.0	3362.27	3362.270020

[4 rows x 10 columns]

In [32]: `parquet_file_name_export = 'S3_NAME_DIR/data_output/OnlineRetail__AWS.parquet'`

Calculation started (calculation_id=ccc36bad-9a1c-a2c7-57e2-bb5d12d648c3) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

In [33]: `## Export as parquet file`
`# sdf_clv.write.mode('overwrite').parquet('./data_output/OnlineRetail__AWS_FRMV.parquet')`
`s3_export_file = s3_bucket + parquet_file_name_export`
`sdf_clv.write.mode('overwrite').parquet(s3_export_file)`

Calculation started (calculation_id=22c36bad-9b64-370b-d462-f45aff9423c5) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| | elapsed time = 00:00s

Calculation completed.

Plot and Report sample

```
In [34]: sdf_ml.createOrReplaceTempView('TB_CLV_SDF_ML')
```

```
ml_rpt_sql = """
WITH TB_CLUSTER AS
(
    select kc_monetary_value as cluster_number
    , count(distinct customer_id) as customer_count
    , avg(clv_sql) avg_clv
    , avg(monetary_value) avg_monetary_value
    -- , count(*) as qty_records
    FROM TB_CLV_SDF_ML
    group by kc_monetary_value
)
SELECT cluster_number
--    , customer_count
    , ROUND( customer_count / (select sum(customer_count) from TB_CLUSTER ) * 100, 2) as percent_of_customers
    , ROUND( avg_clv, 2) as avg_clv
    , ROUND( avg_monetary_value, 2) as avg_monetary_value
FROM TB_CLUSTER tb1
order by avg_clv desc
"""

sdf_ml_rpt = spark.sql(ml_rpt_sql)
# sdf_ml_rpt.printSchema()
sdf_ml_rpt.show()
```

Calculation started (calculation_id=24c36bad-b976-b10e-19a2-4a84cf9265f4) in (session=e8c36baa-8d81-376a-9f0d-b27954f00d24). Checking calculation status...

Progress: 0%| |elapsed time = 00:00s

Calculation completed.

cluster_number	percent_of_customers	avg_clv	avg_monetary_value
5	2.34	1421.17	3358.97
10	2.44	655.37	19804.38
4	3.27	406.42	4364.46
1	2.92	381.8	4220.83
9	1.06	368.59	3599.54
8	22.0	302.7	747.79
11	0.22	292.02	2031.47
3	10.87	265.01	1024.83
7	8.5	214.13	675.68
2	8.47	201.7	1180.36
6	37.91	192.5	703.28

Plot

```
In [8]: ## Local execution only
# sdf_ml_rpt.pandas_api().plot.scatter(x='avg_monetary_value', y='avg_clv', color='cluster_number')
```

Optimization in Spark - considerations

Spark 1.x : Catalyst Optimizer and Tungsten Project (CPU, cache and memory efficiency, eliminating the overhead of JVM objects and garbage collection)

Spark 2.x : Cost-Based Optimizer (CBO) to improve queries with multiple joins, using table statistics to determine the most efficient query execution plan

Spark 3.x : Adaptive Query Execution (AQE) is an optimization technique in Spark SQL that use runtime statistics to choose the most efficient query execution plan, which is enabled by default since Apache Spark 3.2.0

- <https://spark.apache.org/docs/latest/sql-performance-tuning.html>
- three major features in AQE: including coalescing post-shuffle partitions, converting sort-merge join to broadcast join, and skew join optimization

This notebook use Spark 3.x and Adaptive Query Execution (AQE)

In []: