

Loading data and Preprocessing

```
In [2]: import pandas as pd

df_customer = pd.read_excel("C:/Users/Dave/Desktop/GDA_-_CIC_Case_Study_Rubric_and_Data_Sets/Customer_data.xlsx")
df_sales = pd.read_excel("C:/Users/Dave/Desktop/GDA_-_CIC_Case_Study_Rubric_and_Data_Sets/Sales_data.xlsx")
df_store = pd.read_excel("C:/Users/Dave/Desktop/GDA_-_CIC_Case_Study_Rubric_and_Data_Sets/Store_data.xlsx")

In [3]: print(df_customer.shape)
print(df_sales.shape)
print(df_store.shape)

(customer, s)
(403542, 6)
(87, 9)

In [4]: df_customer.head(10)

Out[4]:
```

	INSTNC_ID	CUSTOMER_ID	CITY	STATE	ZIP_CD	CUSTOMER_FIRST_PURCHASE	CUSTOMER_LAST_PURCHASE	CUSTOMER_PRIMARY_STORE_ID
0	172010167	129001864	AUSTIN	TX	78701	2021-01-04	2025-02-12	425
1	1760514191	127787349	AUSTIN	TX	78701	2021-01-07	2025-02-22	428
2	1961170347	59863515	AUSTIN	TX	78701	2021-01-02	2025-02-20	768
3	1973432551	1480074402	AUSTIN	TX	78701	2021-01-28	2025-02-22	768
4	200787243	157340401	AUSTIN	TX	78701	2021-01-04	2025-02-23	404
5	2205939131	131855205	AUSTIN	TX	78701	2021-01-28	2025-02-13	91
6	2224640375	1310887758	AUSTIN	TX	78701	2021-01-01	2025-02-23	428
7	223666233	1146718342	AUSTIN	TX	78701	2021-01-05	2025-02-14	61
8	224680179	150450368	AUSTIN	TX	78701	2021-01-01	2025-02-23	774
9	2258954439	169431318	AUSTIN	TX	78701	2021-06-29	2025-02-12	768

```
In [5]: df_sales.head(10)

Out[5]:
```

	CORP_ID	CUSTOMER_ID	ID_SHOPPING_TRANSACTION	DATE_ID	ITEM_SALES	DOLLAR_SALES
0	68	7787085	6801198078061	2025-01-19	4.191	22.26
1	425	1367615838	428011980772261	2025-01-19	11.870	112.35
2	29	1664896025	230011980772901	2025-01-19	11.875	53.26
3	404	1509734824	4040119807114616	2025-01-19	3.000	20.86
4	768	1453948029	76812284328849	2024-11-28	9.800	64.02
5	768	1325541042	76812284394467	2024-12-23	15.835	87.45
6	754	1539490902	754122384144617	2024-12-23	11.000	29.74
7	591	1534805484	59112284426517	2024-11-28	13.000	69.05
8	591	1744414479	59112284424913	2024-11-28	6.000	23.40
9	768	2202382511	76812238436573	2024-12-23	40.310	89.00

```
In [6]: df_store.head(10)

Out[6]:
```

	LOB	REGION	ID_RGN	DISTRICT	CITY	STR_FORMAT	STR_SEGMENT	CORP_ID	STORE
0	SA F/D	CNTRL TEXAS	3	AUSTIN WEST	BEE CAVE	Non-Plus	Up	404	AUS21
1	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	BASTROP	Non-Plus	Core	582	BASTROP
2	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	LULING	Non-Plus	Core	464	LULING
3	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	AUSTIN	Non-Plus	Up	639	AUS30
4	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	AUSTIN	Non-Plus	Up	465	AUS01
5	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	BUDA	Non-Plus	Up	477	BUDA
6	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	WIMBERLEY	Non-Plus	Up	708	WIMBERLEY
7	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	AUSTIN	Non-Plus	Core	229	AUS02
8	SA F/D	CNTRL TEXAS	3	AUSTIN SOUTH	LOCKHART	Non-Plus	Core	445	LOCKHART
9	SA F/D	CNTRL TEXAS	3	AUSTIN WEST	LAKEWAY	Non-Plus	Up	714	LAKEWAY

Removing (-) in Cusomter_ID

```
In [8]: df_negative = df_customer[df_customer["CUSTOMER_ID"] <= 0]
print(df_negative)

INSTNC_ID CUSTOMER_ID CITY STATE ZIP_CD CUSTOMER_FIRST_PURCHASE \
6676 2811455039 -108837950 AUSTIN TX 78703 2025-01-30
13493 235777219 -058832915 AUSTIN TX 78703 2025-01-31
17181 2745739775 -2059304829 AUSTIN TX 78704 2025-01-28
13403 2745739775 -2059304829 AUSTIN TX 78704 2025-01-28
19795 2813252539 -2012403780 AUSTIN TX 78704 2024-08-12
13487 2745737115 -2059303302 AUSTIN TX 78704 2025-01-28
32427 2357450279 -1448050200 AUSTIN TX 78704 2024-02-07
21642 2595247159 -1387213166 AUSTIN TX 78704 2024-09-20
32160 2495792453 -2021844492 AUSTIN TX 78705 2024-09-24
2718424011 -2118213699 AUSTIN TX 78705 2024-09-20
43719 2684623943 -2011348793 AUSTIN TX 78746 2024-08-17
2718424011 -2118213699 AUSTIN TX 78746 2024-09-20
47708 2448841487 -1986448324 AUSTIN TX 78746 2024-09-31
52398 2239433383 -1724871488 AUSTIN TX 78746 2021-12-12
53193 271225131 -249391849 AUSTIN TX 78746 2024-01-24

CUSTOMER_LAST_PURCHASE CUSTOMER_PRIMARY_STORE_ID
6676 2025-02-17 768.0
13493 2025-02-21 768.0
17181 2025-02-23 754.0
13403 2025-02-23 754.0
19795 2024-12-19 639.0
21642 2025-01-30 754.0
32427 2025-02-10 NaN
21643 2025-02-23 768.0
32160 2025-02-17 768.0
27184 2025-02-23 768.0
43719 2024-12-13 21.0
52398 2025-02-23 754.0
53193 2024-12-24 768.0

In [9]: df_customer["CUSTOMER_ID"] = df_customer["CUSTOMER_ID"].abs()

In [10]: print(df_customer["CUSTOMER_ID"] <= 0)

0 False
1 False
2 False
3 False
4 False
5 False
...
54349 False
54350 False
54351 False
54352 False
54353 False
Name: CUSTOMER_ID, Length: 54954, dtype: bool
```

Checking for duplicates

```
In [12]: print(df_customer.duplicated().sum())
print(df_sales.duplicated().sum())
print(df_store.duplicated().sum())

0
0
0

In [13]: duplicate_customer_id = df_customer[df_customer.duplicated(subset=["CUSTOMER_ID", "CITY", "STATE", "ZIP_CD", "CUSTOMER_FIRST_PURCHASE", "CUSTOMER_LAST_PURCHASE", "CUSTOMER_PRIMARY_STORE_ID"])]
print("Number of duplicate customer IDs: ", duplicate_customer_id.shape[0])

Duplicate Rows (Based on CUSTOMER_ID):
INSTNC_ID CUSTOMER_ID CITY STATE ZIP_CD CUSTOMER_FIRST_PURCHASE \
52 2708932013 127787349 AUSTIN TX 78701 2021-01-07
123 2708932013 127787349 AUSTIN TX 78701 2021-01-07
284 256302859 1195571380 AUSTIN TX 78701 2021-01-01
280 2732013015 1411613887 AUSTIN TX 78701 2021-01-09
283 2742328739 1374996126 AUSTIN TX 78701 2021-01-01
54349 2812476260 1239835302 AUSTIN TX 78746 2021-09-10
54350 2812476260 1468941170 AUSTIN TX 78746 2021-09-20
54351 2816805071 1452242588 AUSTIN TX 78746 2021-01-01
54352 2816894919 1234466180 AUSTIN TX 78746 2021-01-01
54353 282045005 1327643327 AUSTIN TX 78746 2021-01-02

CUSTOMER_LAST_PURCHASE CUSTOMER_PRIMARY_STORE_ID
52 2025-02-22 428.0
123 2025-02-21 754.0
284 2025-02-19 768.0
280 2025-02-20 754.0
283 2023-02-23 431.0
...
54349 2023-02-12 421.0
54350 2025-02-23 768.0
54351 2023-02-22 21.0
54352 2023-02-11 21.0
54353 2024-12-24 768.0

[4026 rows x 8 columns]
```

```
In [14]: df_customer.drop_duplicates(subset=["CUSTOMER_ID", "CITY", "STATE", "ZIP_CD", "CUSTOMER_FIRST_PURCHASE", "CUSTOMER_LAST_PURCHASE", "CUSTOMER_PRIMARY_STORE_ID"])

print(df_customer.shape)

(30292, 8)
```

Checking for nulls

```
In [16]: print(df_customer.isnull().sum())

INSTNC_ID 0
CUSTOMER_ID 0
CITY 0
STATE 0
ZIP_CD 0
CUSTOMER_FIRST_PURCHASE 6
CUSTOMER_LAST_PURCHASE 6
CUSTOMER_PRIMARY_STORE_ID 562
dtype: int64

In [17]: print(df_customer["CITY"!="ATX"].isnull().sum())

INSTNC_ID CUSTOMER_ID CITY STATE ZIP_CD CUSTOMER_FIRST_PURCHASE \
43285 2749085679 33516288 AUSTIN NaN 78746 2021-01-01
CUSTOMER_LAST_PURCHASE CUSTOMER_PRIMARY_STORE_ID
43285 2021-02-23 31.0

In [18]: df_customer["STATE"] = fillna("TX", inplace=True)

C:\Users\Dave\AppData\Local\Temp\ipykernel_12244\1009564711.py:11: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
The behavior will change in pandas 3.0. This inplace method will never work because it takes a DataFrame or Series object on which we are setting values whereas we are a copy.
For example, when using df[col].method(where, inplace=True), try using df.method(col= value, inplace=True) or df[col] = df[col].method(where, inplace=True) to perform the operation inplace on the original object.

df_customer["STATE"] = fillna("TX", inplace=True)

In [19]: print(df_customer["CUSTOMER_ID", "CUSTOMER_FIRST_PURCHASE"].isnull().sum())

INSTNC_ID CUSTOMER_ID CITY STATE ZIP_CD CUSTOMER_FIRST_PURCHASE \
21100 2736661383 2052494976 AUSTIN TX 78704 NaN
21156 2495564539 1415713180 AUSTIN TX 78704 NaN
26446 2728675935 2004244832 AUSTIN TX 78704 NaN
21123 2059304829 1381341414 AUSTIN TX 78704 NaN
43749 2715562847 1342253224 AUSTIN TX 78746 NaN
49440 2557848403 191882491 AUSTIN TX 78746 NaN

CUSTOMER_LAST_PURCHASE CUSTOMER_PRIMARY_STORE_ID
21100 NaN
21156 NaN
26446 NaN
21123 NaN
43749 NaN
49440 NaN

In [20]: print(df_customer["CUSTOMER_PRIMARY_STORE_ID"].isnull().sum())

INSTNC_ID CUSTOMER_ID CITY STATE ZIP_CD CUSTOMER_FIRST_PURCHASE \
44 268462319 1462241388 AUSTIN TX 78705 2021-01-07
81 2143900135 1401722100 AUSTIN TX 78705 2021-01-04
110 257431423 1393052413 AUSTIN TX 78701 2023-11-19
294 232077703 132300746 AUSTIN TX 78705 2021-01-24
...
41287 2368316897 1774271484 AUSTIN TX 78746 2022-09-10
54333 2374624453 1410231796 AUSTIN TX 78746 2021-04-21
54334 249770455 1224724150 AUSTIN TX 78746 2021-09-20
54310 5324732887 1336982940 AUSTIN TX 78746 2022-07-23
54394 1770149827 1211993705 AUSTIN TX 78746 2021-01-09

CUSTOMER_LAST_PURCHASE CUSTOMER_PRIMARY_STORE_ID
44 2021-02-18 NaN
110 2024-12-13 NaN
81 2024-11-25 NaN
110 2024-12-24 NaN
294 2025-01-03 NaN
...
54333 2024-11-24 NaN
54308 2024-11-24 NaN
54310 2024-07-13 NaN
54394 2024-12-09 NaN

[542 rows x 8 columns]
```

```
In [21]: df_customer.dropna(subset=["CUSTOMER_FIRST_PURCHASE"], inplace=True)

In [22]: print(df_customer.isnull().sum())

INSTNC_ID 0
CUSTOMER_ID 0
CITY 0
STATE 0
ZIP_CD 0
CUSTOMER_FIRST_PURCHASE 0
CUSTOMER_LAST_PURCHASE 0
CUSTOMER_PRIMARY_STORE_ID 556
dtype: int64
```

```
In [23]: print(df_sales.isnull().sum())

CORP_ID 0
CUSTOMER_ID 0
ID_SHOPPING_TRANSACTION 0
DATE_ID 0
ITEM_SALES 0
DOLLAR_SALES 0
dtype: int64
```

```
In [24]: print(df_store.isnull().sum())

LOB 0
REGION 0
ID_RGN 0
DISTRICT 0
CITY 0
STR_FORMAT 0
STR_SEGMENT 0
CORP_ID 0
STORE 0
dtype: int64
```

```
In [25]: print(df_store[df_store["STR_FORMAT"]!="INSTNC_ID"].isnull().sum())

LOB REGION ID_RGN DISTRICT CITY STR_FORMAT
45 SA F/D CNTRL TEXAS 3 CORE RESERVE LEANER NaN
52 SA F/D CNTRL TEXAS 3 CORE RESERVE BOUND ROCK NaN
56 SA F/D CNTRL TEXAS 3 CORE RESERVE BOUND ROCK NaN

STR_SEGMENT CORP_ID STORE
45 NaN 764 EPFC4
52 NaN 140 EPFC11
56 NaN 763 EPFC3
```

```
In [26]: print(df_store[df_store["STR_SEGMENT"]!="INSTNC_ID"].isnull().sum())

LOB REGION ID_RGN DISTRICT CITY STR_FORMAT
26 SA F/D CNTRL TEXAS 3 AUSTIN NORTH NANCOR Non-Plus
45 SA F/D CNTRL TEXAS 3 CORE RESERVE LEANER NaN
52 SA F/D CNTRL TEXAS 3 CORE RESERVE AUSTIN NaN
56 SA F/D CNTRL TEXAS 3 CORE RESERVE BOUND ROCK NaN
57 SA F/D CNTRL TEXAS 3 AUSTIN SOUTH SAN MARCOS Non-Plus

STR_SEGMENT CORP_ID STORE
26 NaN 811 NANCOR
45 NaN 764 EPFC4
52 NaN 140 EPFC11
56 NaN 763 EPFC3
57 NaN 822 SAN MARCOS2
```

```
In [27]: df_store.dropna(subset=["STR_FORMAT"], inplace=True)

In [28]: print(df_store.isnull().sum())

LOB 0
REGION 0
ID_RGN 0
DISTRICT 0
CITY 0
STR_FORMAT 0
STR_SEGMENT 2
CORP_ID 0
STORE 0
dtype: int64
```

Checking dtype formatting

```
In [30]: df_customer.dtypes

INSTNC_ID int64
CUSTOMER_ID int64
CITY object
STATE object
ZIP_CD object
CUSTOMER_FIRST_PURCHASE datetime64[ns]
CUSTOMER_LAST_PURCHASE datetime64[ns]
CUSTOMER_PRIMARY_STORE_ID float64
dtype: object
```

```
In [31]: df_customer["CUSTOMER_PRIMARY_STORE_ID"] = pd.to_numeric(df_customer["CUSTOMER_PRIMARY_STORE_ID"], errors="coerce").astype("int64")
df_customer.dtypes

INSTNC_ID int64
CUSTOMER_ID int64
CITY object
STATE object
ZIP_CD int64
CUSTOMER_FIRST_PURCHASE datetime64[ns]
CUSTOMER_LAST_PURCHASE datetime64[ns]
CUSTOMER_PRIMARY_STORE_ID int64
dtype: object
```

```
In [32]: df_sales.dtypes

CORP_ID int64
CUSTOMER_ID int64
ID_SHOPPING_TRANSACTION int64
DATE_ID int64
ITEM_SALES float64
DOLLAR_SALES float64
dtype: object
```

```
In [33]: df_store = df_store.round(2)

In [34]: df_store.dtypes

LOB object
REGION object
ID_RGN int64
DISTRICT object
CITY object
STR_FORMAT object
STR_SEGMENT object
CORP_ID int64
STORE object
dtype: object
```

Determine the day difference from First to last day purchase

```
In [36]: df_customer["time_difference_day"] = df_customer["CUSTOMER_LAST_PURCHASE"] - df_customer["CUSTOMER_FIRST_PURCHASE"]
df_customer.drop(columns=["INSTNC_ID"], inplace=True)

In [37]: df_customer.head(10)
```

	CUSTOMER_ID	CITY	STATE	ZIP_CD	CUSTOMER_FIRST_PURCHASE	CUSTOMER_LAST_PURCHASE	CUSTOMER_PRIMARY_STORE_ID	time_difference_day
0	129001864	AUSTIN	TX	78701	2021-01-04	2025-02-12	425	1500
1	127787349	AUSTIN	TX	78701	2021-01-07	2025-02-22	428	1607
2	59863515	AUSTIN	TX	78701	2021-01-02	2025-02-20	768	1510
3	1480074402	AUSTIN	TX	78701	2021-01-28	2025-02-22	768	1486
4	157340401	AUSTIN	TX	78701	2021-01-04	2025-02-23	404	1471
5	131855205	AUSTIN	TX	78701	2021-01-28	2025-02-13	91	1477
6	1310887758	AUSTIN	TX	78701	2021-01-01	2025-02-23	428	1514
7	11467342	AUSTIN	TX	78701	2021-01-05	2025-02-14	61	1501
8	155458368	AUSTIN	TX	78701	2021-01-01	2025-02-23	774	1514
9	169431318	AUSTIN	TX	78701	2021-06-29	2025-02-12	768	1324

Create new DF with the ZIP_CD [78701, 78703, 78704]

```
In [38]: zip_codes = ["78701", "78703", "78704"]
df_filtered = df_customer[df_customer["ZIP_CD"].isin(zip_codes)]

In [40]: df_filtered.head(10)
```

	CUSTOMER_ID	CITY	STATE	ZIP_CD	CUSTOMER_FIRST_PURCHASE	CUSTOMER_LAST_PURCHASE	CUSTOMER_PRIMARY_STORE_ID	time_difference_day
0	129001864	AUSTIN	TX	78701	2021-01-04	2025-02-12	425	1500
1	127787349	AUSTIN	TX	78701	2021-01-07	2025-02-22	428	1607
2	59863515	AUSTIN	TX	78701	2021-01-02	2025-02-20	768	1510
3	1480074402	AUSTIN	TX	78701	2021-01-28	2025-02-22	768	1486
4	157340401	AUSTIN	TX	78701	2021-01-04	2025-02-23	404	1471
5	131855205	AUSTIN	TX	78701	2021-01-28	2025-02-13	91	1477
6	1310887758	AUSTIN	TX	78701	2021-01-01	2025-02-23	428	1514
7	11467342	AUSTIN	TX	78701	2021-01-05	2025-02-14	61	1501
8	155458368	AUSTIN	TX	78701	2021-01-01	2025-02-23	774	1514
9	169431318	AUSTIN	TX	78701	2021-06-29	2025-02-12	768	1324

```
In [41]: df_filtered.shape

Out[41]: (120270, 9)
```

Look at Sales from the Store 754 & 768

```
In [42]: store = [754, 768]
sales_store = df_sales[df_sales["CORP_ID"].isin(store)]

In [43]: sales_store.shape

Out[43]: (162842, 6)
```

```
In [44]: sales_customer_id = pd.merge(sales_store, df_customer, on="CUSTOMER_ID", how="inner")

In [46]: sales_customer_id.head(10)
```

	CORP_ID	CUSTOMER_ID	ID_SHOPPING_TRANSACTION	DATE_ID	ITEM_SALES	DOLLAR_SALES	CITY	STATE	ZIP_CD	CUSTOMER_FIRST_PURCHASE	CUSTOMER_LAST_PURCHASE	CUSTOMER_PRIMARY_STORE_ID	time_difference_day
0	768	1459480829	76812284328849	2024-12-23	9.80	64.02	AUSTIN	TX	78703	2021-01-01	2025-02-23	768	1514
1	768	1325541042	76812284394467	2024-12-23	15.84	87.45	AUSTIN	TX	78704	2021-01-05	2025-02-22	61	1512
2	754	1539490902	754122384144617	2024-12-23	11.01	29.74	AUSTIN	TX	78704	2021-01-21	2025-01-21	754	1432
3	768	2202382511	76812238436573	2024-12-23	40.31	89.00	AUSTIN	TX	78704	2024-10-03	2025-01-01	754	120
4	768	1683296126	76811284328863	2024-11-28	8.91	10.60	AUSTIN	TX	78703	2021-01-02	2025-02-20	768	1510
5	768	1523224455	768012385020384	2024-12-23	32.30	121.59	AUSTIN	TX	78703	2021-01-22	2025-02-18	768	1488
6	768	1683250396	768122384363019	2024-12-23	42.00	162.70	AUSTIN	TX	78746	2022-12-06	2025-01-29	768	1596
7	768	1473701581	768112843400394	2024-11-28	15.47	30.90	AUSTIN	TX	78746	2021-10-06	2025-01-21	742	1706
8	768	1720494220	768122384363022	2024-12-23	36.55	87.23	AUSTIN	TX	78703	2021-12-08	2025-02-20	768	1110
9	768	1262208294	768122384366075	2024-12-23	16.00	75.24	AUSTIN	TX	78703</				