

(i) Invoice (Bill ID) → 6 digit "xxxxxx"

→ each bill has multiple items

## Single customer for unique invoices

↳ Frequency Of Purchase.

## Item per invoice

↳ Average basket size (Per transaction Avg. ~~Amount~~ ~~Value~~)

## Some invoice starts with "Cxxxxxx" are returned orders with Qty → negative.

| Work around return

| Flag return (tag final frequency if it matters)

| Net spend (Remove return if it creates difference)

| Value of item returned.

(ii) Stock Code: (Unique product ID)

## Variety a customer bought.

~~Will think more later.~~

(iii) Description (Product specification)

Can't see any use...

Have 1 idea (Advance ~~!!~~)

NLP use stock Categorise product as "Fashion", "Decor", "Hygiene".  
then encode and use... Let see later !!

Categorically No Use.

(iv) Quantity: (units of product)

## Total Qty by a customer (bought)

## Avg. qty per order. (invoice)

## Detect returns

## buyer category → Bulk, avg, light.

Trend of buying overtime, Decrease, const, Increase.

### (v) InvoiceDate

DD-MM-YYYY, HH:MM,

↓ ↓ No use as such!!  
Recency very useful.

Extract Months → Group it into Quarters over the Year.

! Check the trend

- Based on orders (Qty, invoices)
- Amt. spent

Category: Increasing, Constant, Decreasing.

Few terms:

# Recency: Pivot Date (Let say end of year)  
- Last Purchase.

→ Lower the number recent the customer.

# First Purchase - Last Purchase = Tenure.

~~Period~~

Total Purchase / Tenure = Purchase Frequency.  
(Qty)

# Avg (Gap between each purchase)

→ "Frequent" or "infrequent" buyers.

### (vi) Price (Per unit Product)

# Total Spend → Qty × Price.

# Avg Spend per Invoice

# Spend Pattern over time

Observation: For -ve Qty somewhere the price → 0  
Somewhere same as at what they bought.

Mostly no effect, but Yes 'bulk pricing' is there. ↴ no effect  
Small discount on high Qty.

(vii) Customer-ID (unique identifier for customers).

Basically we have 2 Years of Data

2009-2010 ; 2010-2011

# Note: Observation Says:

few 2010 data in Year 1, copy also present in Year 2

∴ Redundancy (Handle it)

Year 1

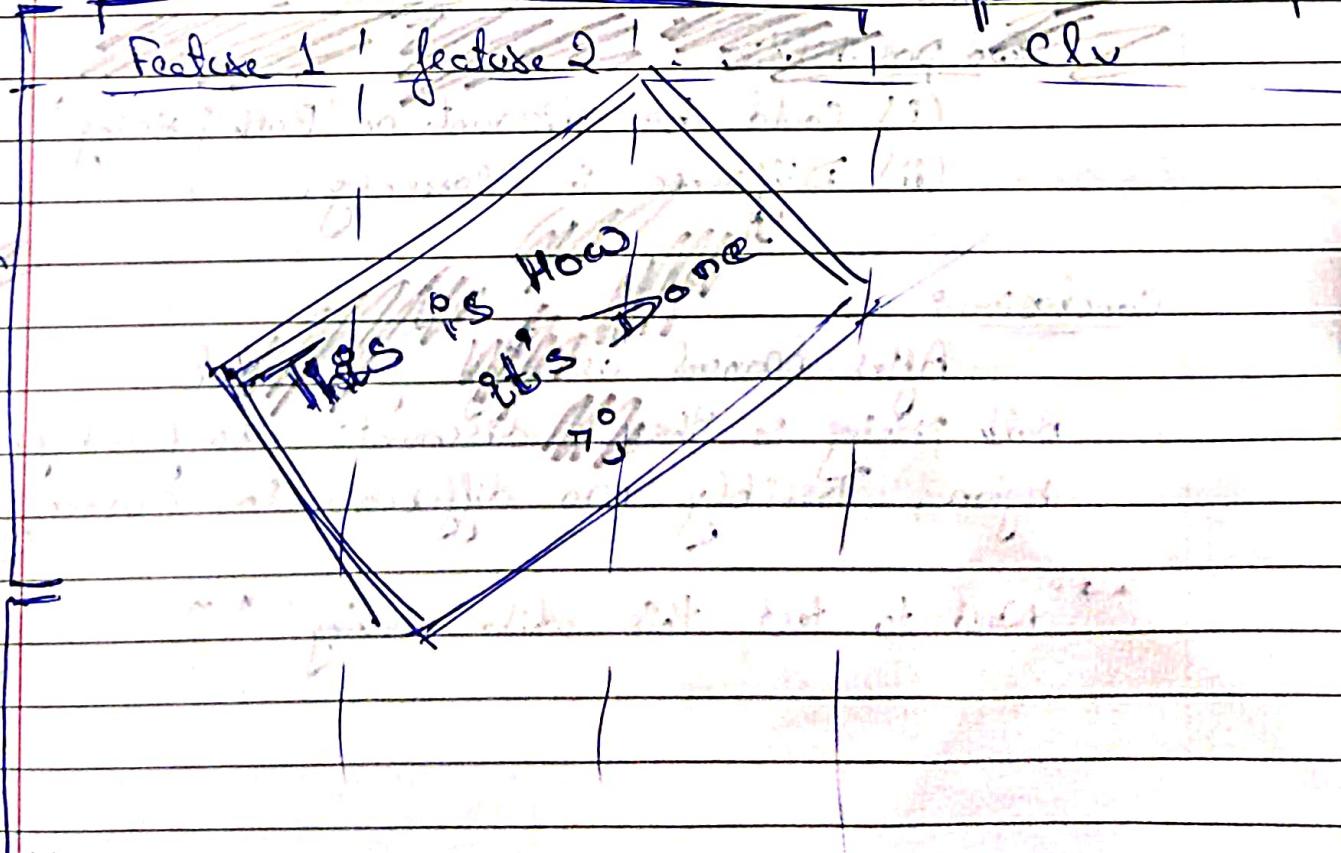
Year 2.

- Group by customer ID → Some customer Group from Year 1
- Find features → Calculate chv of Year 2
- Train model to learn it → Map each customer to its corresponding chv of Year 2

Join Year 1 and Year 2 data as:

Year 1

Year 2.



### (iii) Country : (Country of the customer)

Observation :

# Can categorize based on some 'Continent' or 'Region' wise,

Need to find correlation with CLV to say whether it impose any impact or not !.

### OBSERVATIONS :

# Manual Data Assessing

Let's go through the CSV and try observing things and note it down, other than what we have discovered yet !.

#### Observation 1 :

Currency of 'price' not mentioned.

→ Found some countries have different pricing for same product.

Two possibilities

(i) Could be discount or 'Bulk Pricing'

(ii) Difference in currency.

#### Conclusion :

After Manual assessing found that Bulk pricing is there, discounts also found for region, Possibly no difference in currency...

\*\* "Need to test this while working" \*\*



## Observation 2 :

Distribution :

## Year 1 : (2009-2010)

5,24,461 → Rows

8 → Columns (Features)

But; Out of 5,24,461 : Customer ID

are blank 'NULL'.

→ Most Blank Customer ID's are associated with

- (i) Invalid product description - '?' '???' or even null description

(ii) Null Stock Code

(iii) UnPiced → 0 (Returns)

(iv) Qty → -ve (Returns)

## Year 2 : (2010-2011)

5,41,900 → Rows

8 → Columns

Customer ID's are Blank.

Rest observation same as Year 1.

Can we do something about Customer ID...

→ If consider 'Country' → Multiple Invoice {can be of different persons}

→ Considering grouping Invoice as a single,

→ might miss other orders of same person

→ No exact CLV...

→ won't be able to map 2 customers over  
multiple years.

Therefore: Let's think more,

else last option 'dropNA' !!

Observation 3 :

## Unusual StockID

includes : Gifts, PADS, POST, M, S, SPXXX, TESTXXX

→ Gifts are i guess free items (complimentary) things.

- Gift with → Price zero(0) are free.

- Gift with → Some price are vouchers.

at end doesn't matter as none of the column

with Gift Category is associated with CustomerID  
have to drop it anyway.

→ M → Manual entries, looks artificial.

→ Mostly ~~entries~~ looks artificial.  
→ Looked deeply, best to drop.

→ PADS → It is an item? Associated with cushion?  
thus, no need to drop, it's not usual.

→ POST → Postage charges (Service charges)  
and contributes to a customer's expense.  
thus, will keep it...

→ Samples (S) → Drop (null CustomerID)

→ SPXXX → 2 valid entry

→ TESTXXX → Not real Sales, thus drop it.

Description tell its a  
Test entry possibly created by the  
one who made this dataset.

Very few entries, no effect if dropped.

## # Finalised Features :

- Customer ID
- TotalSpend (Imp)
- Purchase Frequency
- OrderHabits
- Recency
- Avg Purchase Gap
- Trend
- Churn
- Return Rate

→ CLV (Target)

## Customer ID : No use for this

## TotalSpend : → "Log-TotalSpend"

- Skewed ; Normalised (using log Transform)

## Purchase Frequency

- Skewed ;
- Normalised (Log Transformation, Box-Cox)

Reduced : Skewness

0.9

Skewness

0.08

But : Distribution is not Normal

Due to missing and more zeroes

- Binning ; Converted to Labels (Categories)

Shows : better trend, will go for it !

→ "PurchaseFreq-Label."

- IsFrequentBuyer : binary feature if Label is High or SuperBuyer

## ## OrderHabit : (Avg. Qty /- invoice)

- Highly Skewed.
- Log transformation {Bell curve achieved}
- Label - OrderHabit
  - = Light Buyer, Regular, Moderate, Heavy, Bulk
  - Gives more information.
  - quant...

Will try both features ;

## ## Recency :

- Bit skewed (Right)
- long right tail
- Log transformation  $\Rightarrow$  Skewness reduced  
but the distribution is not bell curve  
lot of missing zero values.
- BoxCox Handles 'J' well but yet not satisfied distribution.

- Best  $\Rightarrow$  Recency to Recency QRin

~~V.Rec Rec Mid Old Dormant~~

## ## Average Purchase Gap :

- Skewed.

- log transform. (Good) {less skewed, but - long tail }
- Gap\_Label (Bining).

$\hookrightarrow$  Shows better trend with "Churn, trend, C/V" . . .



## Trend → Trend-Reffined.

classified 'Flag' → New, seasonal based on features like their clu and spend trend.

Trend is discretized with 5 value (categories)

##

Churn logic:

- last 2 Quarters no spend.

- Hard, Not trustworthy.

o Pseudo churn

o

Churn (New)

- If inactive in last 2 Quarters (Pseudo churn == 1)
- Purchase frequency is  $\neq 2$ . (Not single buyer)
- clu < median clu.

Then 1 - churn

o - Not churn.

Beffer and Strong Churn

{ we already know they have produced value in clu, hence definitely not churn }

but featured it to incorporate few things ..

## Refund rate:

Lots of zeros, Can't perform log transformation.

o Binning.

/ 10  
no occs per. High

Trend is not much good!

High referrer also have High CLV.

will try and decide while No referrer have completely to keep or drop. low clv.



## # Final Set of features

- Totalspend : Log-Totalspend.
- PurchaseFrequency : PurchaseFreg\_Label.
- OrderHabit : Log-OrderHabit, Label-OrderHabit.
- Recency : Recency\_QBin.
- Avg-PurchaseGap : Log-PurchaseGap , Gap\_Label.
- Trend\_Refined.
- Churn
- ReturnRate\_Label.
- CLV. — — — Target\_Label.