Data Cleansing

Objective: Clean Up Unnecessary Data Before Exploring The Data Source

Number of Rows: 12403

Number of Fields: 23

Number of Customers: 100

Understand the Field Data

	-1-14
status	object
card_present_flag	float64
account	object
long_lat	object
txn_description	object
merchant_id	object
first_name	object
balance	float64
date	object
gender	object
age	int64
merchant_suburb	object
merchant_state	object
extraction	object
amount	float64
transaction_id	object
customer_id	object
merchant_long_lat	object
movement	object
dtype: object	

Data schema

Validate Missing Data on Each Field

```
Missing Percentage is status
card_present_flag
                    35.92
bpay biller code
                   92.65
                     0.00
account
currency
                   0.00
                     0.00
long lat
txn_description
                     0.00
merchant_id
                    35.92
merchant_code
                    92.67
                     0.00
first_name
balance
                     0.00
date
                     0.00
                     0.00
gender
                     0.00
age
merchant_suburb
                    35.92
merchant_state
                    35.92
extraction
                     0.00
                     0.00
amount
transaction_id
                     0.00
country
                   0.00
customer_id
                     0.00
merchant_long_lat
                    35.92
                     0.00
movement
```

Jupyter Noteebok (AWS SageMaker)

Eliminate 4 Unnecessary Columns

- Currency

```
AUD 12043
Name: currency, dtype: int64

df.currency.unique()

array(['AUD'], dtype=object)
```

- Country
- Bpay Biller Code

- Merchant Code

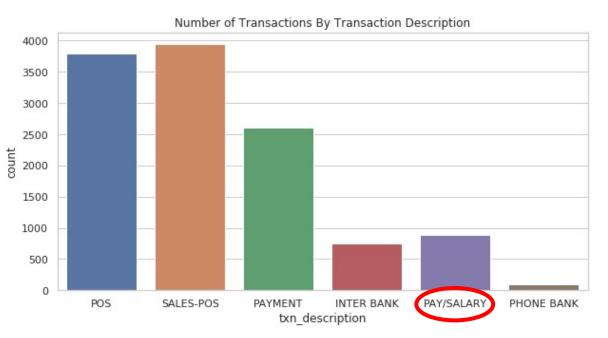
Data Discovery

Objective: Understand The Basic Data Set

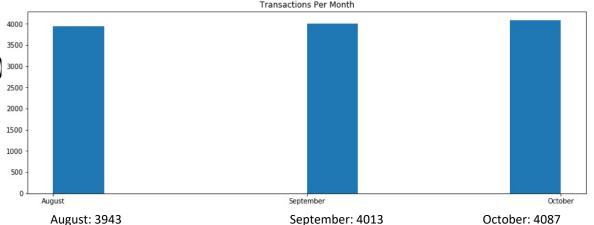
All Transactions in August 2018 - October 2018 (Seasonal) 3000

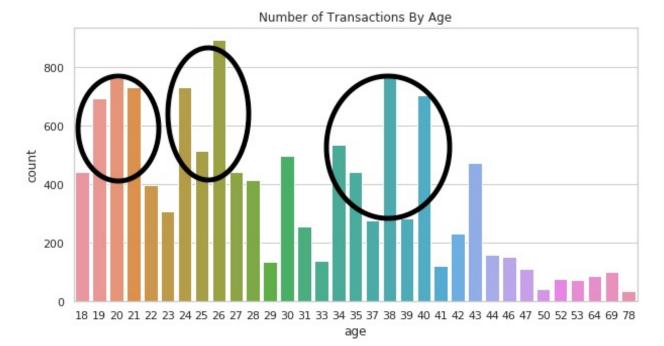
of Transactions By Gender:

Male: 6285 (9% Higher) Female: 5758

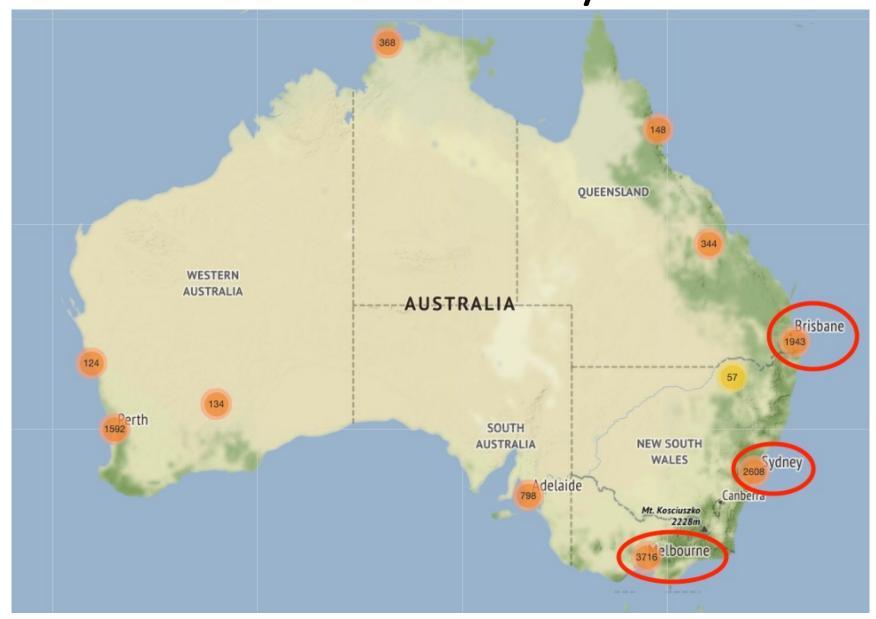


7.33% is Salary Payment (100% Credit Transactions)





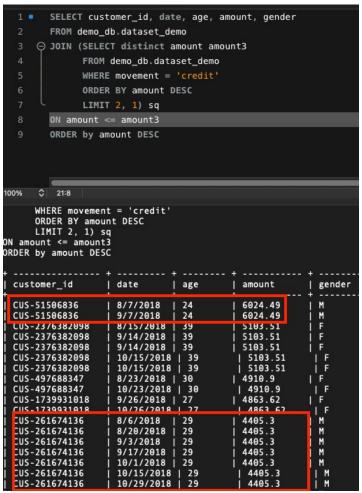
Data Discovery

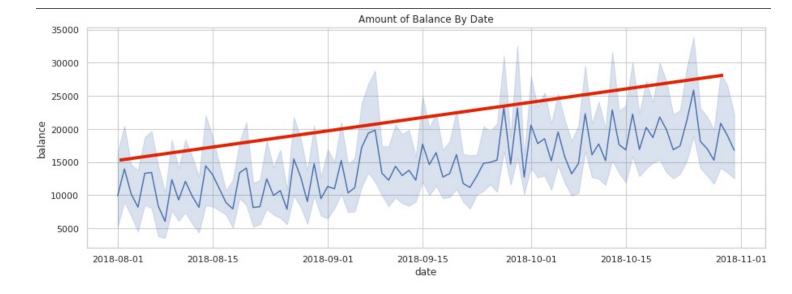


Data Insights

Objective: Analyze Data Patterns

Customers Receive Salary on Different Cycle





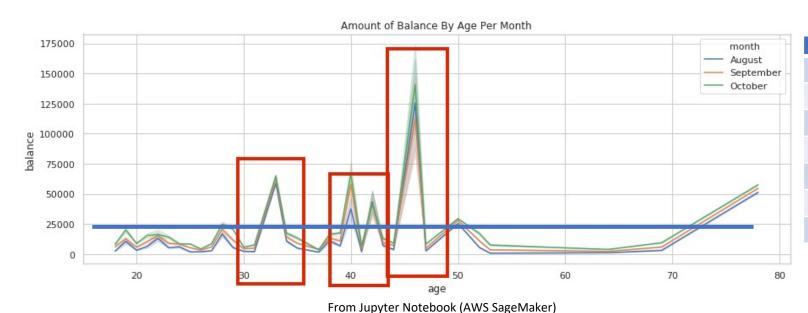
Balance Increase from AUD 150,000 to 250,000

Customers Have Been Saving (67 % Increase) Within 3 Months

Annual Salary Calculation Is Required

From MySQL (AWS RDS)

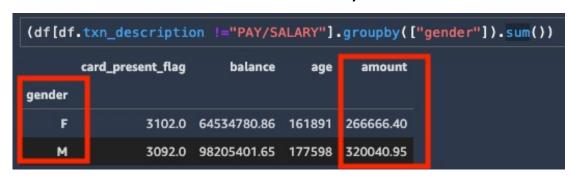
Data Insights



Age	Total Salary	Total Spending	% of Spending
10 to 19	\$175548	\$49649.279	28.28 4
20 to 29	\$635198.39 1	\$230161.259	36.23 2
30 to 39	\$449660.8 3	\$177494.85	39.47 1
40 to 49	\$335513.429	\$107910.359	32.16 3
50 to 59	\$49175.43 <mark>2</mark>	\$12859.42	26.15
60 to 69	\$22091.679	\$7175.96	32.48
70 to 79	\$9389.039	\$1456.21	15.51

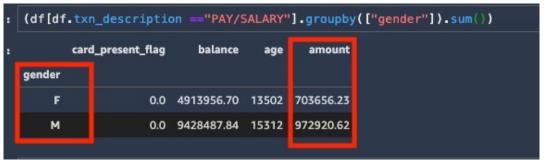
Age of 10 - 30 Spends > Age of 40 - 70

Total Amount of Balance By Gender



Men Have 20% Higher Balance

Total Amount of Salary By Gender



Men Have 38% Higher Salary
At This Stage, Gender Has a Low Impact Data Field

Data Validation

So What Now?

More Data Modelling Required For Data Examination

Reason: I. Mitigate Analytical Error (Risk)

II. Reduce Time & Cost Against Scope and Resource From Re-Cycling Previous Phases

Objective: Constant Validation Required For Accurate Solution

Approach: Use Predictive Analysis > Analyze historical data to make an accurate guess of an unknown subject

Data Heatmap

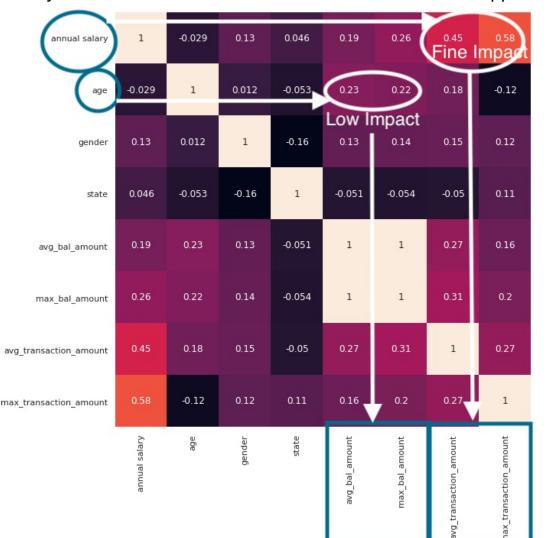
```
for count, cs id in enumerate(customers):
   df cus = df[(df.customer_id == cs_id) & (df.txn_description == 'PAY/SALARY')]
   pay_period = pd.to_datetime(df_cus.date).diff().mean() / pd.Timedelta('86400s')
   pay_amount = df_cus.amount.mean()
   daily_rate = pay_amount / pay_period
   annual_rate = daily_rate * 365
                                               When Transaction's Type is Salary/Payment,
                                                           Fetch Customers' Salary
   age = df_cus.age.mean()
                                                                        Age
                                                                      Gender
   gender = df_cus["gender"].mode()[0]
   df_cus_profile.loc[count, ["customer_id", "annual salary", "age", "gender" ]] = [id, annual_rate, age, gender]
for count, cs id in enumerate(customers):
 df cus = df[(df.customer id == cs id)]
                                                               Fetch Customers' Transaction of
 state = df_cus["merchant_state"].mode()[0]
                                                                                State
 # calculate the avg balance amount
                                                                   Avg./Max. Balance Amount
 avg_bal_amount = df_cus["balance"].mean()
                                                                Avg./Max. Transaction Amount
 # calculate the max balance amount
 max_bal_amount = df_cus["balance"].max()
 avg_trans_amount = df_cus["amount"].mean()
 # calculate the max trans amount
 max_trans_amount = df_cus["amount"].max()
 df_cus_profile.loc[count, ["state", "avg_bal_amount", "max_bal_amount", "avg_transaction_amount",
                           "max_transaction_amount"]] = [state,
                                                        avg_bal_amount, max_bal_amount,
                                                        avg_trans_amount, max_trans_amount]
```

Data Heatmap (Correlation Heatmap)

- 0.2

A Visual Exploratory Data Analysis Tool

Objective: Encounter Relational Variables From A Visual Approach



Range from 0 to 1 Closer to 1 Means Higher Correlation

High Correlation: Annual Salary x Max. Trans. Amount

Reason: Highest Amount (Numeric) is Salary

This Impacts Avg. Transaction Amount As Well

0.58 Has A Medium Correlation

Low Correlation: Age and Avg. (Max.) Balance Amount Reason: Different % of Spending By Age Range

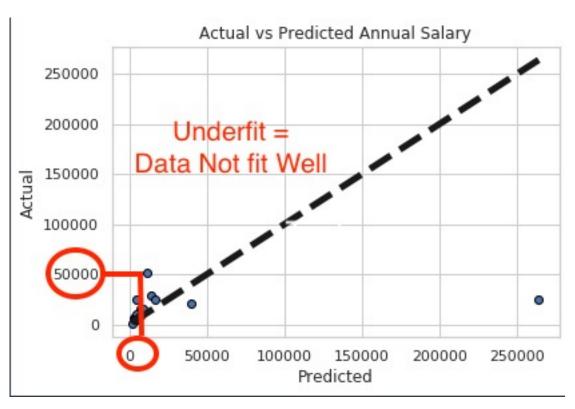
0.23 Has A Low Correlation

Linear Regression

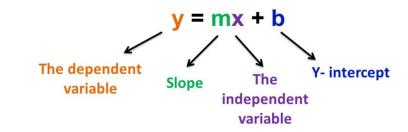
A Statistical Tool to Encounter The Relationship Between The Two Variables

Uses Algorithms in Machine Learning To Test and Predict The Outcome

Steps To Take: Build -> Train (80%) -> Test (20%) The Model

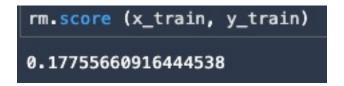


X Axis - Making Model To Predict Data



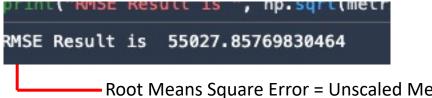
Y: Avg. Balance Amount

X: Age, Annual Salary, Avg. Trans. Amount



R^2 (Coefficient Determination = Measurement of How Close The Data Fits (How Well X Variable Explains Y Variable)

It Ranges From 0 to 1 (Scaled)
Closer to 1 Means Higher Correlation



 Root Means Square Error = Unscaled Measurement (Square Root of Residual – Measure How Far From The Regression Line The Data Points Are)

Lower the RMSE, The Better The Model Is – 55027 Has Performed Poorly

Decision Tree Regression

A Tool to Capture Data and Divide into Groups to find the Values Relationship More used to capture Non-Linear Data Shape (opposite to Linear Regression)

Steps To Take: Build -> Train (80%) -> Test (20%) The Model

Y: Avg. Balance Amount

X: Age, Annual Salary, Avg. Trans. Amount

model.score(x_train, y_train)
0.6563793326807585

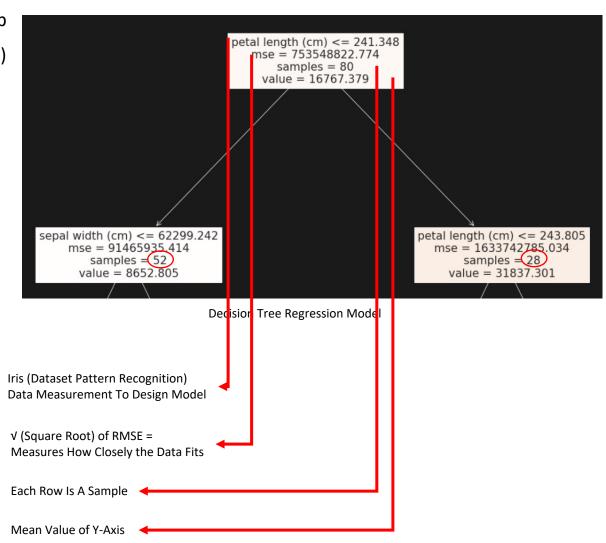
Measurement of how close the data fits 0.65 is an OK measurement

RMSE Result is ", np.sqrt(metrics.me

RMSE Result is 54387.679192649055

Root Means Square Error = Unscaled Measurement

RMSE is 54387, A Similar Result From Linear (Poorly Performed)



Conclusion

Conducted Customers' Spending Pattern Between Aug. 18 – Oct. 18 - A Seasonal Transactional Dataset

Encountered Customers Have Been Building Savings (67 % Increase of Savings)

Group of Age 10 - 30 Have Been More Active On Spending Than The Age of 40 - 70

Business Recommendation: Target Group of Age of 40 - 70 To Be More Proactive with Number of Transactions

Data Recommendation: Inadequate Analysis -> Insufficient Data

Get and Analyze Customers' Transactional Dataset of 2020 For Further Understanding About The Trend

Attempt to Fetch More Data Fields Ex: (Bank Loan – 'Accepted', 'Rejected')

Categorical Model

Higher Chance to Understand Customers' Trend of Spending

Classification Model

Grouping Data Set Into Classes to Figure Out Patterns For a Particular Event

Ex: Segment Customers Based on Purchasing History For Targeted Marketing (i.e Cashback/Rewards if Spent X Amount)

AWS Architecture

Glue: Serverless Data Integration Service (ETL)

