

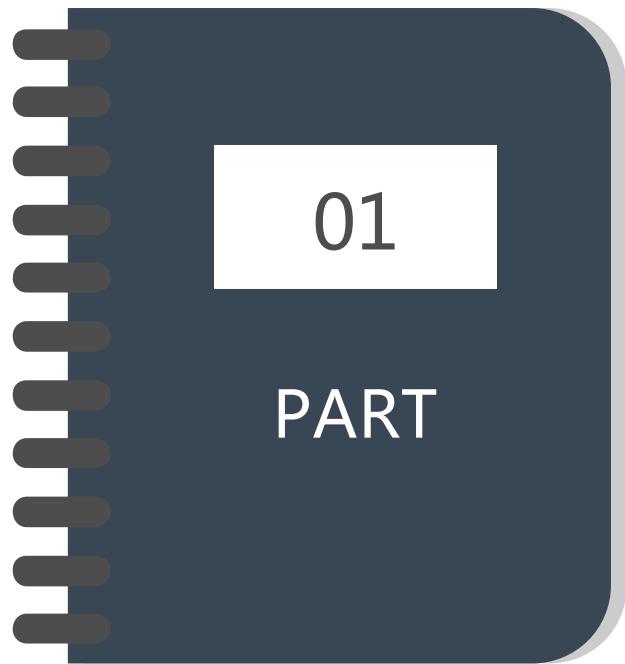
Google Merchandise Store

- REVENUE PREDICTION
- MARKETING STRATEGY



AGENDA

- 1 Business Case
- 2 Feature Engineering
- 3 Modeling
- 4 Clustering
- 5 Strategy



Business Case

Problem & Objective

The 80/20 rule has proven true for many businesses—only a small percentage of customers produce most of the revenue. But major part of the revenue comes from only 20% of the customers.



MARKETING BUDGET ALLOCATION

Marketing budget can better be utilized If they target only those users who are most likely to purchase a product in the future.



REVENUE PREDICTION

Analyze a Google Merchandise Store customer dataset to predict revenue per customer

Business Case

Feature Engineering

Modeling

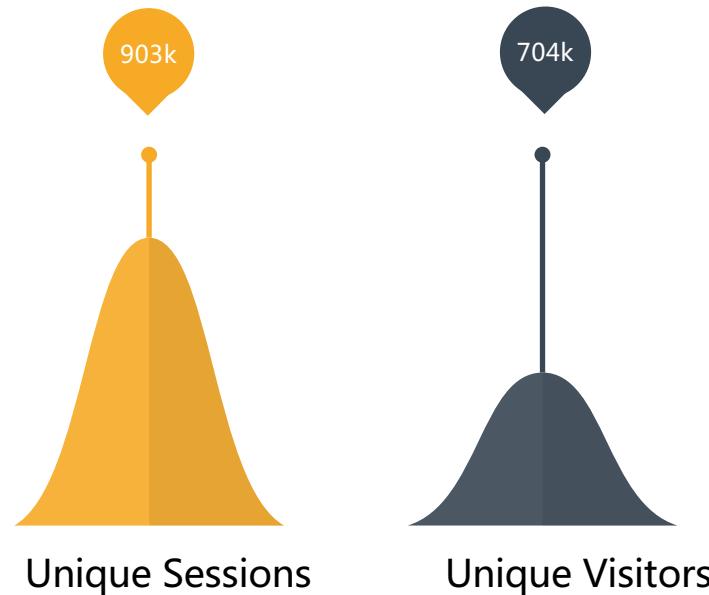
Clustering

Strategy

Data Source

Our Data is from [Google Analytics](#), GA is a tool for tracking online visitor behavior, such as their pageviews, how long do they stay and their transaction revenue.

The time spans of our data is from Aug 1, 2016 to Aug 1, 2017, it records 903,653 sessions, 704,353 unique visitors.



Unique Sessions

Unique Visitors

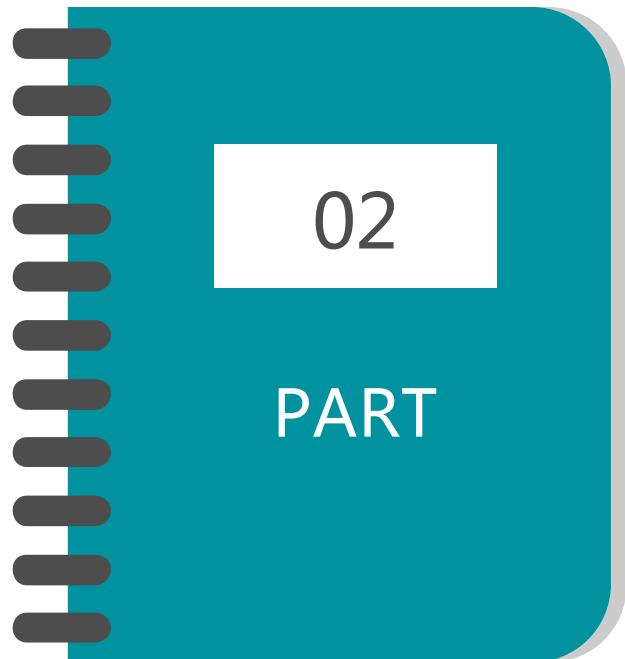
Business Case

Feature Engineering

Modeling

Clustering

Strategy

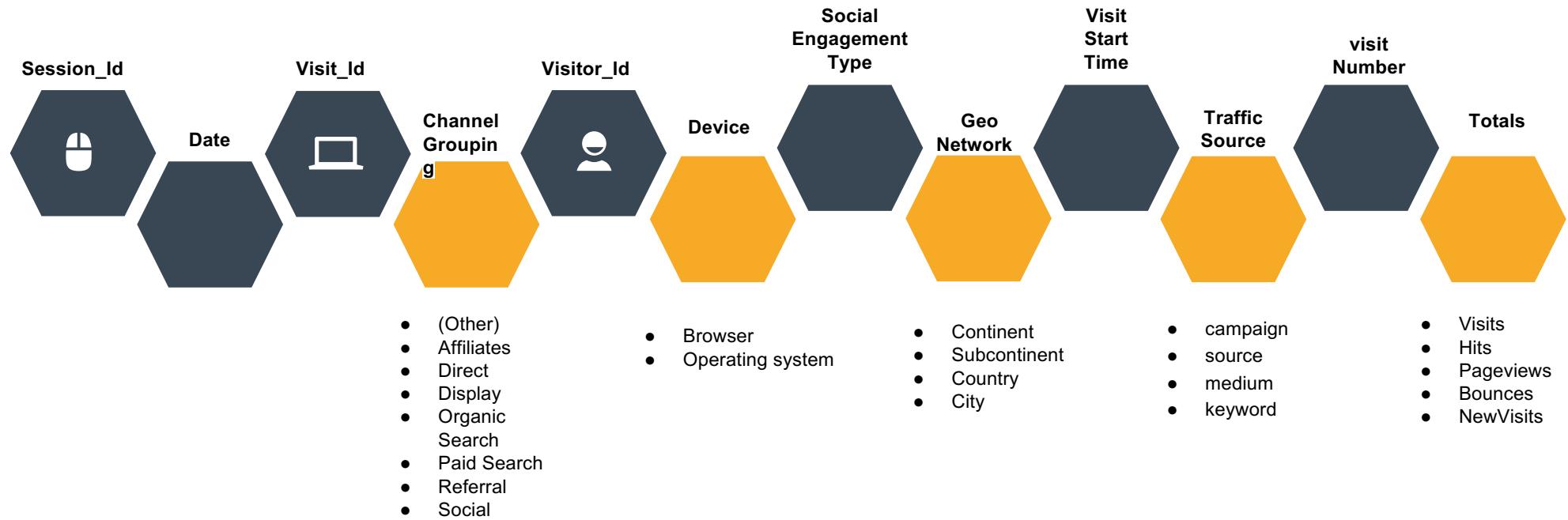


Feature Engineering

- Data Description
- Data Exploring
- Data Transformation
- EDA

Data Description

The raw data contains 12 columns, 5 of them are in the nested Json format.



Business Case

Feature Engineering

Modeling

Clustering

Strategy

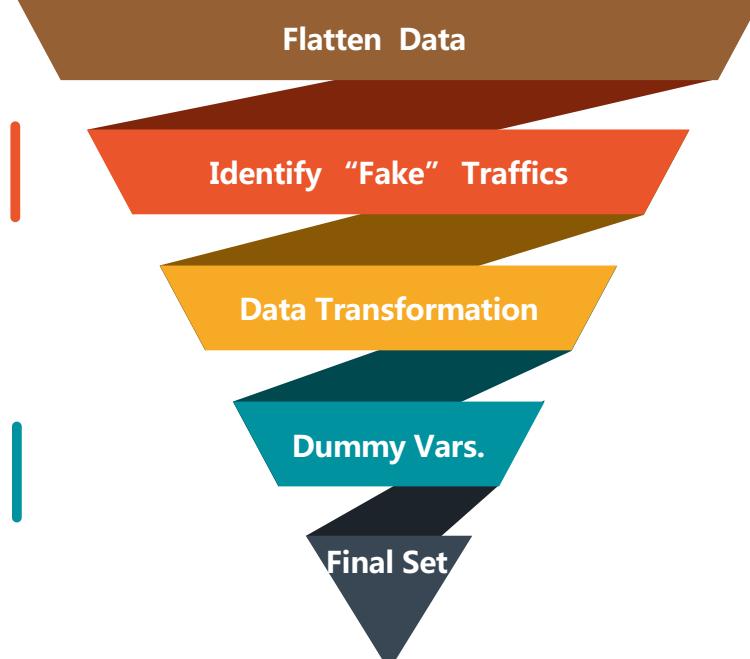
Feature Engineering

- ## Identify "Fake" traffics
- Define "fake" traffics shows strong internal testing purpose, such as traffics from Google Analytics, Github, third party traffic tracking tools.
 - Verify with transaction values
 - Deleted 33k "fake" traffics

Dummy Variables

- Many of our columns carried categorical data, therefore, we created dummy variables.
- eg.: Weekdays, Visit hour, Subcontinent, Devices etc.
- Create campaign identifier if this traffic comes from Google Adwords

Session level dataset: 903,653 rows x 95 columns
Visitor level dataset: 704,353 rows x 95 columns



Flatten Data

- Flatten the JSON columns to 59 columns.
- Delete variables that has more than 75% missing values (expect Transaction Value).
- Drop variables that have same values.
- Reduced to have 29 columns

Data Transformation

- Filling missing values using mode. Such as filled Pageview with 1. Revenue with 0
- Construct time-related variables

Final Set

- Created session level dataset
- Created visitor level dataset group by visitor id, take corresponding aggregate functions to each columns

Business Case

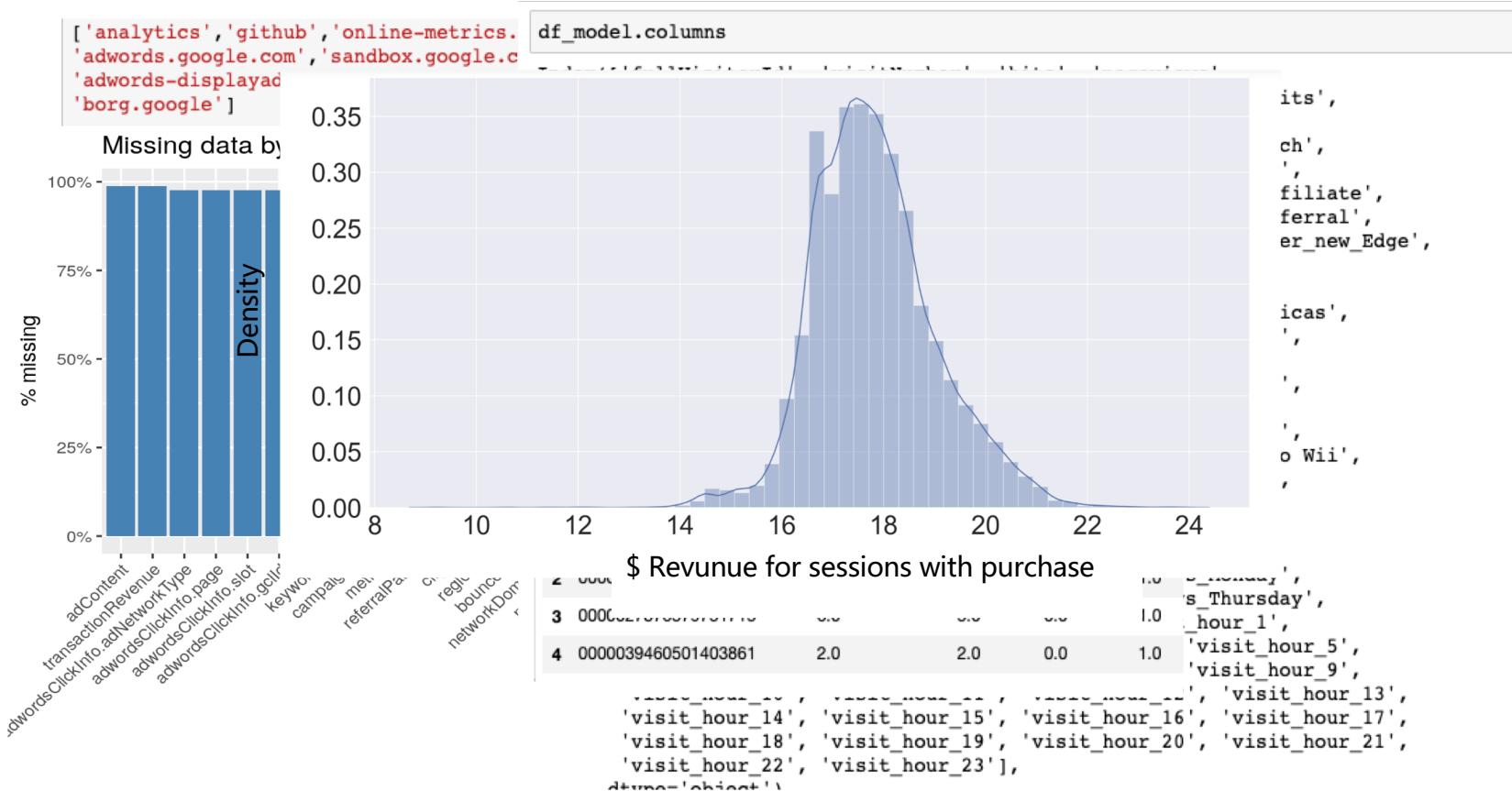
Feature Engineering

Modeling

Clustering

Strategy

EDA - Feature Engineering



columns to 59 columns.
that has more than 75%
expect Transaction Value).
that have same values.
29 columns

mation

values using mode. Such as
with 1. Revenue with 0
related variables

level dataset
level dataset group by
corresponding aggregate
h columns

Business Case

Feature Engineering

Modeling

Clustering

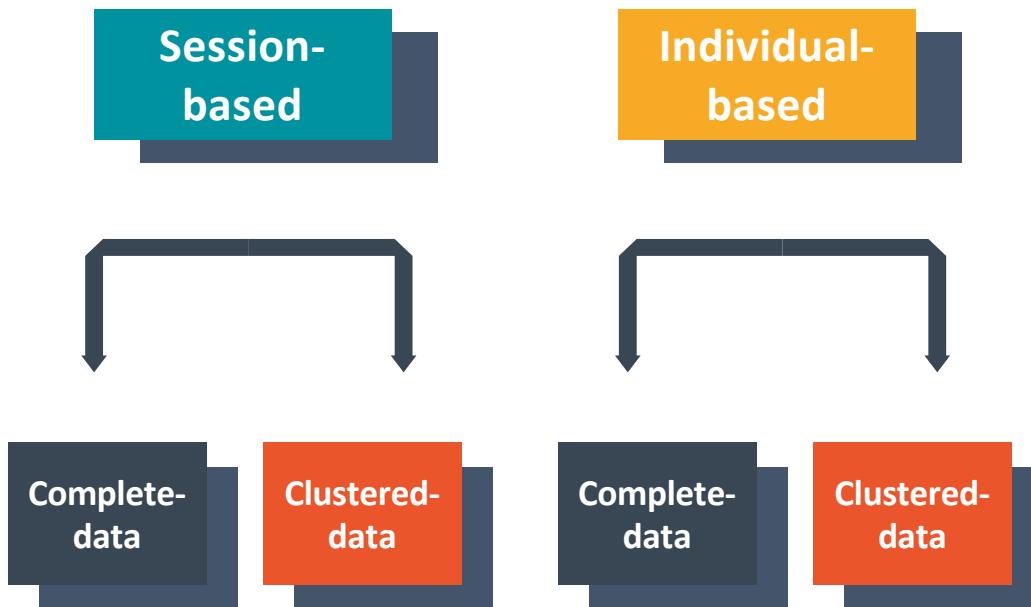
Strategy



Modeling

- ML Models
- Clustered
- Session/Individual
- Balanced

Hierarchical Modeling Approaches



Approach-1: Session-based vs. Individual-based Modeling

To serve for different purposes

Approach-2: Complete-data vs. Clustered-data Modeling

To improve models' performances

Business Case

Feature Engineering

Modeling

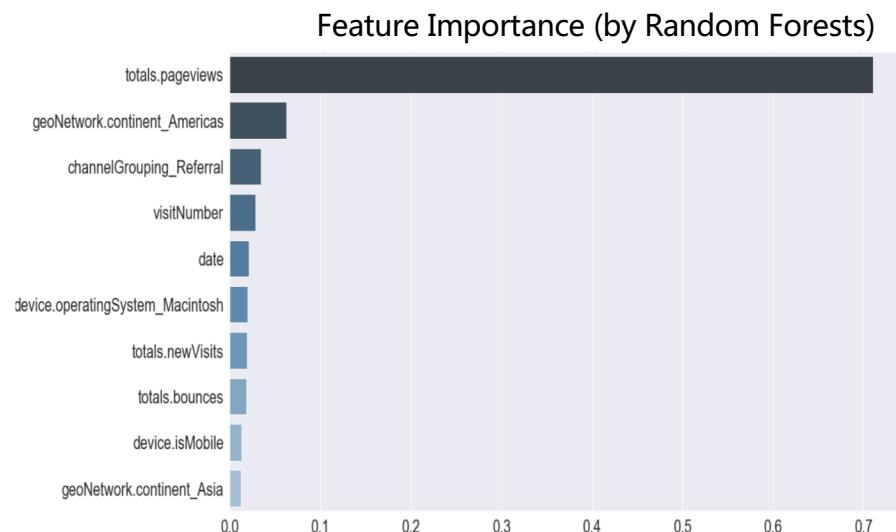
Clustering

Strategy

Approach-1: Session-based vs. Individual-based Modeling

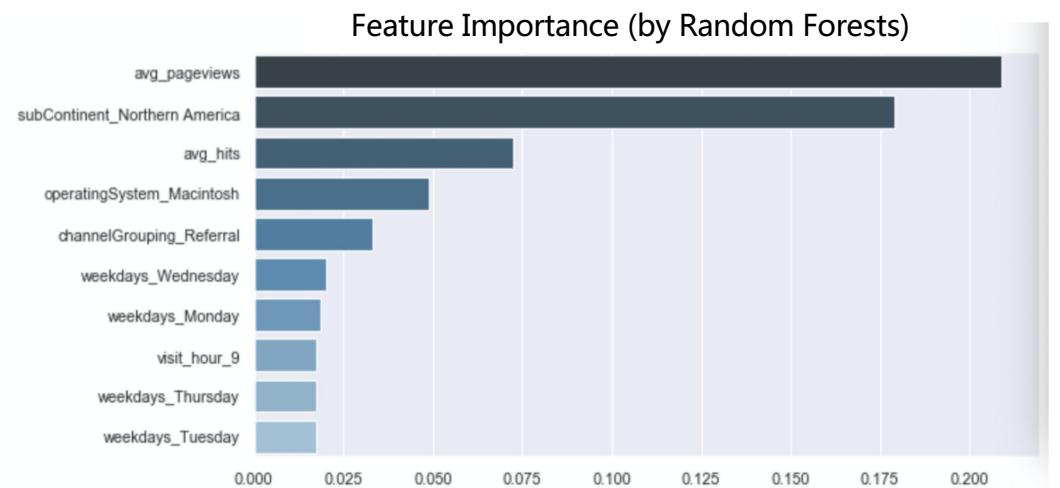
Session-based Modeling

- Original dataset to predict revenue for each session



Individual-based Modeling

- Aggregated features based on unique Individual ID (e.g. sum transaction level, average pageview)
- Most new created features' weights swell compared to original features that they are derived from



Business Case

Feature Engineering

Modeling

Clustering

Strategy

Candidate Models

	Pro	Con
Random Forest	Average result resulting in unbiased result	Took time to tune the parameters
Gradient Boosting	Corrected errors made by each sample to enhance results	May result in over-fitting
XGB Boosting	Resolved over-fitting by putting a penalized terms	Time-consuming to deliver the result
Ridge Regression	Less time to finalize result and remove multicollinearity	Results are, most of time, not as good as Machine Learning Model
Lasso Regression	Less time to finalize result and remove multicollinearity	Like Ridge Regression, is parametric method and has assumption of residual
Linear Regression	Good interpretability	Several assumptions about linear regressions should be met

Business Case

Feature Engineering

Modeling

Clustering

Strategy

Models' Performance Assessment

Session-based Models Performances

	GradientBoosting	LinearRegression	RandomForest	Ridge	XGBRegressor
All_Customer (RMSE)	1.82	2.01	1.82	2.01	1.92

- Comparison Between Session-based and Individual-based Models

- Session-based models generally outperform Individual-based models, but the differences are tiny
- It means a little amount of bias may be introduced with the creation of aggregated features

Individual-based Models Performances

	GradientBoosting	Lasso	RandomForest	Ridge	XGBRegressor
All_Customer (RMSE)	2.022375	2.370443	2.253886	2.263846	2.027313

- Comparison Within Session-based or Individual-based Models

- Tree models outperform GLM models
- Boosting result outperforms Random Forest

Business Case

Feature Engineering

Modeling

Clustering

Strategy

Approach-2: Complete vs. Clustered Modeling

Complete Modeling

- Complete dataset for all models
- **Pro:** Can be compared among candidate models directly using a single RMSE
- **Con:** Patterns of different sections of data may be hard to capture by a single model



Clustered Modeling

- Clustered datasets using the selected features
- **Pro:** Respectively modeling for different clusters to choose the best fit for each
- **Con:** Time-consuming and Heavily relying on how fit the data is clustered



Business Case

Feature Engineering

Modeling

Clustering

Strategy

Models' Performance Assessment

Complete Modeling

- (Same as the previous page)

- Session-based models

	GradientBoosting	LinearRegression	RandomForest	Ridge	XGBRegressor
All_Customer (RMSE)	1.82	2.01	1.82	2.01	1.92

- Individual-based models

	GradientBoosting	Lasso	RandomForest	Ridge	XGBRegressor
All_Customer (RMSE)	2.022375	2.370443	2.253886	2.263846	2.027313

Clustered Modeling

- RMSE varies a lot among clusters but changes little within a single cluster
- Individual-based Clustering improves model fitting results, except Cluster_4
- In cluster_4, Standard Deviation of response is much higher than other clusters due to the small sample size of this cluster

- Session-based models

	Linear Regressor	Random Forest Regressor	Ridge Regressor	XGBRegressor	GradientBoostingRegressor
1	2.173270	2.001115	2.173344	2.116877	2.011489
2	1.035442	0.990795	1.035533	1.010749	0.984474
3	3.944547	3.680210	3.944839	3.811918	3.657345
4	13.468780	8.968768	13.355924	13.633525	10.950509

- Individual-based models

	GradientBoosting	Lasso	RandomForest	Ridge	XGBRegressor
Cluster1 (RMSE)	1.899737	1.939845	1.939845	1.900526	1.766069
Cluster2 (RMSE)	0.069714	0.137477	0.091165	0.000181	0.076158
Cluster3 (RMSE)	1.351296	1.458767	1.417745	1.422674	1.350185
Cluster4 (RMSE)	16.370706	15.842980	16.278821	16.031220	16.852300

Business Case

Feature Engineering

Modeling

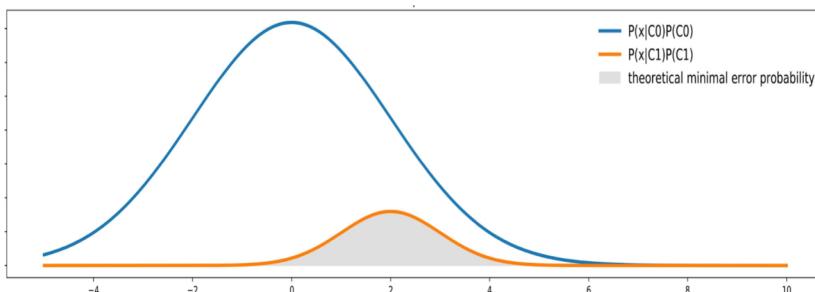
Clustering

Strategy

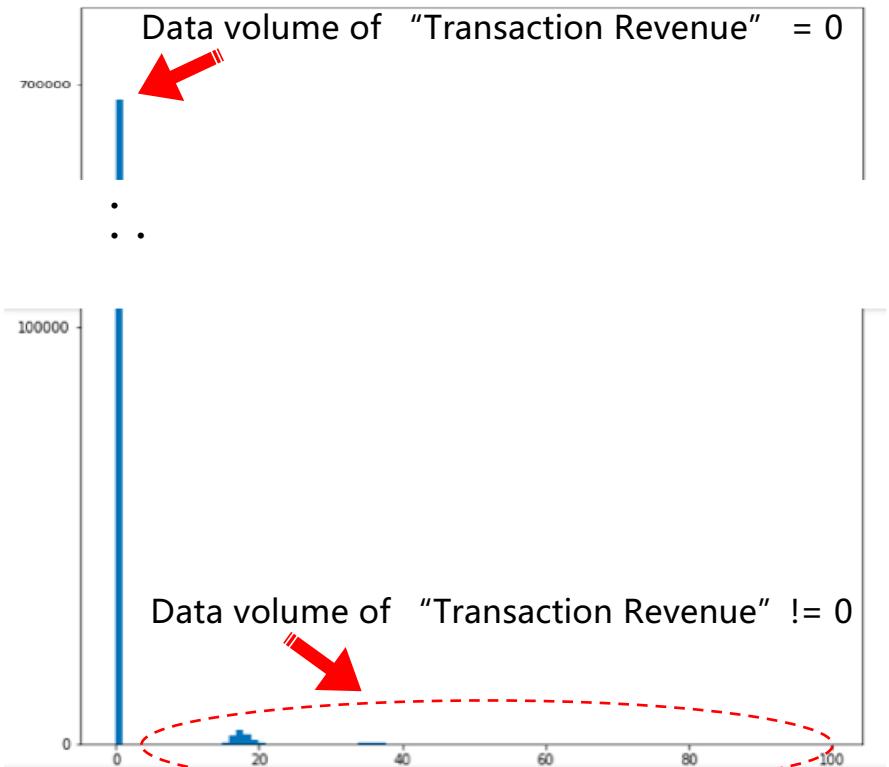
Unsolved Challenge : Data Imbalance

Problems caused by Data Imbalance

- Except the long tails, observations close to the Majority's mean are inclined to fall into its SD range
- Too small theoretical minimum error probability (area under the minority curve) to teach machines



	Zero	Non-Zero
Target Variable “Transaction Revenue”	98.6%	1.4%



Business Case

Feature Engineering

Modeling

Clustering

Strategy

Approach-3: Two-steps Modeling with Balanced Re-sampling

Step-1: Classification for Transaction

- To predict individuals who will have transactions
- Models trained with full dataset
- Given the extreme imbalance between Majority ("TR" = 0) and Minority ("TR" != 0), probability of prediction as Majority outweighs too much
- **Balanced Re-sampling is needed**



Step-2: Regression for Transaction Revenue

- To predict how much the transactions will be
- Models trained with sub-dataset of Minority ("TR" != 0)
- Process similar to the previous regression models

Comparison between Models

- Under-sampling better benefits model fitting than Over-sampling, but the diff is tiny
- LR and SVM overperform Tree models, although the diffs are not significant
- Overfitting is the root cause

	Over-sampling	Under-sampling
Random Forest	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Decision Tree	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
SVM	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Logistic Regression	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Business Case

Feature Engineering

Modeling

Clustering

Strategy

Models' Performance Assessment

	predicted "TR" = 0	predicted "TR" != 0
True label class 1	Predicted label class 1	Predicted label class 2
"TR" = 0	correct true positive for class 1	wrong false positive for class 2
True label class 2	wrong false positive for class 1	correct true positive for class 2

$$\text{accuracy} = \frac{\text{orange} + \text{blue}}{\text{orange} + \text{yellow} + \text{blue} + \text{green}}$$

$$\text{class 1 precision} = \frac{\text{orange}}{\text{orange} + \text{yellow}}$$

$$\text{class 2 precision} = \frac{\text{blue}}{\text{blue} + \text{green}}$$

$$\text{class 1 recall} = \frac{\text{orange}}{\text{orange} + \text{green}}$$

$$\text{class 2 recall} = \frac{\text{blue}}{\text{blue} + \text{yellow}}$$

Model Interpretation

- **CONSERVATIVE** for "TR = 0"

low Recall + high Precision : the model misses some real values, but those recognized values are highly reliable

- **LIBERAL** for "TR != 0"

high Recall + low Precision : most (if not all) real values are recognized, but counterfeits are mixed inside

Tradeoff

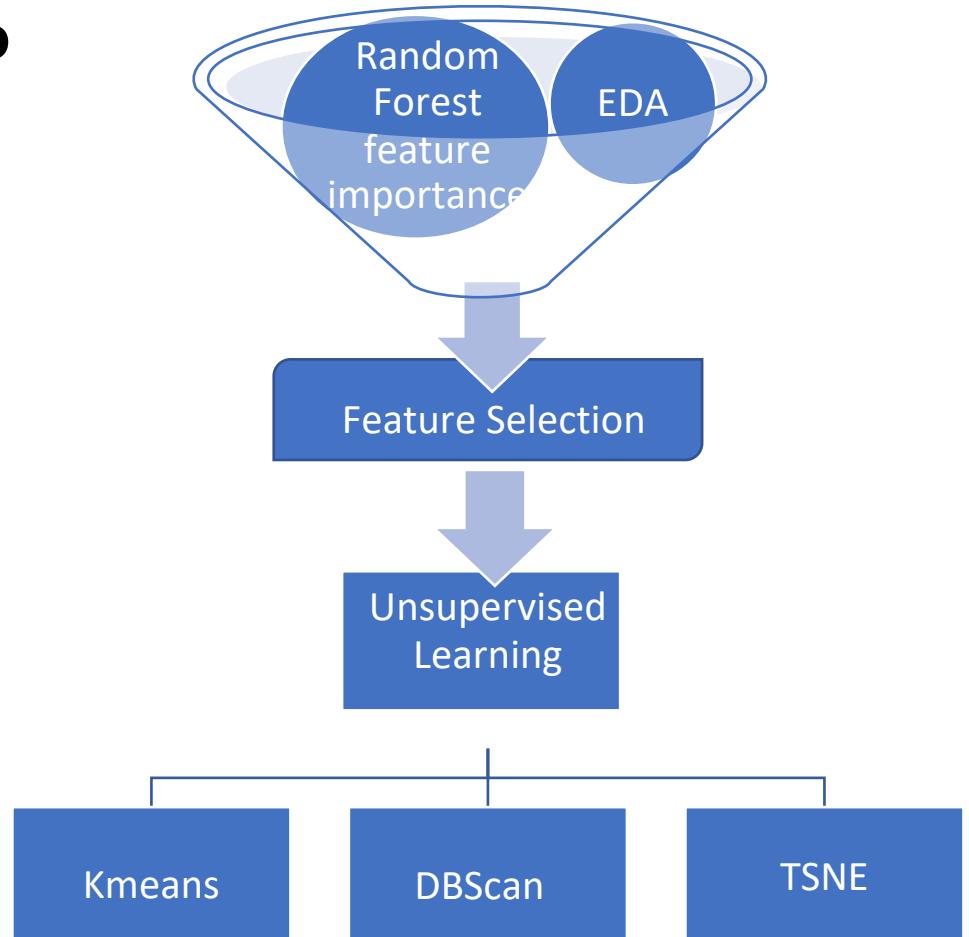
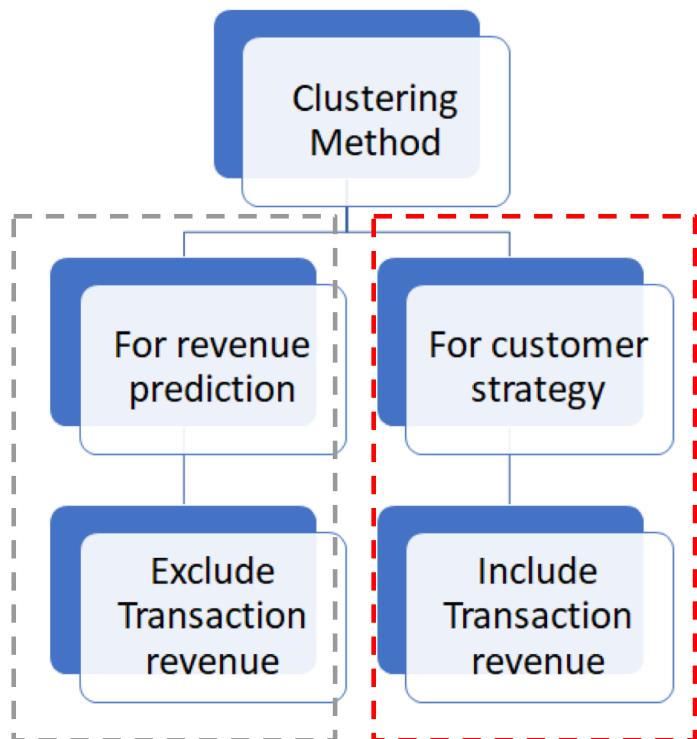
- **Targetted Customers : "TR != 0"**

	precision	recall	f1-score
0	1.00	0.93	0.96
1	0.16	0.98	0.28
accuracy			0.93
macro avg	0.58	0.95	0.62
weighted avg	0.99	0.93	0.95



Clustering

Customer Clustering for Two Purposes



Business Case

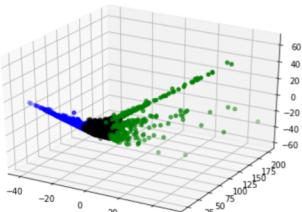
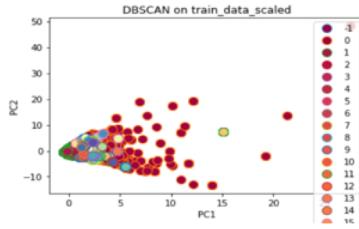
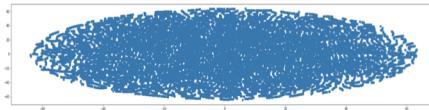
Feature Engineering

Modeling

Clustering

Strategy

Clustering Method - Models Comparison

Metrics	Kmeans		DBSCAN		TSNE	
Number of clusters:	5 clusters		> 20 clusters		1 cluster	
Visualization:						
Conclusion:	<p>Pro:</p> <ul style="list-style-type: none"> - Clear clustering - <p>Con:</p> <ul style="list-style-type: none"> - need of determination of number of cluster first 		<p>Pro:</p> <ul style="list-style-type: none"> - No need to determine the number of cluster initially <p>Con:</p> <ul style="list-style-type: none"> - Works better with non-round shape - Curse of Dimensionality 		<p>Pro:</p> <ul style="list-style-type: none"> - Dimension reduction - Good for exploration purpose <p>Con:</p> <ul style="list-style-type: none"> - No clear pattern - Points are relative to each other without concrete interpretation. 	

Business Case

Feature Engineering

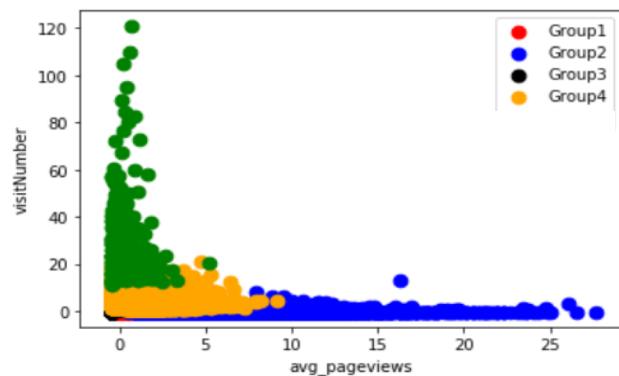
Modeling

Clustering

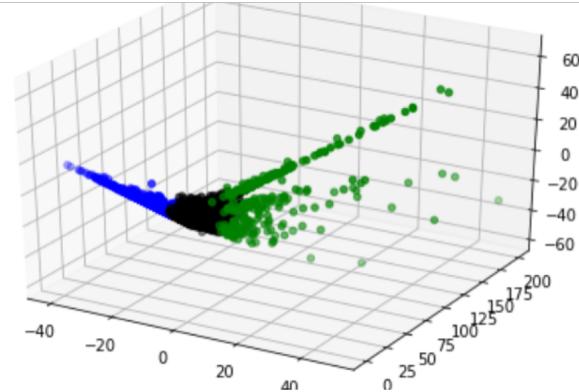
Strategy

Customer Clustering - Using the Cluster Centroids

Actual Features Graph - 2D



Principle Components Graph - 3D



Customer Type	ColorCode	TransRevenue	Pageview	VisitNumber	Bounces	Referral
One-time buyer	Red/Blue	Low	Low	Low	Low	Low
Easy-be-influenced	Orange	High	Low	High	High	High
Budget buyer	Black	Medium	High	Low	Low	Low
Technophile	Green	High	Medium	Medium	Medium	Medium

Business Case

Feature Engineering

Modeling

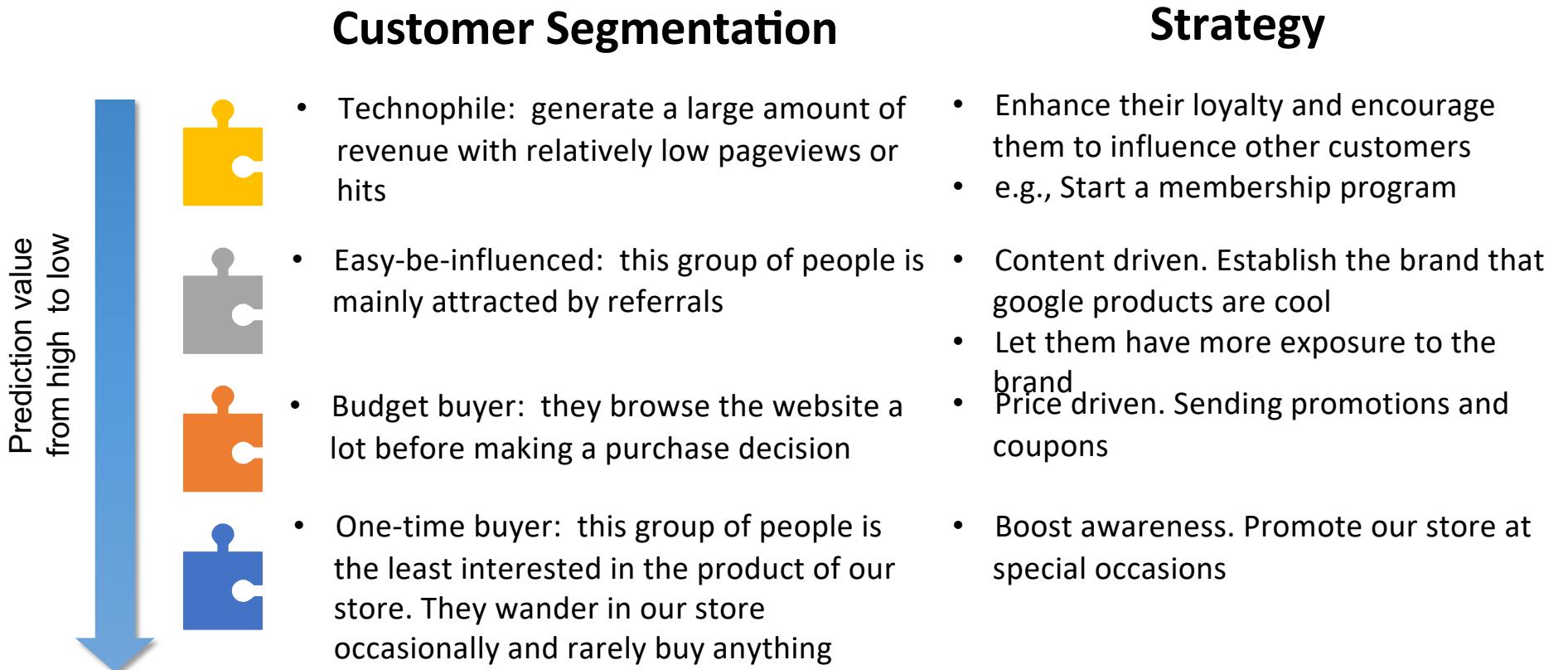
Clustering

Strategy



Strategy

Segmented Marketing Strategy



Business Case

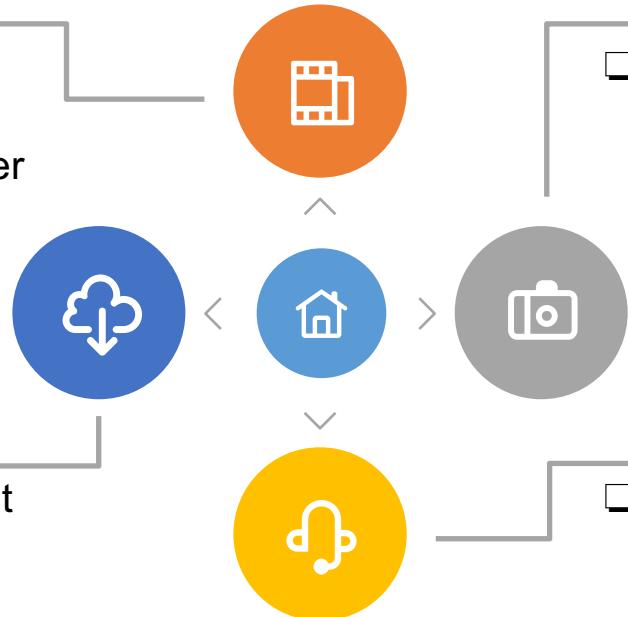
Feature Engineering

Modeling

Clustering

Strategy

With Machine Learning techniques, Google Merchandise Store can

- Identify *likely* high-value customers by analyzing existing high-value customer purchasing behavior
 - Focus product development on high-value customer products
 - Predict revenue based on customer purchasing history
 - Retain high-value customers by different marketing strategies
- 

Business Case

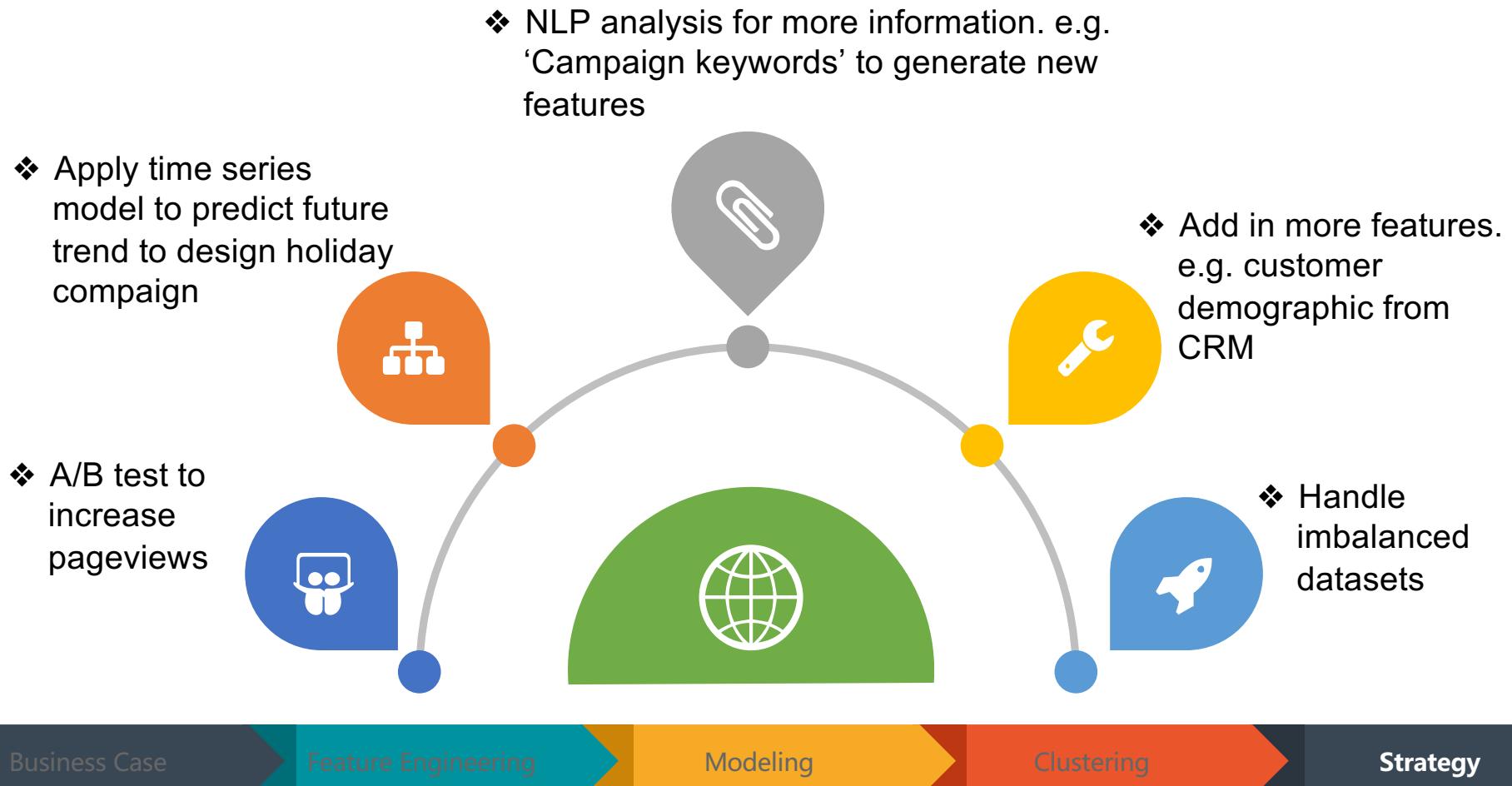
Feature Engineering

Modeling

Clustering

Strategy

Future Plans



THANK YOU

