



Departemen Statistika
IPB University



SAS Data Science Competition

Campaign Analysis
for the sake of effectiveness



Biodata



Nama

: Farah Qotrunnada

Usia

: 20 tahun

Pekerjaan

: Mahasiswa matematika tingkat 4

Institusi

: Institut Teknologi Bandung

Domisili

: Lampung

Ketertarikan :

- Text mining (on going skripsi “Automated Essay Scoring”)
- Customer Analytics
- Machine Learning

Riwayat projek :

- Supervised dan Unsupervised Kaggle
(<https://github.com/farahqotrunnada/DTI>)
- Menganalisis wajah Pendidikan Indonesia dengan machine learning (<https://www.youtube.com/watch?v=3vDYdSH3tak>)

CONTENT

1. Problem Statement

2. Workflow Penyelesaian Masalah

Problem Statement

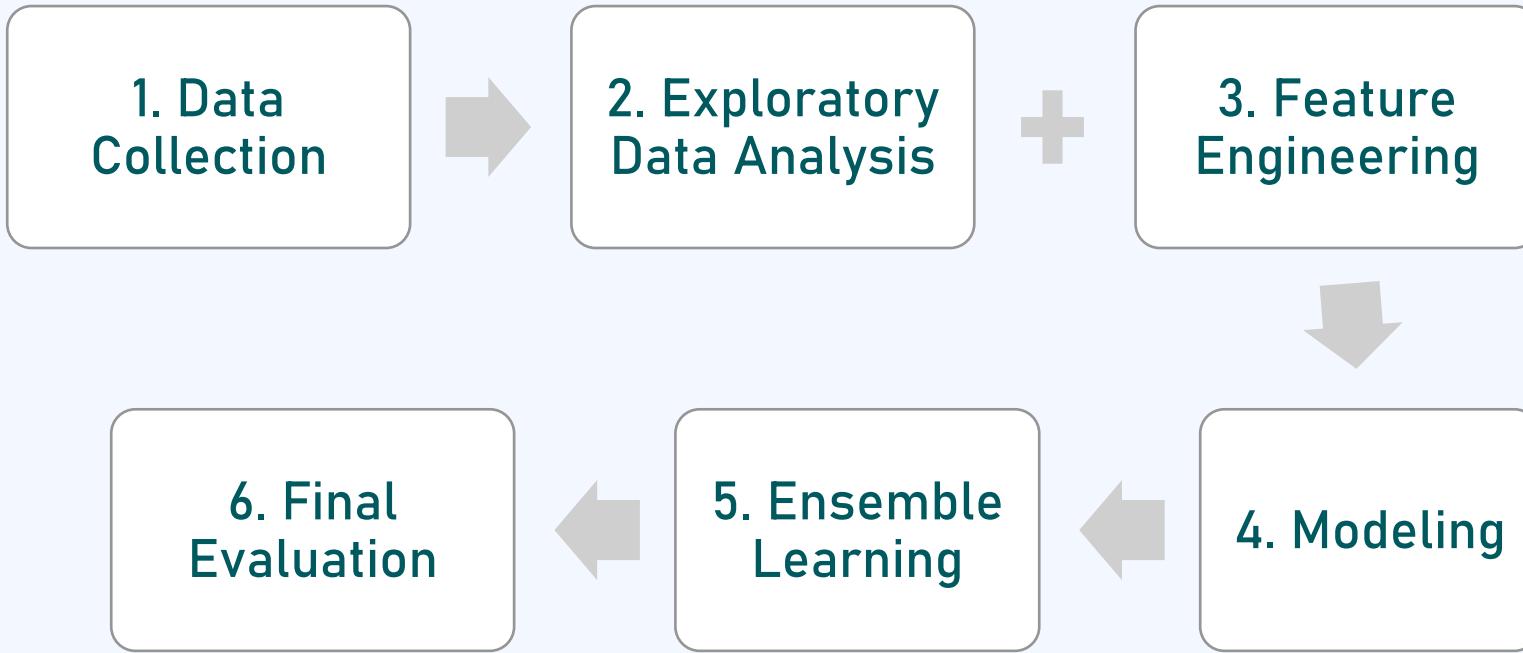
1. Masalah

Analisis kondisi karakteristik pelanggan yang mengambil promo pada initial campaign

2. Tujuan

1. Mengekstrak **informasi dan insight** dari data
2. Mengetahui **member mana** yang akan diberikan promo supaya **campaign efektif**
3. Saran bagi manajer SAS Mart terkait **keefektifan pemberian promo**

WORKFLOW PENYELESAIAN MASALAH



1. Data Collection

Sumber Data :

Soal Semifinal SAS Data Science Competition 2020

Attribute Information :

Mendefinisikan masing-masing variabel yang diketahui

2. Exploratory Data Analysis + 3. Feature Engineering

Permasalahan/hidden pattern yang ditemukan pada Exploratory Data Analysis langsung dieksekusi dengan Feature Engineering

i. Statistika Deskriptif

I. Statistika Deskriptif

	count	mean	std	min	25%	50%	75%	max
member_id	5000.0	4.916268e+05	2.886199e+05	182.0	240132.50	491138.0	740317.50	999588.0
visit_last_1mo	5000.0	3.036200e+00	1.759065e+00	0.0	2.00	3.0	4.00	10.0
visit_last_2mo	5000.0	3.737200e+00	2.375387e+00	0.0	2.00	3.0	5.00	19.0
visit_last_3mo	5000.0	3.503800e+00	2.195028e+00	0.0	2.00	3.0	5.00	14.0
spending_last_1mo	5000.0	4.523516e+05	2.659213e+05	0.0	269549.50	410076.0	585913.75	2721152.0
spending_last_2mo	5000.0	4.506261e+05	2.756227e+05	0.0	264233.75	403767.0	587337.25	2968310.0
spending_last_3mo	5000.0	4.487855e+05	2.762052e+05	0.0	262589.00	403364.5	580590.75	3227694.0
age	5000.0	2.958000e+01	5.531629e+00	18.0	26.00	29.0	33.00	59.0
monthly_income	5000.0	5.304118e+06	2.291426e+06	596700.0	3621450.00	4934000.0	6645525.00	18517900.0
buy_groceries	5000.0	1.896000e-01	3.920237e-01	0.0	0.00	0.0	0.00	1.0
buy_toiletries	5000.0	7.412000e-01	4.380193e-01	0.0	0.00	1.0	1.00	1.0
buy_food	5000.0	8.934000e-01	3.086349e-01	0.0	1.00	1.0	1.00	1.0
buy_electronic	5000.0	7.940000e-02	2.703891e-01	0.0	0.00	0.0	0.00	1.0
buy_clothes	5000.0	9.600000e-02	2.946207e-01	0.0	0.00	0.0	0.00	1.0
buy_home_appliances	5000.0	2.430000e-01	4.289380e-01	0.0	0.00	0.0	0.00	1.0
recency_last_visit	5000.0	1.610140e+01	1.151175e+01	0.0	7.00	15.0	23.00	117.0
response	5000.0	4.318000e-01	4.953765e-01	0.0	0.00	0.0	1.00	1.0

Hal-hal yang dapat di highlight dari statistika deskriptif

- a. Tidak ada data yang duplicate, karena member_id 100% unik



- b. Tidak ada data null dan blank, karena jumlah data pada tiap feature adalah 5000 sesuai dengan deskripsi soal

- b. Range data antar feature berbeda. Contohnya dapat dilihat bahwa range spending adalah sekitar 10^5 , income 10^6 , dan recency last visit 10^1-10^2 . Maka akan dilakukan standardisasi. Lebih jauh lagi semua numerical feature akan distandardisasi, karena jika variable memiliki transformasi yang sama maka algoritma akan lebih baik dalam membaca input.

d. Terdapat data yang inconsistent. Rasa penasaran terhadap bagaimana keadaan customer yang visit last 1, 2, dan 3 month nya bernilai 0, serta spending last 1, 2, dan 3 month nya juga 0, membuat hasil. Dapat dilihat dibawah bahwa customer-customer tersebut memiliki minimal 1 feature buy suatu product yang bernilai 1. Hal ini tidaklah masuk akal bahwa terdapat customer yang tidak pernah mengunjungi SAS mart namun dapat membeli produk (asumsi toko tidak melayani online buying). Maka feature engineering yang dilakukan adalah drop inconsistent data tersebut (14 row)

④ me...	④ gen...	④ visit...	④ visit...	④ visit...	④ spe...	④ spe...	④ spe...	④ buy...						
884...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
58521	Male	0	0	0	0	0	0	0	0	1	0	0	0	0
985...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
535...	Female	0	0	0	0	0	0	0	0	1	0	0	0	1
413...	Male	0	0	0	0	0	0	0	1	1	0	0	0	0
405...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
851...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
168...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
259...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
849...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
605...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
401...	Male	0	0	0	0	0	0	0	1	1	0	0	0	0
474...	Female	0	0	0	0	0	0	0	1	1	0	0	0	0
490...	Female	0	0	0	0	0	0	0	1	0	0	1	1	1

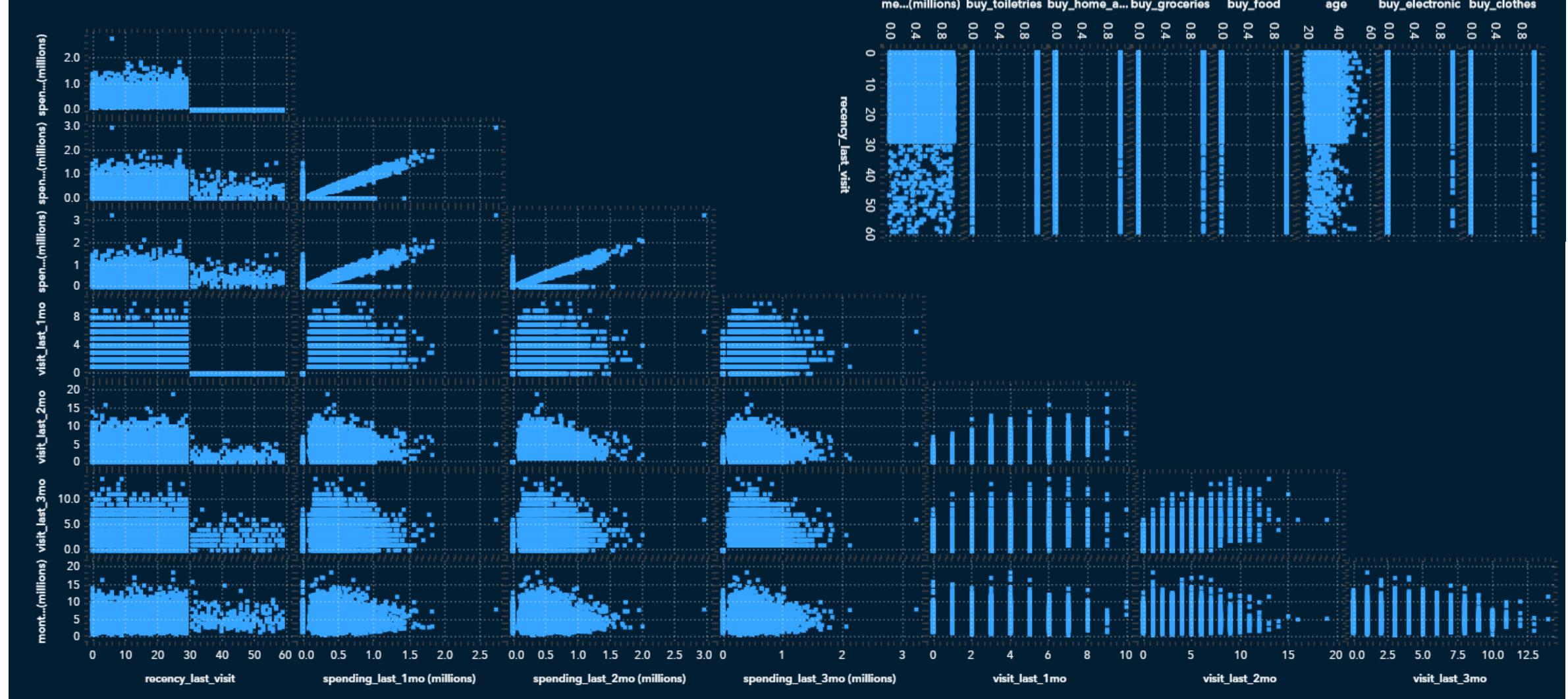
2. Exploratory Data Analysis

+

3. Feature Engineering

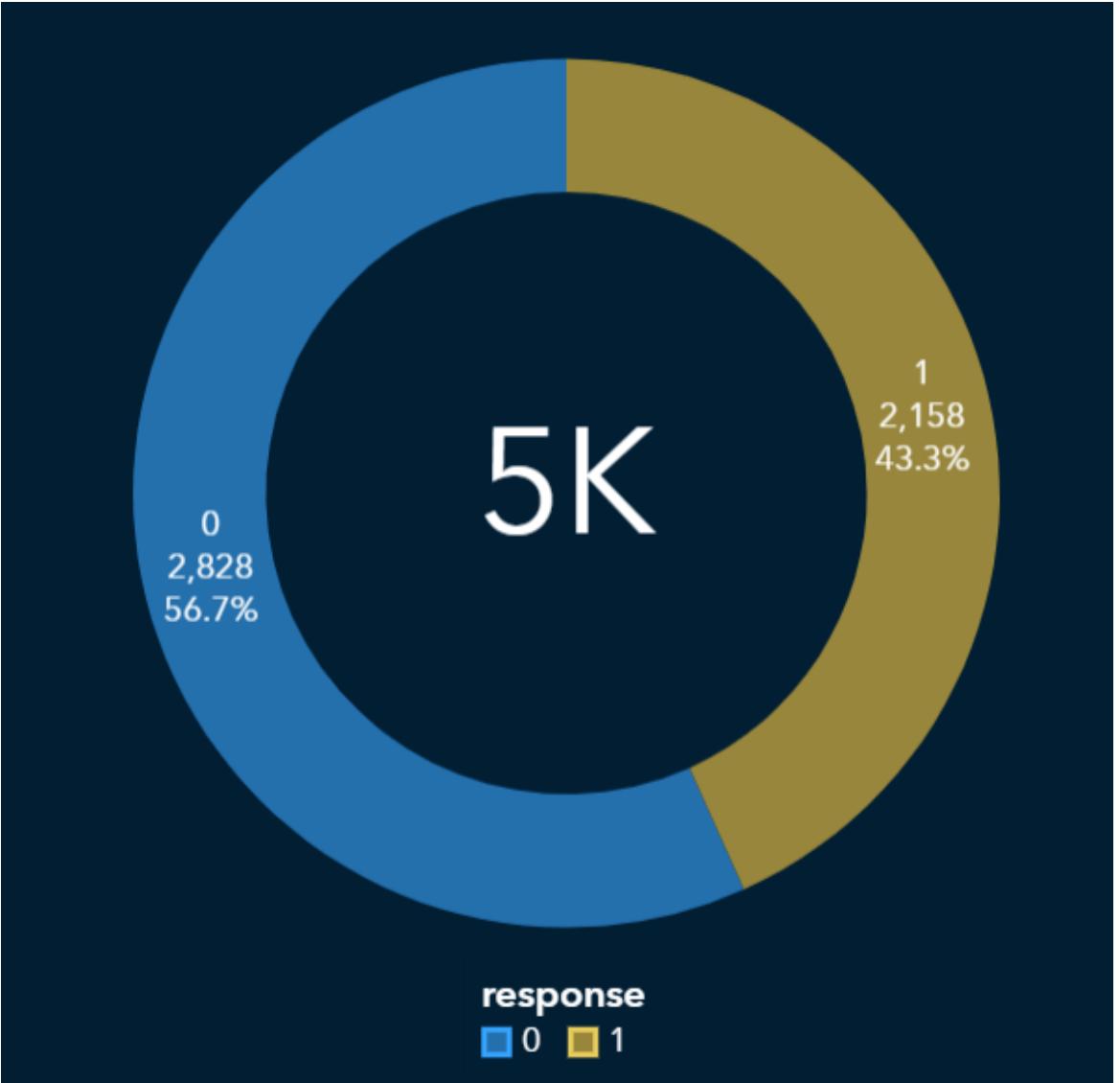
ii. Visualisasi

Scatter Plot of Selected Measures



- Fokus pada feature recency_last_visit, seperti ada suatu perbedaan yang cukup terlihat antara recency_last_visit <=30 dan >30 pada plot bersama setiap feature lain. Maka kita dapat mencoba untuk membuat feature baru recency_convert dengan recency_last_visit <=30 di assign 1 dan <30:0
- Spending last 1,2,3 month saling memiliki trend yang positif

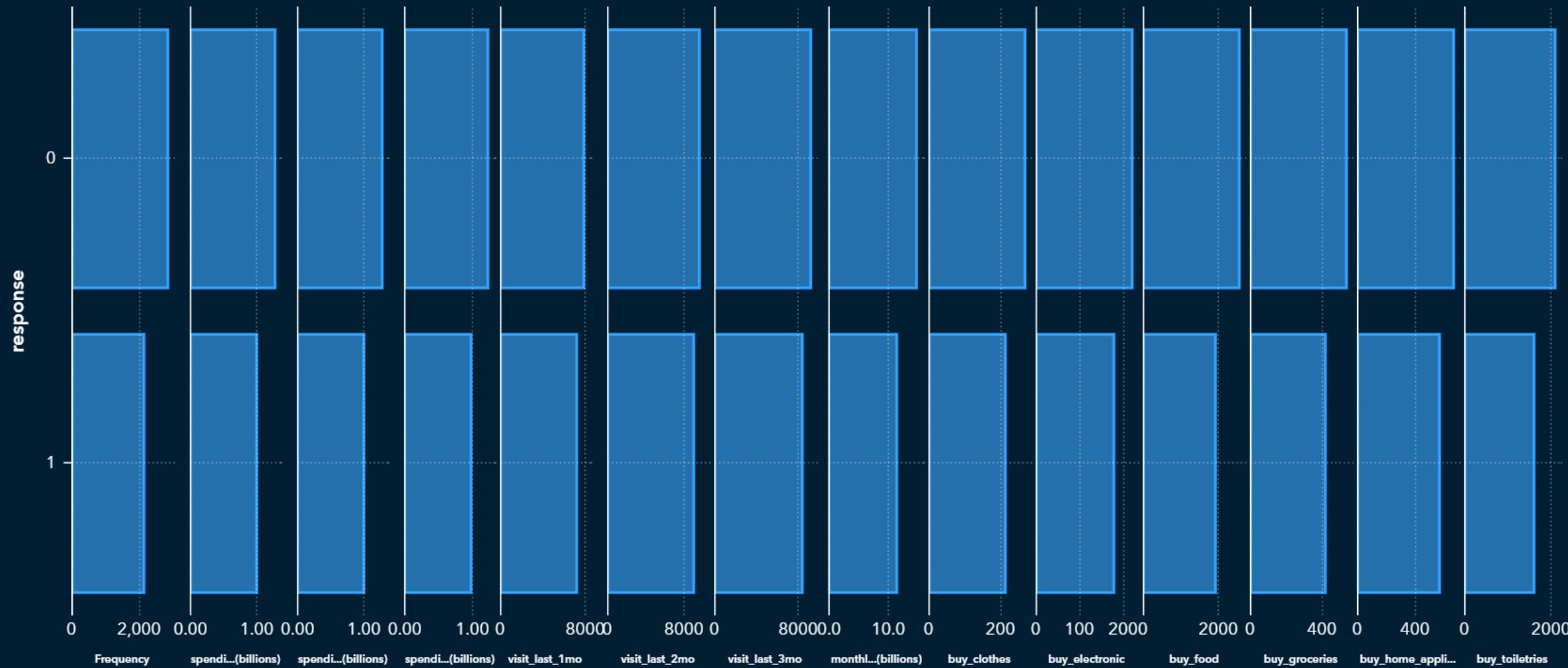
Frequency of response



Dapat dilihat bahwa perbandingan antar kelas pada response kita adalah 56.7:43.3 maka respons kita cukup seimbang, sehingga tidak diperlukan balancing kelas.

Karena hanya terdapat perbedaan yang sangat sedikit sekali bukannya 1:5, atau 1:1000. Dan teknik balancing sendiri dapat menghasilkan hasil yang kurang baik seperti undersample (dapat menghilangkan observasi yang ternyata penting) atau oversample (dapat menimbulkan bias)

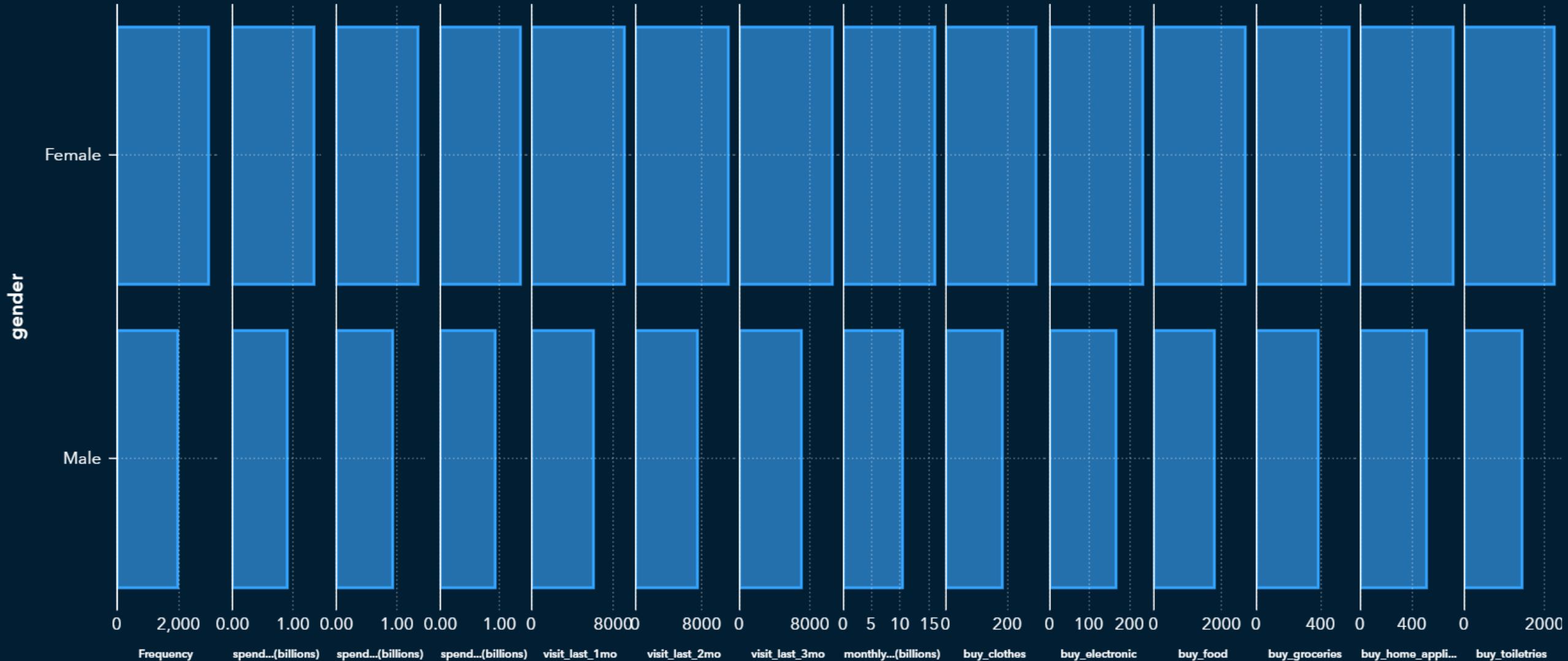
Frequency, spending_last_1mo, spending_last_2mo, spending_last_3mo, visit_last_1mo, visit_last_2mo, visit_last_3mo, monthly_income, buy_clothes, buy_electronic, buy_food, buy_groceries, buy_home_appliances, buy_toiletries by response

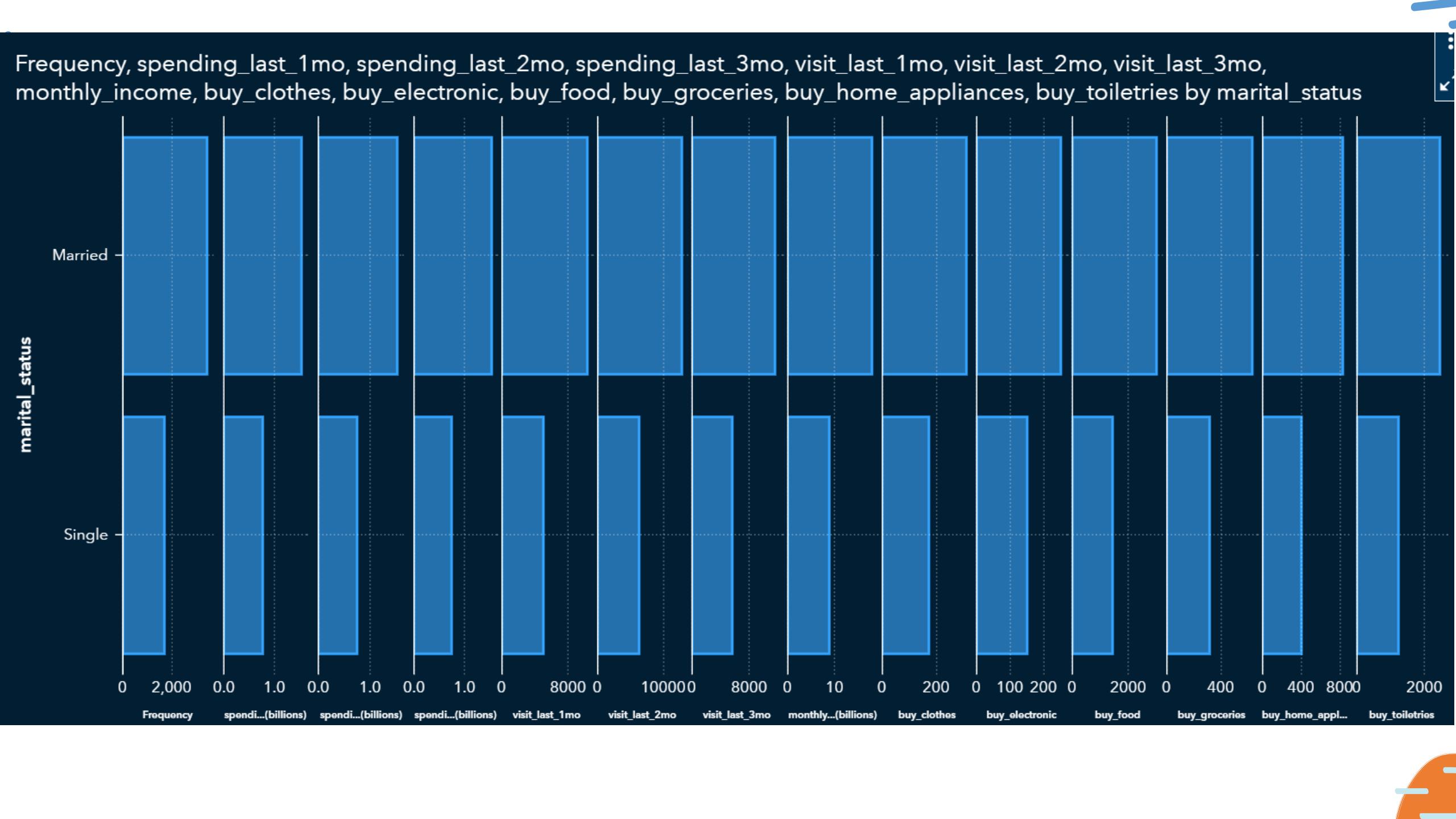


Perhatikan bahwa feature buy produk adalah biner 0 dan 1 yang artinya, total buy pada barplot tersebut menyatakan banyaknya customer yang melakukan buy , maka dari barplot pada bagian feature buy suatu produk, kita bisa mendapat insight bahwa justru lebih banyak customer-response-no yang membeli produk clothes, food, dan setiap produk di atas. Sehingga sepertinya feature ini kurang important terhadap response.

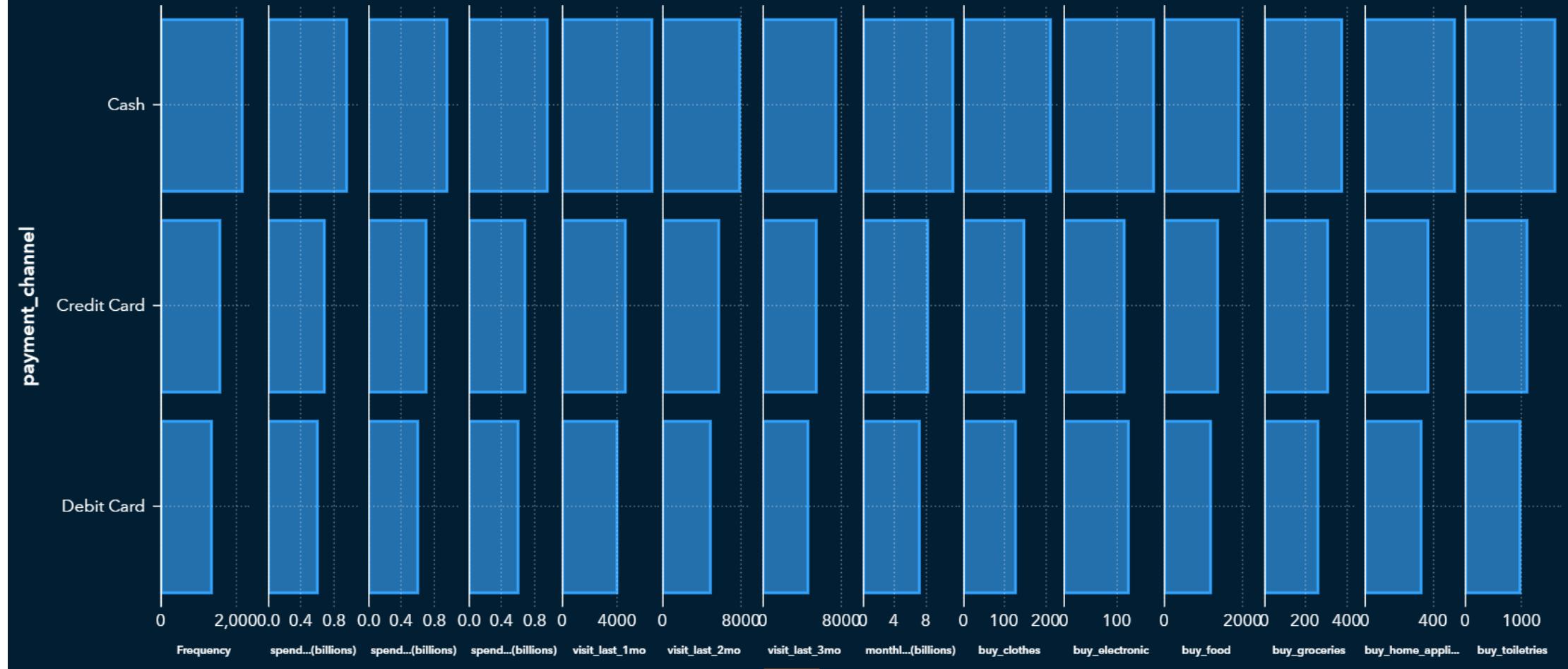
BAR PLOT

Frequency, spending_last_1mo, spending_last_2mo, spending_last_3mo, visit_last_1mo, visit_last_2mo, visit_last_3mo, monthly_income, buy_clothes, buy_electronic, buy_food, buy_groceries, buy_home_appliances, buy_toiletries by gender



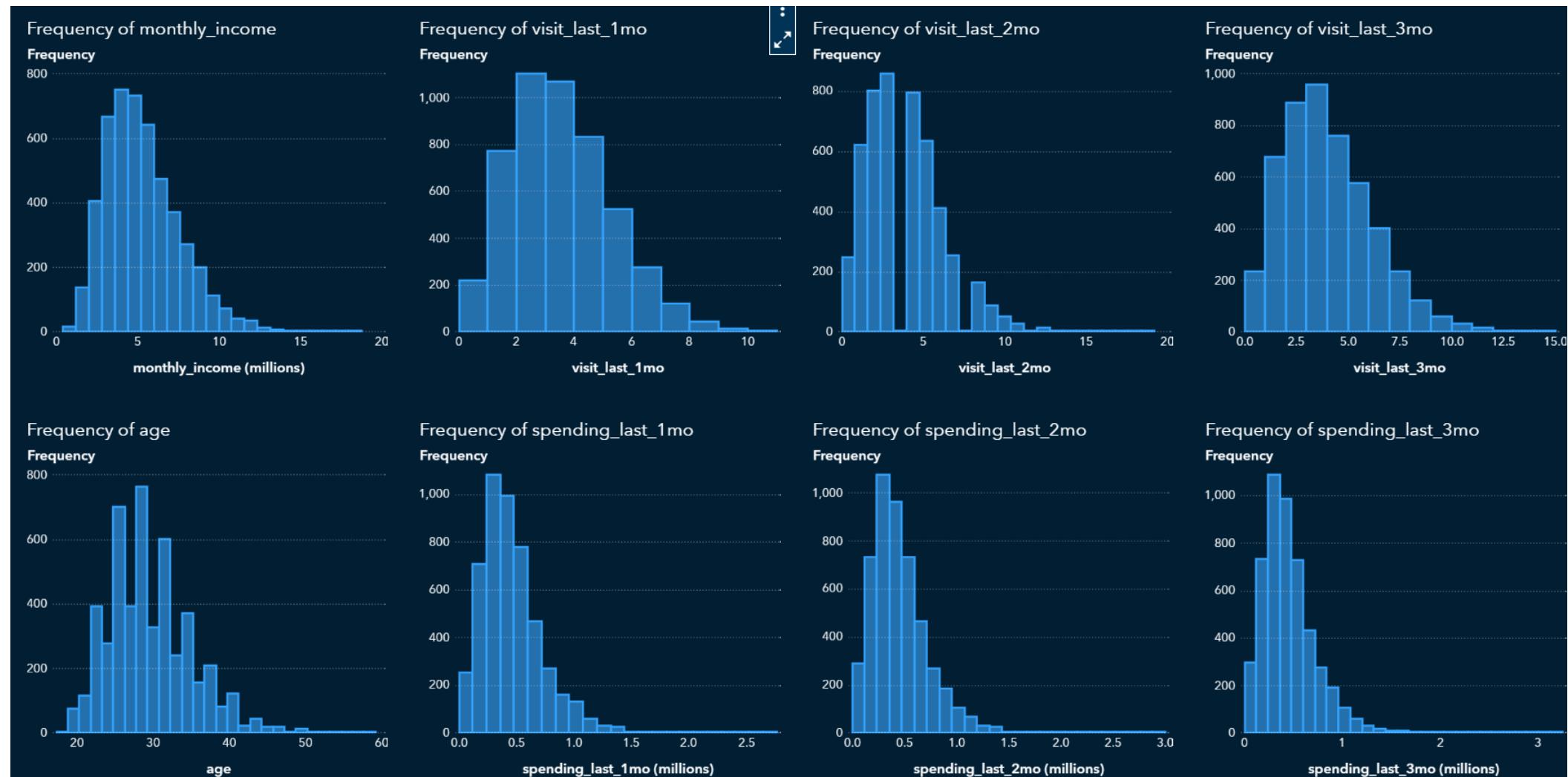


Frequency, spending_last_1mo, spending_last_2mo, spending_last_3mo, visit_last_1mo, visit_last_2mo, visit_last_3mo, monthly_income, buy_clothes, buy_electronic, buy_food, buy_groceries, buy_home_appliances, buy_toiletries by payment_channel



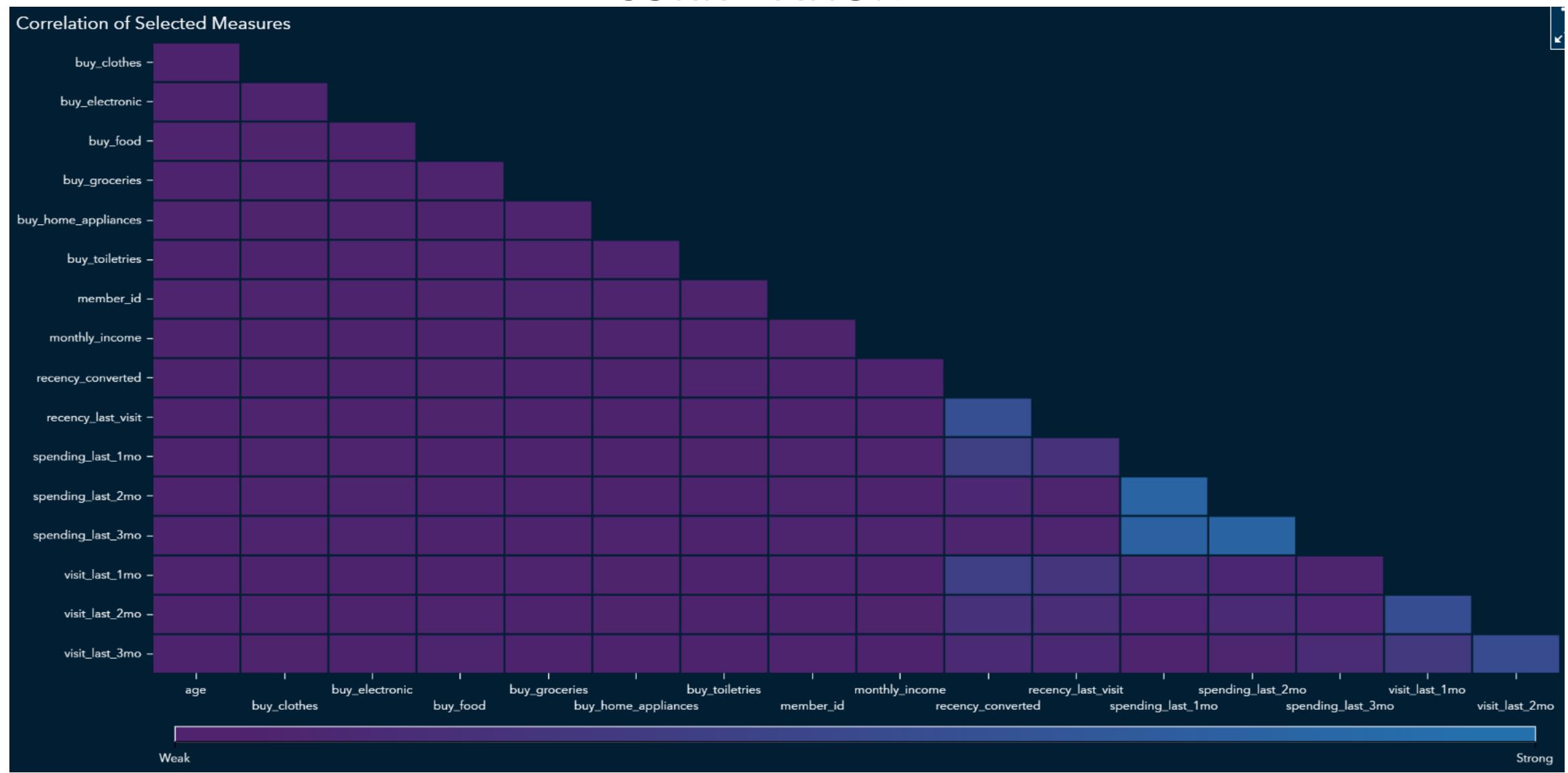
Baik gender, marital_status, maupun payment_channel, terdapat satu kelas pada masing-masing feature tersebut yang memiliki frekuensi yang lebih besar, bahkan untuk setiap nilai total feature lainnya. Contohnya Cash memiliki frekuensi, total spend, visit, income dst yang lebih besar dari kelas lainnya. Maka kita bisa mencoba untuk melakukan pemodelan dengan groupby feature kategorikal

HISTOGRAM



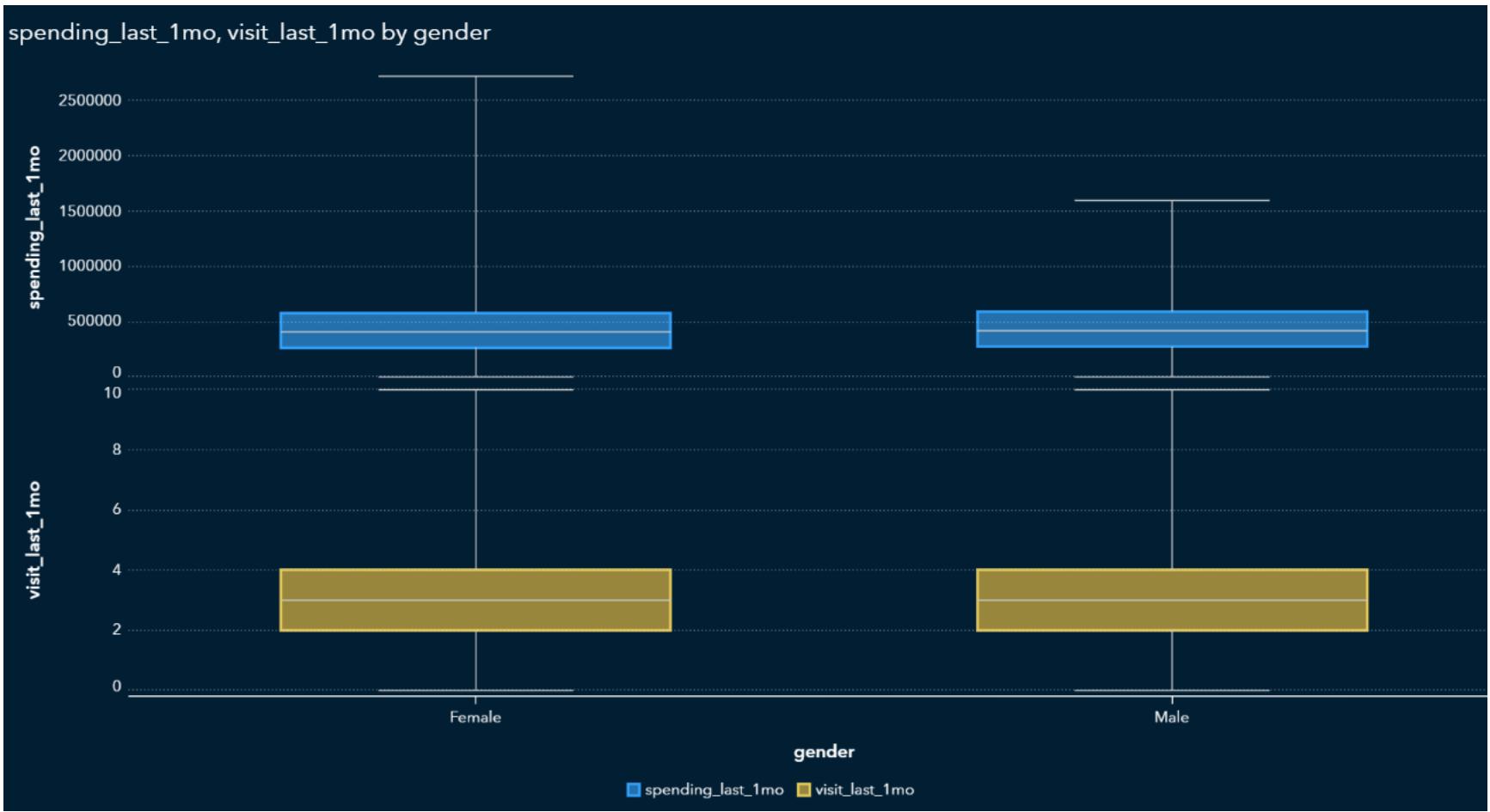
Dari pengamatan, distribusi dari numerical feature **tidak normal dan sedikit skew positif**. Jadi standardisasi yang dilakukan di awal telah tepat sehingga akan menjadikan distribusi dari feature menjadi cukup normal. Hal ini penting karena algoritma model memiliki asumsi bahwa distribusi data adalah normal.

CORRELATION



Spending last 1,2,3 month memiliki korelasi yang kuat, hal ini mengkonfirmasi adanya trend pada scatter plot sebelumnya. Hal ini juga dapat mengindikasikan bahwa jika salah satu feature tersebut memiliki variable importance yang tinggi, maka feature lainnya yg memiliki korelasi tinggi dengannya juga akan memiliki variable importance yang tinggi pula

BOX-PLOT



Untuk efisiensi waktu, tidak semua boxplot saya masukkan. Namun insight yang perlu ditekankan adalah baik dipisah berdasarkan kategorikal feature maupun tidak, data kita secara umum tidak terdapat pencilan, karena tidak ada observasi yang diluar upper maupun lower whisker .

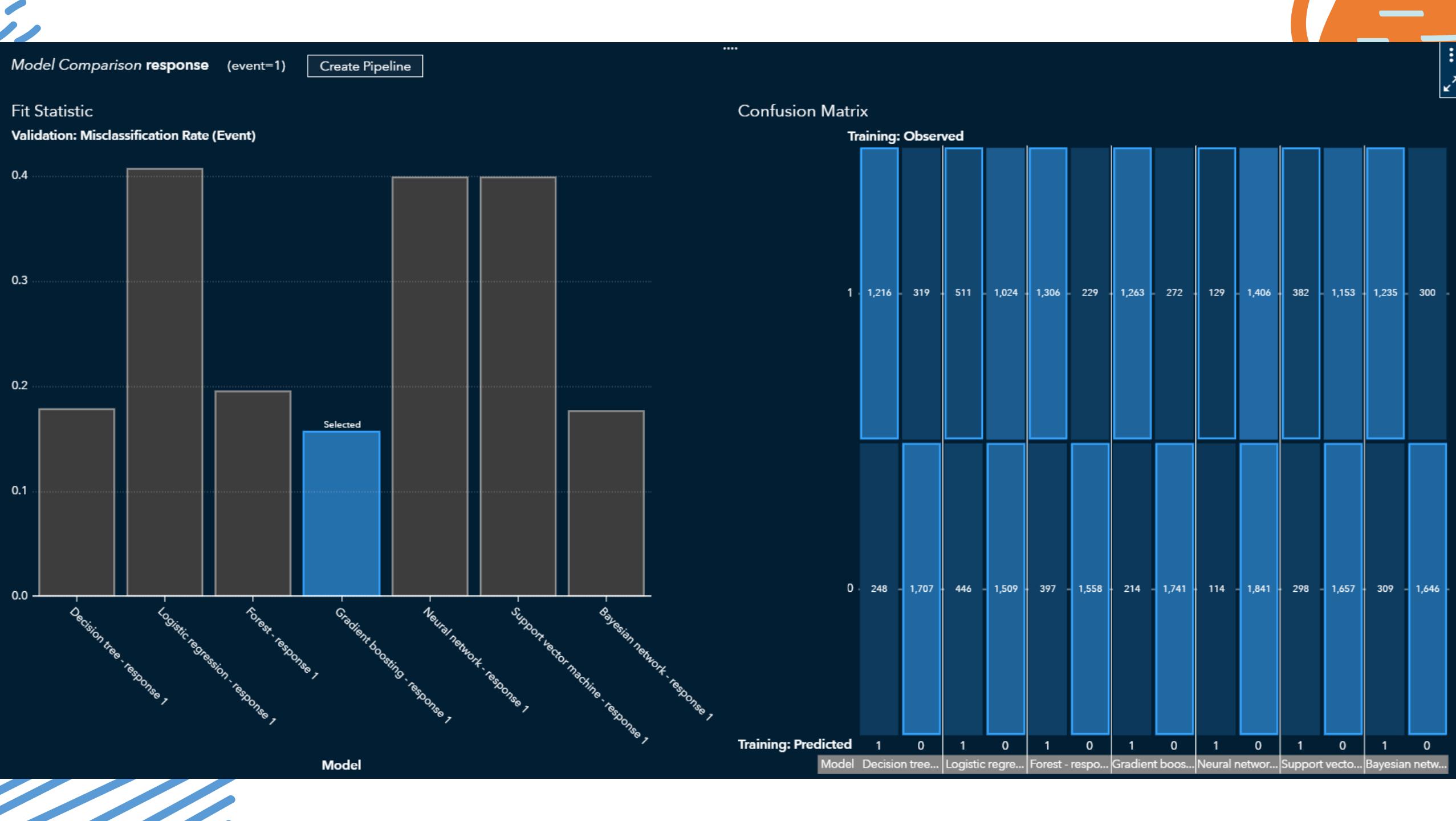
4. Modelling + 5. Ensemble Learning

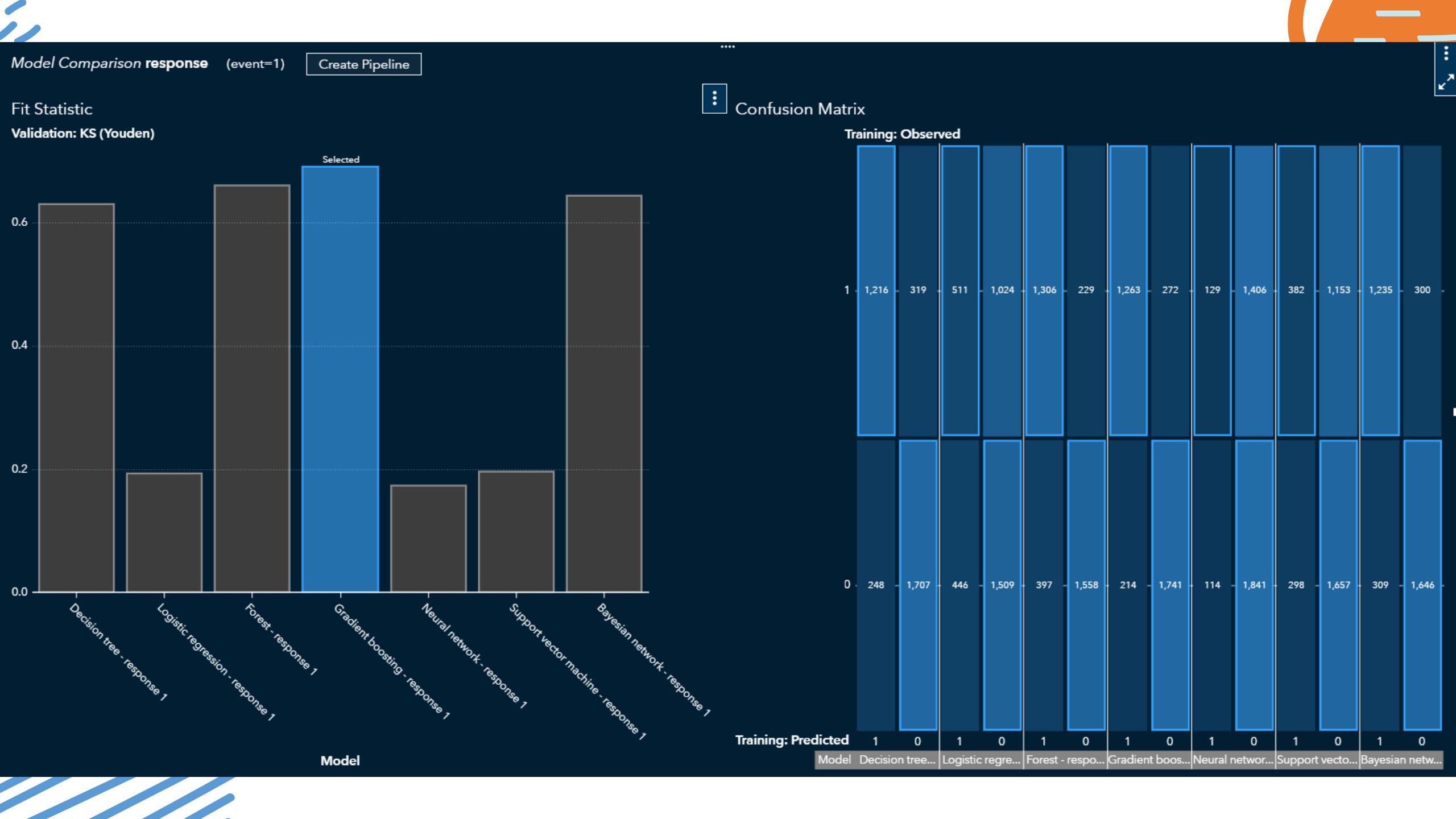
Ensemble Learning yang digunakan

- Bagging: Random Forest*
- Boosting: Gradient Boosting*



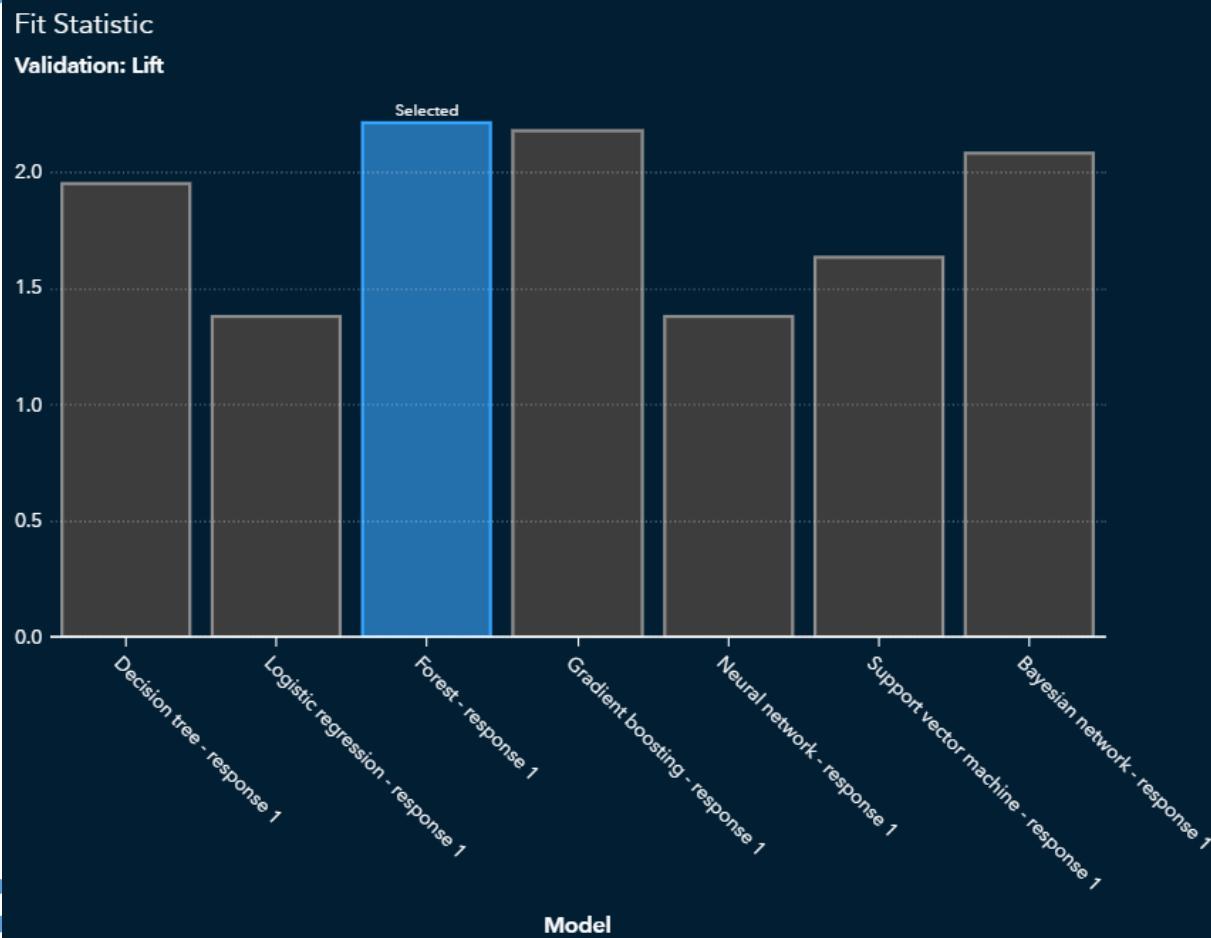
Model Comparison





Model Comparison response (event=1)

Create Pipeline



Confusion Matrix

Training: Observed

Model	Category 0	Category 1
Decision tree	248	1,216
Logistic regression	1,707	319
Forest - response	446	511
Gradient boosting	1,509	1,024
Neural network	397	1,306
Support vector machine	1,558	229
Bayesian network	214	1,263
	1,741	272
	114	129
	1,841	1,406
	298	382
	1,657	1,153
	309	1,235
	1,646	300

DATA WITH STANDARDIZATION

Selected	Model	KS (Youden)	Lift	Misclassification Rate...	Visualization Type	Number Of Observ...	Perce ntile	Predicti on cutoff	C Statistic	Cumulativ e % Captured	Cumulative % Events	Cumulative Lift	F1 Score	FDR	FPR	Gain	Gamma	Gini	Tau
Yes	Training:...	0.717	2.1759	0.139	Gradient Boosting	3,490	5	0.5	0.925	10.879	95.429	2.1759	0.839	0.145	0.109	1.176	0.856	0.851	0.419
Yes	Validation:...	0.690	2.1830	0.156	Gradient Boosting	1,496	5	0.5	0.890	10.915	90.667	2.1830	0.812	0.186	0.132	1.183	0.786	0.780	0.379
No	Training: ...	0.665	1.9108	0.162	Decision Tree	3,490	5	0.5	0.853	9.554	83.804	1.9108	0.811	0.169	0.127	0.911	0.835	0.705	0.348
No	Training:...	0.649	2.0586	0.174	Unknown	3,490	5	0.5	0.871	10.293	90.286	2.0586	0.802	0.200	0.158	1.059	0.751	0.741	0.365
No	Validation:...	0.644	2.0867	0.176	Unknown	1,496	5	0.5	0.864	10.433	86.667	2.0867	0.792	0.222	0.164	1.087	0.737	0.728	0.354
No	Validation:...	0.630	1.9541	0.178	Decision Tree	1,496	5	0.5	0.845	9.770	81.159	1.9541	0.784	0.209	0.147	0.954	0.827	0.690	0.336
No	Training: F...	0.667	2.1086	0.179	Forest	3,490	5	0.5	0.891	10.543	92.478	2.1086	0.807	0.233	0.203	1.109	0.804	0.782	0.385
No	Validation:...	0.660	2.2109	0.195	Forest	1,496	5	0.5	0.876	11.054	91.824	2.2109	0.783	0.270	0.223	1.211	0.774	0.752	0.366
No	Validation:...	0.173	1.3804	0.398	Neural Network	1,496	5	0.5	0.614	6.902	57.333	1.3804	0.179	0.369	0.044	0.380	0.248	0.229	0.111
No	Validation:...	0.196	1.6372	0.399	Support Vector M...	1,496	5	0.5	0.617	8.186	68.000	1.6372	0.342	0.454	0.148	0.637	0.239	0.235	0.114
No	Validation:...	0.193	1.3804	0.406	Logistic Regression	1,496	5	0.5	0.625	6.902	57.333	1.3804	0.405	0.481	0.220	0.380	0.258	0.250	0.121
No	Training:...	0.178	1.1987	0.416	Support Vector M...	3,490	5	0.5	0.619	5.993	52.571	1.1987	0.345	0.438	0.152	0.199	0.242	0.237	0.117
No	Training: L...	0.167	1.1987	0.421	Logistic Regression	3,490	5	0.5	0.615	5.993	52.571	1.1987	0.410	0.466	0.228	0.199	0.238	0.231	0.114
No	Training:...	0.165	1.2638	0.436	Neural Network	3,490	5	0.5	0.615	6.319	55.429	1.2638	0.145	0.469	0.058	0.264	0.248	0.230	0.113

DATA WITHOUT STANDARDIZATION

Selected	Model	KS (Youden)	Lift	Misclassification Rate (Events)	Visualization Type	Number Of Observ...	Perce ntile	Predicti on cutoff	C Statistic	Cumulative % Captured	Cumulative % Events	Cumulative Lift	F1 Score	FDR	FPR	Gain	Gamma	Gini	Tau
Yes	Trainin...	0.714	2.2062	0.142	Gradient Boosting	3,490	5	0.5	0.925	11.031	96.000	2.2062	0.835	0.152	0.114	1.206	0.854	0.849	0.418
Yes	Validat...	0.701	2.1417	0.151	Gradient Boosting	1,496	5	0.5	0.901	10.709	90.667	2.1417	0.823	0.182	0.136	1.142	0.808	0.803	0.393
No	Validat...	0.668	2.1314	0.161	Unknown	1,496	5	0.5	0.891	10.657	90.228	2.1314	0.809	0.184	0.134	1.131	0.799	0.783	0.383
No	Validat...	0.664	1.9921	0.162	Decision Tree	1,496	5	0.5	0.858	9.960	84.331	1.9921	0.805	0.177	0.125	0.992	0.838	0.716	0.350
No	Trainin...	0.652	1.9035	0.168	Decision Tree	3,490	5	0.5	0.847	9.518	82.830	1.9035	0.802	0.173	0.126	0.904	0.829	0.694	0.341
No	Trainin...	0.660	2.0465	0.189	Forest	3,490	5	0.5	0.889	10.233	89.054	2.0465	0.798	0.255	0.227	1.047	0.790	0.778	0.383
No	Validat...	0.689	2.1740	0.198	Forest	1,496	5	0.5	0.897	10.870	92.031	2.1740	0.787	0.276	0.243	1.174	0.805	0.795	0.389
No	Trainin...	0.584	2.0126	0.204	Unknown	3,490	5	0.5	0.862	10.063	87.577	2.0126	0.763	0.227	0.171	1.013	0.741	0.724	0.356
No	Trainin...	0.184	1.3526	0.407	Support Vector Machine	3,490	5	0.5	0.625	6.763	58.857	1.3526	0.359	0.425	0.149	0.353	0.254	0.249	0.123
No	Trainin...	0.179	1.3526	0.415	Logistic Regression	3,490	5	0.5	0.621	6.763	58.857	1.3526	0.419	0.461	0.228	0.353	0.250	0.243	0.120
No	Validat...	0.175	1.1339	0.416	Support Vector Machine	1,496	5	0.5	0.604	5.669	48.000	1.1339	0.341	0.481	0.173	0.134	0.212	0.208	0.102
No	Validat...	0.148	1.1024	0.420	Neural Network	1,496	5	0.5	0.601	5.512	46.667	1.1024	0.224	0.483	0.099	0.102	0.216	0.203	0.099
No	Trainin...	0.175	1.3263	0.423	Neural Network	3,490	5	0.5	0.623	6.632	57.714	1.3263	0.225	0.440	0.086	0.326	0.261	0.246	0.121
No	Validat...	0.157	1.1969	0.424	Logistic Regression	1,496	5	0.5	0.601	5.984	50.667	1.1969	0.399	0.500	0.245	0.197	0.208	0.202	0.099

6. Final Evaluation

HASIL PERBANDINGAN

Diperoleh hasil bahwa data dengan standardisasi memperoleh performansi lebih baik dibandingkan dengan data tanpa standardisasi, dengan hasil sebagai berikut

Urutan berdasarkan validation misclassification rate (dari terendah)

1. Gradient Boosting 0.156
2. Bayesian Network 0.174

Urutan berdasarkan validation KS (Youden)/ROC (dari tertinggi)

1. Gradient Boosting 0.696
2. Random Forest 0.672

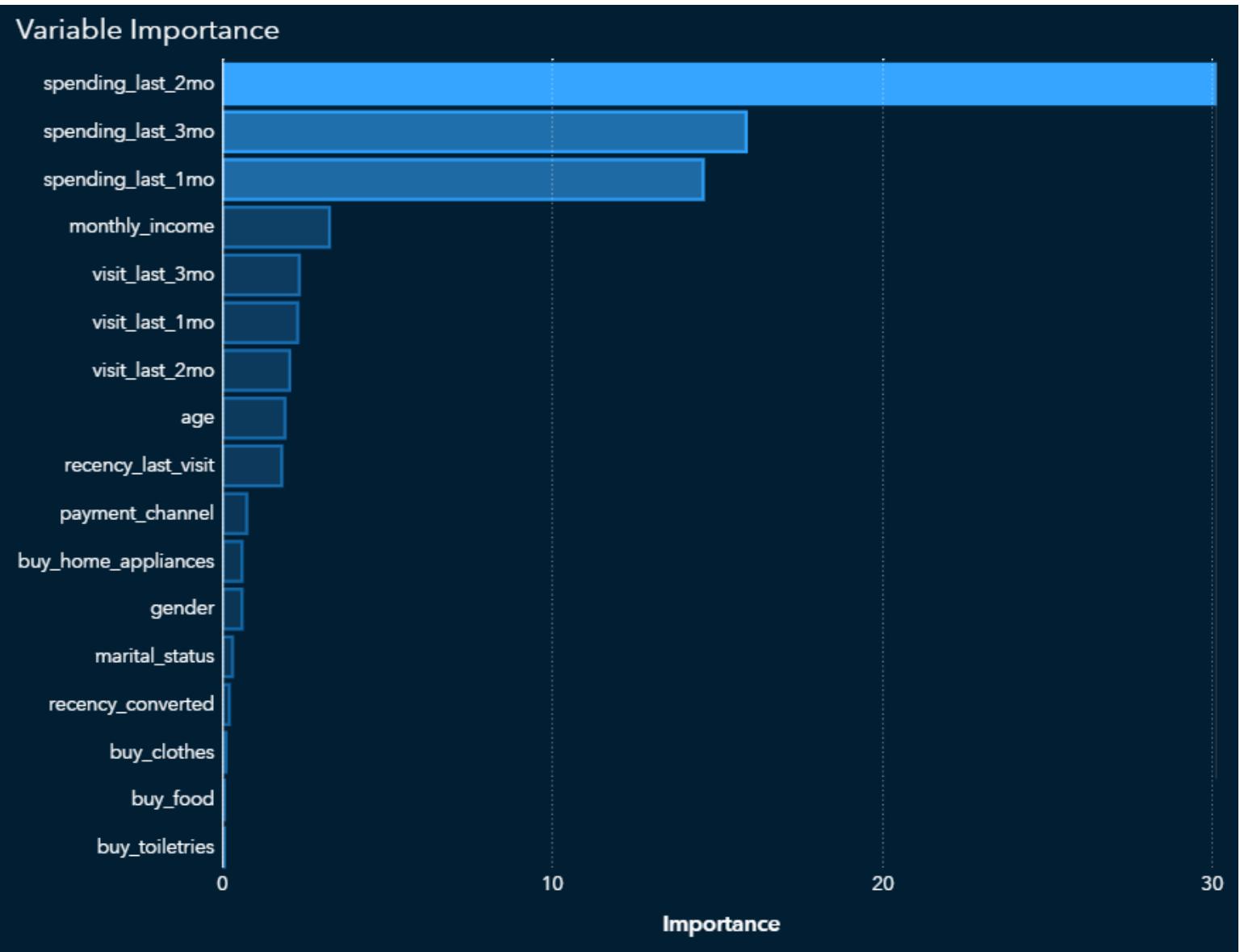
Urutan berdasarkan validation lift (dari tertinggi)

1. Random Forest 2.1109
2. Gradient Boosting 2.0955

Feature yang berpengaruh

Berdasarkan 3 model yang terbaik yang telah diperoleh sebelumnya

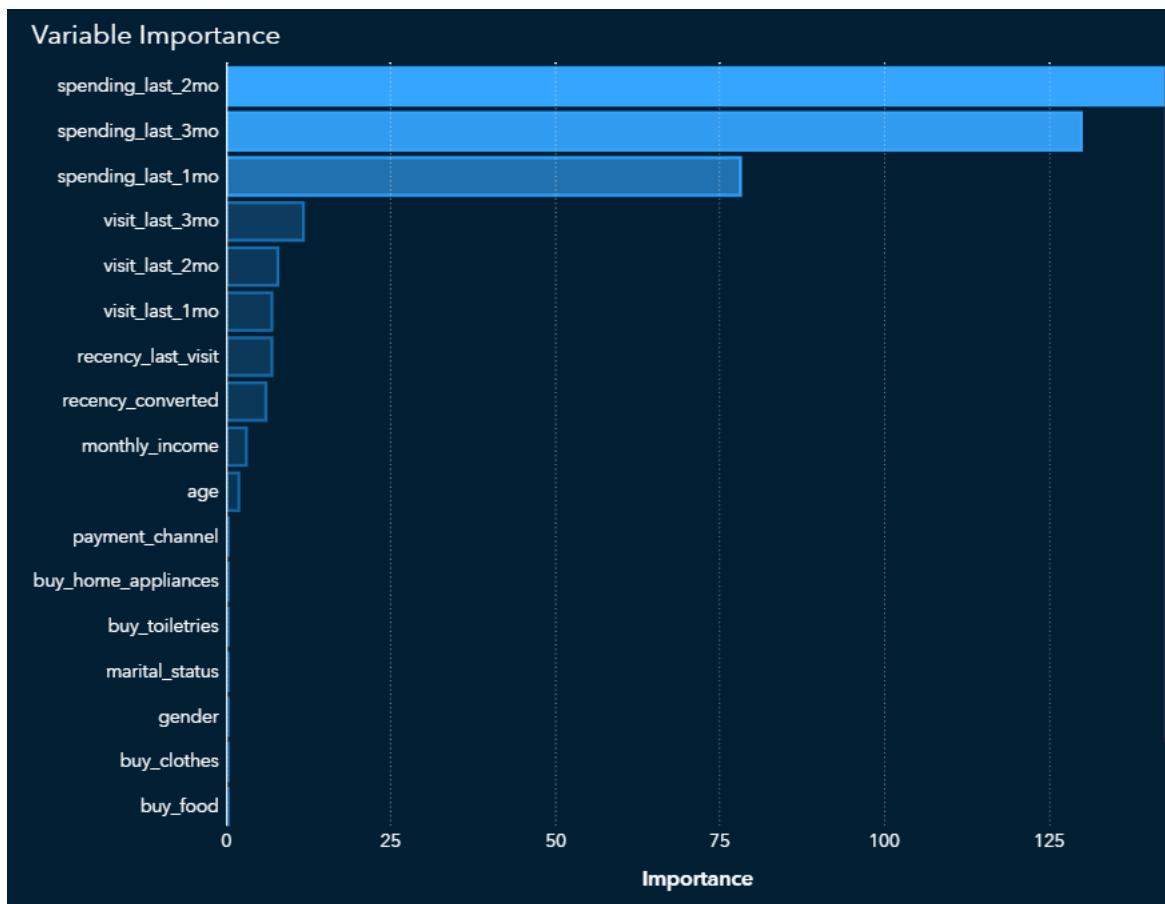
GRADIENT BOOSTING



GRADIENT BOOSTING

Variable	Importance	Standard Deviation
spending_last_2mo	30.0598	32.3957
spending_last_3mo	15.8926	20.5770
spending_last_1mo	14.5848	15.7909
monthly_income	3.2631	1.0975
visit_last_3mo	2.3181	1.9672
visit_last_1mo	2.2799	1.5432
visit_last_2mo	2.0241	2.1607
age	1.9186	0.9524
recency_last_visit	1.8144	0.8563
payment_channel	0.7656	0.9883
buy_home_appliances	0.6171	0.8355
gender	0.5819	0.7852
marital_status	0.3171	1.6149
recencyConverted	0.2077	0.0000
buy_clothes	0.1266	0.7492
buy_food	0.0824	0.3424
buy_toiletries	0.0717	0.5515
buy_groceries	0.0642	0.2191
buy_electronic	0.0162	0.0000

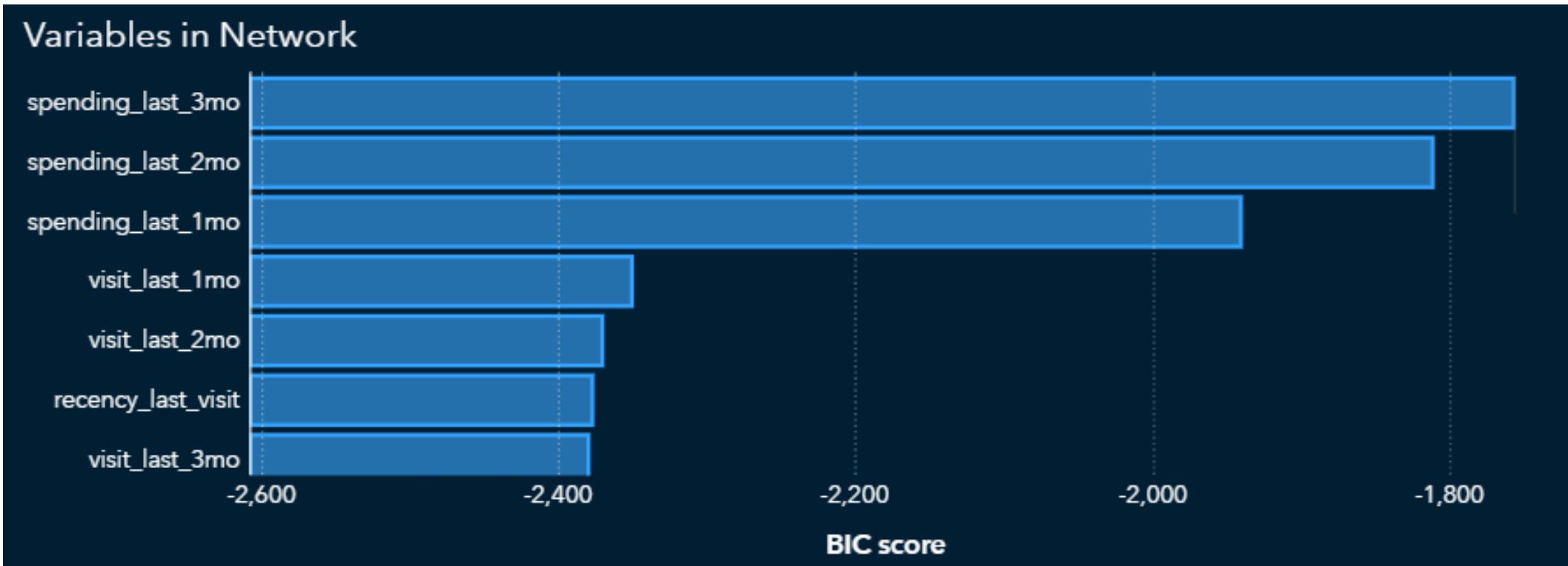
RANDOM FOREST



RANDOM FOREST

Variable	Importance	Standard Deviation
spending_last_2mo	142.3150	74.1511
spending_last_3mo	129.8972	75.4678
spending_last_1mo	78.0636	60.0479
visit_last_3mo	11.7781	14.3169
visit_last_2mo	7.9524	11.0073
visit_last_1mo	7.0251	9.4369
recency_last_visit	6.8535	8.9703
recencyConverted	5.9465	12.5150
monthly_income	3.1356	1.4613
age	1.9958	1.4178
payment_channel	0.4151	0.8754
buy_home_appliances	0.3527	1.3250
buy_toiletries	0.3405	1.2417
marital_status	0.2444	1.1094
gender	0.1735	1.1447
buy_clothes	0.1656	0.7666
buy_food	0.1611	0.7427
buy_electronic	0.1232	0.5487
buy_groceries	0.0740	0.7188

BAYESIAN NETWORK



BAYESIAN NETWORK

Variable	Selected	Chi-Square	Pr > Chi-Square ▲	G-Square	Pr > G-Square	Degrees of Freedom	Mutual Information	Conditional Variable
spending_last_3mo	Yes	1,212.3862	<0.00001	1,320.7739	<0.00001	7	0.5613	
spending_last_2mo	Yes	1,072.8396	<0.00001	1,212.2129	<0.00001	7	0.5417	
spending_last_1mo	Yes	839.2428	<0.00001	956.3465	<0.00001	7	0.4896	
visit_last_1mo	Yes	137.6315	<0.00001	152.7537	<0.00001	9	0.2069	
visit_last_2mo	Yes	110.2149	<0.00001	114.5064	<0.00001	9	0.1797	
visit_last_3mo	Yes	93.8316	<0.00001	95.9640	<0.00001	9	0.1647	
recency_last_visit	Yes	82.1227	<0.00001	100.8591	<0.00001	9	0.1688	
gender	No	1.3445	0.24623	1.3437	0.24637	1	0.0196	
marital_status	No	0.3458	0.55649	0.3456	0.55662	1	0.0100	
age	No	6.7231	0.66592	6.7277	0.66545	9	0.0439	
payment_channel	No	0.4963	0.78025	0.4960	0.78038	2	0.0119	
monthly_income	No	3.5336	0.93935	3.5215	0.94000	9	0.0318	
recencyConverted	No	0.0000	0.99984	0.0000	0.99988	1	0.0000	
buy_toiletries	No	0.0000	0.99998	0.0000	0.99998	1	0.0000	
buy_home_appliances	No	0.0000	0.99998	0.0000	0.99998	1	0.0000	
buy_electronic	No	0.0000	0.99998	0.0000	0.99998	1	0.0000	
buy_groceries	No	0.0000	0.99998	0.0000	0.99998	1	0.0000	
buy_clothes	No	0.0000	0.99998	0.0000	0.99998	1	0.0000	
buy_food	No	0.0000	0.99999	0.0000	0.99999	1	0.0000	

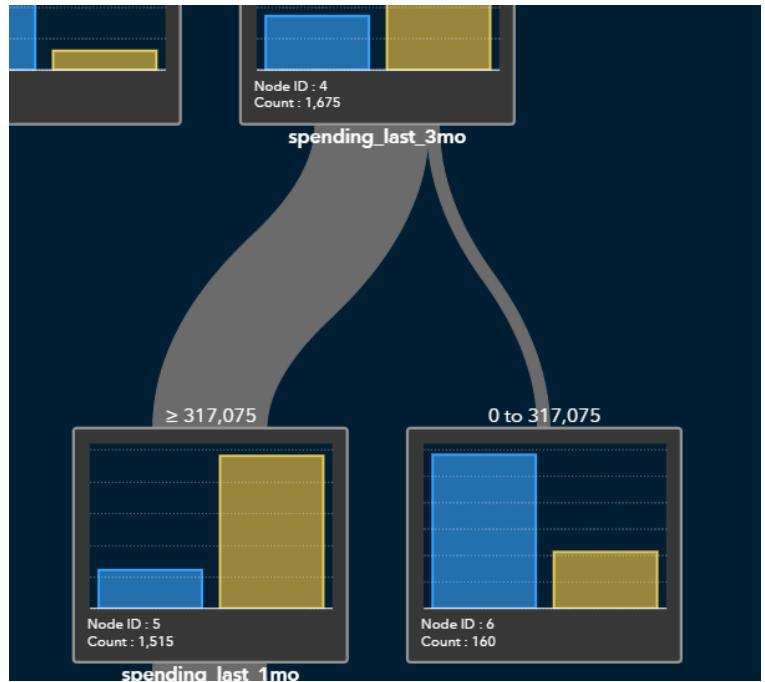
KESIMPULAN

Berdasarkan 3 model yang terbaik yang telah diperoleh sebelumnya

1. Informasi dan insight dari data dapat dilihat dalam EDA
2. Feature yang paling penting secara berurutan kurang lebih adalah spending_last baik bulan ke-3, ke-2, maupun ke-1, kemudian visit_last baik bulan ke-3, ke-2, maupun ke-1, serta recency last visit maka strategi untuk memaksimumkan jumlah customer yang aktif adalah dengan fokus pada sedikit feature tersebut secara berurutan. Hal ini sejalan dengan prinsip Pareto 80/20 yang mengungkapkan bahwa 20% input mempengaruhi 80% output, jadi perusahaan bisa melakukan efisiensi dengan fokus pada sedikit feature namun memperoleh hasil yang lebih besar.

Untuk menentukan member mana yang mengambil promo, jika ada data member maka bisa dilakukan prediksi. Namun secara umum, dari decision tree bisa dilihat bahwa kelompok member yang mengambil promo adalah member yang memiliki feature-feature pada no.2 yang tinggi (batas tinggi dapat mempertimbangkan decision tree pada slide berikutnya)

3. Sebaiknya manajer SAS Mart berkolaborasi dengan tim marketing untuk melakukan promosi seperti email atau whatsapp kepada member yang memiliki nilai tinggi feature pada no.2





Departem
IPB University



Terimakasih!

“Data are just summaries of thousands of stories - tell a few of those stories to make the data meaningful.”

-Chip and Dan Heath