

✓ Approach 1: The AI Architect

My Mission: is to use BigQuery's built-in generative AI to architect intelligent business applications that generate creative content, summarize complex information, and forecast future trends directly within the data warehouse — for example, by building a hyper-personalized marketing engine that creates unique emails for every customer based on their individual purchase history and preferences


[Open in Colab](#)

Google Cloud Colab Enterprise logo
[Open in Colab Enterprise](#)


[Open in Vertex AI Workbench](#)

[Open in BigQuery Studio](#)

[View on GitHub](#)

Share to:



My Mission: Use BigQuery's built-in generative AI to architect intelligent business applications and workflows. Build solutions that can generate creative content, summarize complex information, or even forecast future trends directly within your data warehouse.

Your Toolbox (Must use at least one):

- **Generative AI in SQL:**
 - `ML.GENERATE_TEXT` : The classic function for large-scale text generation.
 - `AI.GENERATE` : Generate free-form text or structured data based on a schema from a prompt.
 - `AI.GENERATE_BOOL` : Get a simple True/False answer about your data.
 - `AI.GENERATE_DOUBLE` : Extract a specific decimal number from text.
 - `AI.GENERATE_INT` : Extract a specific whole number from text.
 - `AI.GENERATE_TABLE` : Create a structured table of data from a single prompt.
 - `AI.FORECAST` : Predict future values for time-series data with a single function.
- **Generative AI in BigFrames (Python):**
 - `bigframes.ml.llm.GeminiTextGenerator` : Leverage the power of Gemini models in your Python workflows.
 - `bigframes.DataFrame.ai.forecast()` : Run powerful forecasting models directly on your DataFrames.

Inspiration:

- Build a **Hyper-Personalized Marketing Engine**: Generate unique marketing emails for every customer based on their individual purchase history and preferences.
- Create an **Executive "Insight" Dashboard**: Develop a dashboard that automatically ingests raw support call logs and transforms them into summarized, categorized, and actionable business insights.

✓ Before you begin

✓ Set up your Google Cloud project

The following steps are required, regardless of your notebook environment.

1. [Select or create a Google Cloud project](#). When you first create an account, you get a \$300 free credit towards your compute/storage costs.
2. [Make sure that billing is enabled for your project](#).
3. [Enable the BigQuery, BigQuery Connection, and Vertex AI APIs](#).
4. If you are running this notebook locally, you need to install the [Cloud SDK](#).

```
1 ! gcloud auth login
```



You are running on a Google Compute Engine virtual machine.
 It is recommended that you use service accounts for authentication.

You can run:

```
$ gcloud config set account `ACCOUNT`
```

to switch accounts if necessary.

Your credentials may be visible to others with access to this virtual machine. Are you sure you want to authenticate with your personal account?

Do you want to continue (Y/n)? Y

Go to the following link in your browser, and complete the sign-in prompts:

https://accounts.google.com/o/oauth2/auth?response_type=code&client_id=32555940559.apps.googleusercontent.com&


Once finished, enter the verification code provided in your browser: 4/0AVMBsJhz3IofhcPSvEY-LzmoZd_wMyfoG_mz_M-Wt-

You are now logged in as [idriss.choudjem@gmail.com].

Your current project is [gothic-list-469006-v9]. You can change this setting by running:

```
$ gcloud config set project PROJECT_ID
```

```
1 PROJECT_ID = "gothic-list-469006-v9" # @param {type:"
2
3 # Set the project id
4 ! gcloud config set project {PROJECT_ID}
5 # Dans un bloc Colab (préfixe !), remplace par ton PRO
6 ! gcloud services enable bigquery.googleapis.com bigqu
7
```

PROJECT_ID: "gothic-list-469006-v9" 



Updated property [core/project].

Operation "operations/acat.p2-670703985696-f21c0221-e0c5-4040-aff5-ceb52dbd2d45" finished successfully.

✓ Authenticate to your Google Cloud account

Depending on your Jupyter environment, you may have to manually authenticate. Follow the relevant instructions below.

1. Colab Enterprise or BigQuery Studio Notebooks

- Do nothing as you are already authenticated.

2. Colab, uncomment and run:

```
1 from google.colab import auth
2
3 auth.authenticate_user()
```

3. Créer le dataset pour le modèle:

```
1 %%bigquery marketing_model_dataset --project {PROJECT_ID}
2 CREATE SCHEMA IF NOT EXISTS `marketing_email_model_dataset`
3 OPTIONS (location="US");
```



Job ID 7299d1c9-e06c-4f0a-87f7-ac5fc0250780 successfully executed: 100%

✓ Create BigQuery Cloud resource connection

You will need to create a [Cloud resource connection](#) to enable BigQuery to interact with Vertex AI services.

```
1 !bq mk --connection --location=us --project_id=gothic-list-469006-v9 \
2 --connection_type=CLOUD_RESOURCE gothic-list-469006-v9
```



BigQuery error in mk operation: Already Exists: Connection projects/670703985696/locations/us/connections/gothic-list-469006-v9

▼ Set permissions for Service Account

The resource connection service account requires certain project-level permissions to interact with Vertex AI.

```
1 SERVICE_ACCT = !bq show --format=prettyjson --connection us.gothic-list-469006-v9 | grep bqcx-670703985696-at17 |
2 SERVICE_ACCT_EMAIL = SERVICE_ACCT[-1]
3 print(SERVICE_ACCT_EMAIL)
```

↗ bqcx-670703985696-at17@gcp-sa-bigquery-condel.iam.gserviceaccount.com

```
1 import time
2
3 !gcloud projects add-iam-policy-binding --format=none $PROJECT_ID --member=serviceAccount:$SERVICE_ACCT_EMAIL --ro
4 !gcloud projects add-iam-policy-binding --format=none $PROJECT_ID --member=serviceAccount:$SERVICE_ACCT_EMAIL --ro
5
6 # wait 60 seconds, give IAM updates time to propagate, otherwise, following cells will fail
7 time.sleep(60)
```

↗ Updated IAM policy for project [gothic-list-469006-v9].
Updated IAM policy for project [gothic-list-469006-v9].

▼ I - Build a Hyper-Personalized Marketing Engine

Generate text and structured data with `AI.GENERATE` and `AI.GENERATE_TABLE`

These functions allow you to leverage the power of large language models (LLMs) directly within BigQuery to generate new text content, including creating content with a specified schema. Both functions work by sending requests to your choice of a [generally available](#) or [preview](#) Gemini model, and then returning that model's response.

▼ Using `ML.GENERATE_TEXT` : for Hyper-Personalized Marketing email

Let's use the `ML.GENERATE_TEXT` function to generate unique marketing emails for every customer based on their individual purchase history and preferences from the `bigquery-public-data.thelook_ecommerce` datasets.

In today's highly competitive e-commerce landscape, customer retention and campaign optimization are critical success factors. Traditional marketing approaches often fall short because they fail to address the unique preferences and purchase histories of individual customers. By leveraging BigQuery's built-in generative AI capabilities, businesses can move beyond generic campaigns and instead generate hyper-personalized marketing emails tailored to each customer. This approach enables the creation of intelligent workflows and business applications that not only deliver engaging content, but also summarize insights, predict future trends, and drive stronger customer relationships. The result is a Hyper-Personalized Marketing Engine designed to maximize engagement, loyalty, and long-term growth.

▼ 🔑 Building of a Parametric User Table for Behavior Analysis

the goal here is creating a **parametric table** that consolidates key user-level indicators from the public [The Look eCommerce](#) dataset. It builds an enriched customer base by calculating `total orders`, `average basket`, `total revenue`, `recency of last purchase`, and `purchase frequency (monthly and quarterly)`. It also tracks product preferences by counting items purchased, identifying the top three products per user, and determining the top three products per country. Additionally, it captures each user's most recent event. Finally, the table aggregates these insights to provide a comprehensive view of each customer, combining personal details, purchasing behavior, product preferences, last interactions, and `market-level trends`.

```


1 %%bigquery df_parametric --project {PROJECT_ID}
2
3 CREATE OR REPLACE TABLE `gothic-list-469006-v9.ETL.parametric` AS
4 WITH base_with_freq AS (
5   SELECT
6     u.id AS user_id,
7     u.first_name,
8     u.email,
9     u.country,
10    COUNT(o.order_id) AS total_orders,
11    AVG(oi.sale_price) AS avg_basket,
12    SUM(oi.sale_price) AS total_revenue,
13    DATE_DIFF(CURRENT_DATE(), DATE(MAX(o.created_at)), MONTH) AS months_since_last_order,
14    COUNT(o.order_id) / NULLIF(DATE_DIFF(CURRENT_DATE(), DATE(MIN(o.created_at)), MONTH), 0) AS freq_orders_per_m
15    COUNT(o.order_id) / NULLIF(DATE_DIFF(CURRENT_DATE(), DATE(MIN(o.created_at)), MONTH) / 3, 0) AS freq_orders_p
16  FROM `bigquery-public-data.thelook_ecommerce.users` u
17  LEFT JOIN `bigquery-public-data.thelook_ecommerce.orders` o
18    ON u.id = o.user_id
19  LEFT JOIN `bigquery-public-data.thelook_ecommerce.order_items` oi
20    ON o.order_id = oi.order_id
21  GROUP BY u.id, u.first_name, u.email, u.country
22 ),
23
24 # -- Comptage des produits par utilisateur
25 product_counts AS (
26   SELECT
27     u.id AS user_id,
28     p.name AS product_name,
29     COUNT(*) AS product_count
30  FROM `bigquery-public-data.thelook_ecommerce.users` u
31  JOIN `bigquery-public-data.thelook_ecommerce.orders` o
32    ON u.id = o.user_id
33  JOIN `bigquery-public-data.thelook_ecommerce.order_items` oi
34    ON o.order_id = oi.order_id
35  JOIN `bigquery-public-data.thelook_ecommerce.products` p
36    ON oi.product_id = p.id
37  WHERE p.name IS NOT NULL
38  GROUP BY u.id, p.name
39 ),
40
41 # -- Top 3 produits par utilisateur
42 top_products_by_user AS (
43   SELECT
44     user_id,
45     ARRAY_AGG(product_name ORDER BY product_count DESC LIMIT 3) AS top_products_user
46  FROM product_counts
47  GROUP BY user_id
48 ),
49
50 # -- Dernier évènement par utilisateur
51 latest_event AS (
52   SELECT
53     user_id,
54     MAX(event_type) AS last_event
55  FROM `bigquery-public-data.thelook_ecommerce.events`
56  GROUP BY user_id
57 ),
58
59 # -- Top 3 produits par pays
60 top_products_by_country AS (
61   SELECT
62     country,
63     ARRAY_AGG(product_name ORDER BY product_count DESC LIMIT 3) AS top_products_country
64  FROM (
65    SELECT
66      u.country,
67      p.name AS product_name,
68      COUNT(*) AS product_count
69    FROM `bigquery-public-data.thelook_ecommerce.orders` o
70    JOIN `bigquery-public-data.thelook_ecommerce.order_items` oi
71      ON o.order_id = oi.order_id
72    JOIN `bigquery-public-data.thelook_ecommerce.products` p



```




```

73     ON oi.product_id = p.id
74     JOIN `bigquery-public-data.thelook_ecommerce.users` u
75     ON o.user_id = u.id
76     WHERE p.name IS NOT NULL
77     GROUP BY u.country, p.name
78 )
79 GROUP BY country
80 )
81
82 # -- Final : joindre les infos utilisateur avec leurs top produits, évènements et pays
83 SELECT
84     b.user_id,
85     b.first_name,
86     b.email,
87     b.country,
88     IFNULL(b.total_orders, 0) AS total_orders,
89     IFNULL(b.avg_basket, 0) AS avg_basket,
90     IFNULL(b.total_revenue, 0) AS total_revenue,
91     t.top_products_user,
92     l.last_event,
93     c.top_products_country,
94     IFNULL(b.months_since_last_order, 0) AS months_since_last_order,
95     IFNULL(b.freq_orders_per_month, 0) AS freq_orders_per_month,
96     IFNULL(b.freq_orders_per_quarter, 0) AS freq_orders_per_quarter
97 FROM base_with_freq b
98 LEFT JOIN top_products_by_user t
99     ON b.user_id = t.user_id
100 LEFT JOIN latest_event l
101     ON b.user_id = l.user_id
102 LEFT JOIN top_products_by_country c
103     ON b.country = c.country
104 ORDER BY b.user_id
105 # -- LIMIT 10000;
106

```

 Job ID 9eb4f049-3386-4ff2-81b0-ba384ea6f974 successfully executed: 100%

- 1 #  Dataset exported to CSV and Parquet in /content/ETL
- 2 #  Table saved in BigQuery at gothic-list-469006-v9.ETL.parametric

  Dataset exported to CSV and Parquet in /content/ETL
 Table saved in BigQuery at gothic-list-469006-v9.ETL.parametric

Customer Segmentation with Parametric Features using KMeans and

ML.GENERATE_TEXT

The underlying idea of using K-Means combined with "ML.GENERATE_TEXT" is to group data into affinity clusters based on their parametric features, and then reduce the volume of data that needs labeling. By leveraging the centroid, labeling can first be done manually or automatically, before generalizing the centroid labels to the entire cluster related to. This provides a first level of customer profiling by group, before diving into a more granularity level for hyper-personalized marketing

Here's a complete **Hand off** approach for customer segmentation on an e-commerce dataset by loading user-level parametric features from BigQuery, handling missing data and scaling features, reducing dimensionality via PCA, evaluating multiple KMeans clusterings using inertia, silhouette, Calinski-Harabasz, and Davies-Bouldin metrics to identify the optimal number of clusters "k" with a combined score, assigning users to segments and calculating centroids, saving the segmented dataset back to BigQuery, and generating diagnostic plots to visualize cluster evaluation metrics, with the business goal of providing actionable customer segments for targeted marketing, and personalized recommendations.

```

1 import numpy as np
2 import pandas as pd
3 from google.cloud import bigquery

```

```
4 from sklearn.cluster import KMeans
5 from sklearn.metrics import silhouette_score, calinski_harabasz_score, davies_bouldin_score
6 from sklearn.preprocessing import StandardScaler, MinMaxScaler
7 from sklearn.impute import SimpleImputer
8 from sklearn.decomposition import PCA
9 from umap import UMAP
10 import plotly.graph_objects as go
11 from plotly.subplots import make_subplots
12
13
14 # --Connexion BigQuery
15 client = bigquery.Client(project=PROJECT_ID)
16
17
18 # --Charger les données depuis BigQuery
19 query = f"""
20     SELECT *
21 FROM `{PROJECT_ID}.ETL.parametric`
22 """
23 df_segmented = client.query(query).to_dataframe()
24
25
26 # --Features pour clustering
27 features = [
28     "total_orders",
29     "avg_basket",
30     "total_revenue",
31     "months_since_last_order",
32     "freq_orders_per_month",
33     "freq_orders_per_quarter"
34 ]
35 X = df_segmented[features]
36
37
38 # --Imputation des valeurs manquantes + scaling
39 imputer = SimpleImputer(strategy="mean")
40 X_imputed = imputer.fit_transform(X)
41 scaler = StandardScaler()
42 X_scaled = scaler.fit_transform(X_imputed)
43
44
45 # --Réduction dimensionnelle PCA
46 pca = PCA(n_components=min(5, X_scaled.shape[1]), random_state=42)
47 X_reduced = pca.fit_transform(X_scaled)
48
49
50 # --Boucle KMeans pour trouver le meilleur k
51 k_min, k_max = 2, 22
52 k_values = list(range(k_min, k_max+1))
53
54 inertias, silhouettes, calinski, davies = [], [], [], []
55 sample_size = min(1000, X_reduced.shape[0]) # échantillon silhouette si dataset trop grand
56
57 for k in k_values:
58     kmeans = KMeans(n_clusters=k, random_state=42, n_init=10, algorithm="elkan")
59     labels = kmeans.fit_predict(X_reduced)
60
61     inertias.append(kmeans.inertia_)
62     silhouettes.append(silhouette_score(X_reduced, labels, sample_size=sample_size))
63     calinski.append(calinski_harabasz_score(X_reduced, labels))
64     davies.append(davies_bouldin_score(X_reduced, labels))
65
66
67 # --Normalisation & score combiné
68 scaler_norm = MinMaxScaler()
69 inertia_norm = 1 - scaler_norm.fit_transform(np.array(inertias).reshape(-1,1)).flatten()
70 silhouette_norm = scaler_norm.fit_transform(np.array(silhouettes).reshape(-1,1)).flatten()
71 calinski_norm = scaler_norm.fit_transform(np.array(calinski).reshape(-1,1)).flatten()
72 davies_norm = 1 - scaler_norm.fit_transform(np.array(davies).reshape(-1,1)).flatten()
73
74 combined_score = (inertia_norm + silhouette_norm + calinski_norm + davies_norm) / 4
75 best_k = k_values[np.argmax(combined_score)]
```

```

76 print(f"✅ Best k according to combined score : {best_k}")
77
78
79 # --Clustering final
80 kmeans_final = KMeans(n_clusters=best_k, random_state=42, n_init=10, algorithm="elkan")
81 df_segmented["segment_raw"] = kmeans_final.fit_predict(X_scaled)
82 df_segmented["segment"] = ["seg_" + str(i+1) for i in df_segmented["segment_raw"]]
83
84
85 # --Centroïdes par segment
86 centroids = df_segmented.groupby("segment")[features].mean().reset_index()
87 print("✅ centroids Tab :")
88 display(centroids)
89
90
91 # --Tableau enrichi final
92 df_result = df_segmented[[
93     "user_id", "email", "first_name", "segment", "country",
94     "total_orders", "avg_basket", "total_revenue",
95     "top_products_user", "last_event", "top_products_country",
96     "months_since_last_order", "freq_orders_per_month", "freq_orders_per_quarter"
97 ]]
98
99 print("✅ Final DataFrame with segmentation ready to use")
100 display(df_result.head())
101
102
103 # --Sauvegarde vers BigQuery
104 table_id = f"{PROJECT_ID}.ETL.semi_label_parametric"
105 job = client.load_table_from_dataframe(
106     df_result,
107     table_id,
108     job_config=bigquery.LoadJobConfig(write_disposition="WRITE_TRUNCATE"),
109 )
110 job.result()
111 print(f"✅ Results saved in {table_id}")
112
113
114 # --Graphiques diagnostics
115 fig = make_subplots(
116     rows=2, cols=3,
117     subplot_titles=(
118         "Inertia (Elbow)", "Silhouette", "Calinski-Harabasz", "Davies-Bouldin", "Score combiné", ""
119     )
120 )
121 fig.add_trace(go.Scatter(x=k_values, y=inertias, mode='lines+markers', name='Inertia'), row=1, col=1)
122 fig.add_trace(go.Scatter(x=k_values, y=silhouettes, mode='lines+markers', name='Silhouette'), row=1, col=2)
123 fig.add_trace(go.Scatter(x=k_values, y=calinski, mode='lines+markers', name='Calinski-Harabasz'), row=1, col=3)
124 fig.add_trace(go.Scatter(x=k_values, y=davies, mode='lines+markers', name='Davies-Bouldin'), row=2, col=1)
125 fig.add_trace(
126     go.Scatter(
127         x=k_values, y=combined_score, mode='lines+markers+text',
128         text=[f"{v:.2f}" for v in combined_score],
129         textposition="top center", name='Score combiné'), row=2, col=2)
130
131 fig.add_shape(
132     type="line", x0=best_k, x1=best_k, y0=0, y1=1,
133     line=dict(color="red", dash="dash"), row=2, col=2)
134
135 fig.update_layout(
136     height=800, width=1300,
137     title_text="Comparing metrics to determine the optimal number of clusters (K)",
138     showlegend=False
139 )
140 fig.update_xaxes(title_text="k (clusters)")
141 fig.update_yaxes(title_text="metric Values")
142 fig.show()
143

```



Best k according to combined score : 10
centroids Tab :

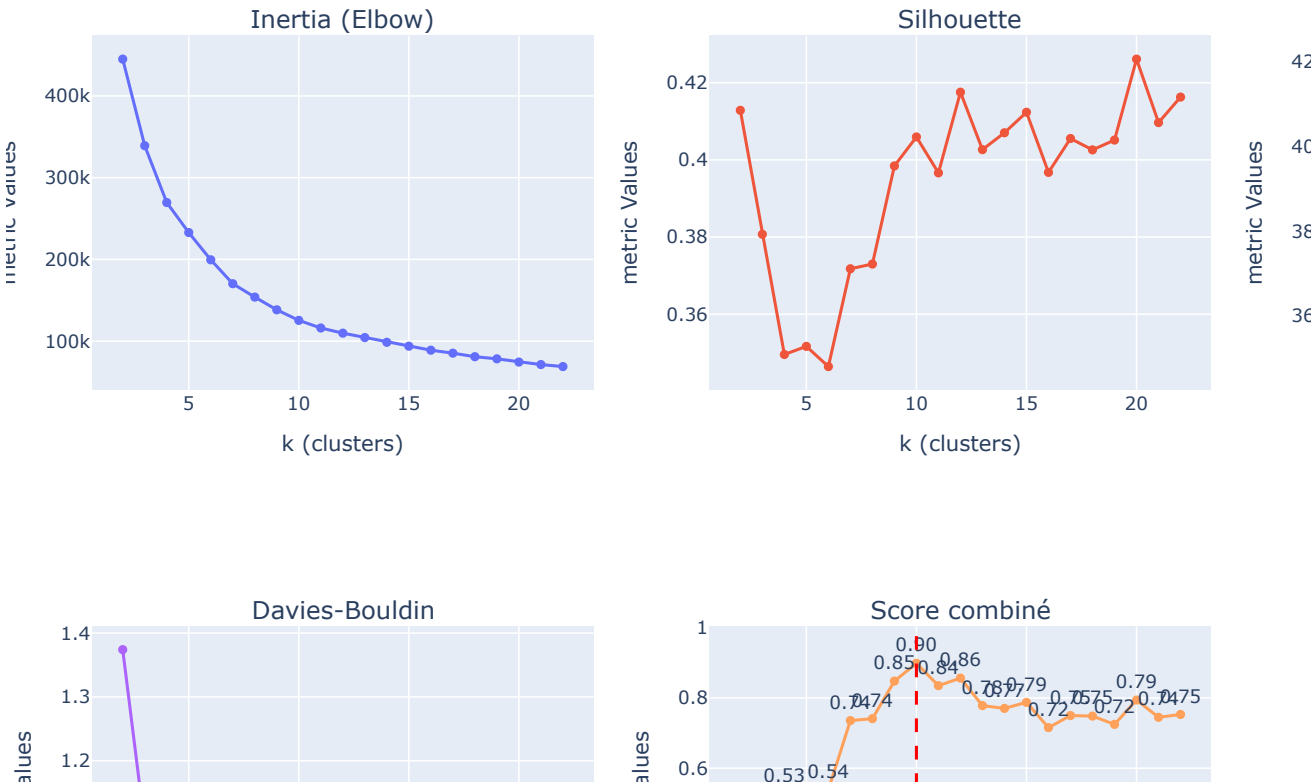
segment	total_orders	avg_basket	total_revenue	months_since_last_order	freq_orders_per_month	freq_orders_pe
seg_1	1.617485	148.680380	228.435312	15.482529	0.130856	
seg_10	6.092105	83.192744	475.947377	9.972588	0.296769	
seg_2	1.878878	47.875371	88.144137	2.035542	0.767615	
seg_3	0.061792	1.191212	1.191212	0.033671	0.001246	
seg_4	1.470852	41.461919	59.944883	11.297364	0.118169	
seg_5	5.010499	59.029332	295.659528	0.503937	4.716535	
seg_6	4.005712	52.644164	204.208923	11.151397	0.188059	
seg_7	1.711806	411.879188	636.499688	16.946181	0.204305	
seg_8	1.56305	46.868280	73.797649	43.569151	0.035904	
seg_9	3.450895	58.016835	200.115293	1.004838	2.035889	

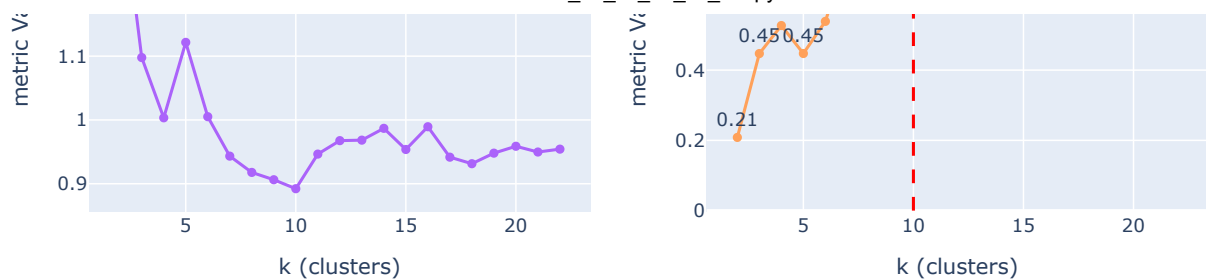
Final DataFrame with segmentation ready to use

user_id	email	first_name	segment	country	total_orders	avg_basket	total_revenue	top_pr
1	reneewilson@example.net	Renee	seg_10	United Kingdom	6	53.201667	319.209999	Women Neck
4	erikzavala@example.net	Erik	seg_1	United States	2	96.995000	193.990000	Sal Slirr
6	williambarber@example.net	William	seg_10	Brasil	7	77.068571	539.480000	[N2N - Sleeve
11	marialloyd@example.org	Maria	seg_8	United States	1	15.670000	15.670000	[Gr
12	samueltomas@example.org	Samuel	seg_8	Brasil	1	9.950000	9.950000	[H Cushion

Results saved in gothic-list-469006-v9.ETL.semi_label_parametric

Comparing metrics to determine the optimal number of clusters (K)





▼ Interactive 3D Visualization of **semi-labeled** Customer Segments using UMAP and PCA

```

1 from umap import UMAP
2 import plotly.graph_objects as go
3 import seaborn as sns
4 from matplotlib import colors as mcolors
5
6
7 # --UMAP 3D + PCA 3D
8 umap_3d = UMAP(n_components=3, n_jobs=-1, random_state=None)
9 X_umap = umap_3d.fit_transform(X_scaled)
10
11 pca_3d = PCA(n_components=3, random_state=42)
12 X_pca = pca_3d.fit_transform(X_scaled)
13
14
15 # --Palette de couleurs fixe
16 colors_palette = ['red', 'green', 'blue', 'olive', 'magenta', 'yellow', 'black',
17                  'orange', 'purple', 'brown', 'pink', 'lime', 'teal', 'navy',
18                  'gold', 'salmon', 'cyan', 'maroon', 'aqua', 'violet']
19
20 segments = df_segmented["segment"].unique()
21 color_map = {seg: colors_palette[i % len(colors_palette)] for i, seg in enumerate(segments)}
22
23
24 # --Graphique 3D interactif UMAP/PCA
25 fig_3d = go.Figure()
26
27 # --UMAP traces
28 for seg in segments:
29     idx = df_segmented["segment"] == seg
30     fig_3d.add_trace(go.Scatter3d(
31         x=X_umap[idx,0],
32         y=X_umap[idx,1],
33         z=X_umap[idx,2],
34         mode="markers",
35         marker=dict(size=6, color=color_map[seg], opacity=0.8),
36         text=df_segmented.loc[idx, "segment"],
37         name=f"{seg} (UMAP)",
38         visible=True
39     ))
40
41 # --PCA traces
42 for seg in segments:
43     idx = df_segmented["segment"] == seg
44     fig_3d.add_trace(go.Scatter3d(
45         x=X_pca[idx,0],
46         y=X_pca[idx,1],
47         z=X_pca[idx,2],
48         mode="markers",
49         marker=dict(size=5, color=color_map[seg], opacity=0.8),
50         text=df_segmented.loc[idx, "segment"],
51         name=f"{seg} (PCA)",
52         visible=False
53     ))
54
55 # --Boutons UMAP/PCA

```

```

56 n = len(segments)
57 fig_3d.update_layout(
58     updatemenus=[dict(
59         type="buttons",
60         direction="right",
61         x=0.5, y=1.1,
62         buttons=[
63             dict(label="UMAP", method="update",
64                 args=[{"visible": [True]*n + [False]*n},
65                     {"scene": dict(xaxis_title="UMAP-1", yaxis_title="UMAP-2", zaxis_title="UMAP-3"),
66                         "title": "Projection UMAP 3D"}]),
67             dict(label="PCA", method="update",
68                 args=[{"visible": [False]*n + [True]*n},
69                     {"scene": dict(xaxis_title="PCA-1", yaxis_title="PCA-2", zaxis_title="PCA-3"),
70                         "title": "Projection PCA 3D"}])
71         ]
72     )]
73 )
74
75 # --Layout final
76 fig_3d.update_layout(
77     title="Interactive 3D projection of segments",
78     scene=dict(xaxis_title="UMAP-1", yaxis_title="UMAP-2", zaxis_title="UMAP-3"),
79     height=700, width=950,
80     legend=dict(
81         title="Segments",
82         orientation="v",
83         yanchor="top",
84         y=1,
85         xanchor="left",
86         x=-0.1
87     )
88 )
89
90 fig_3d.show()
91

```



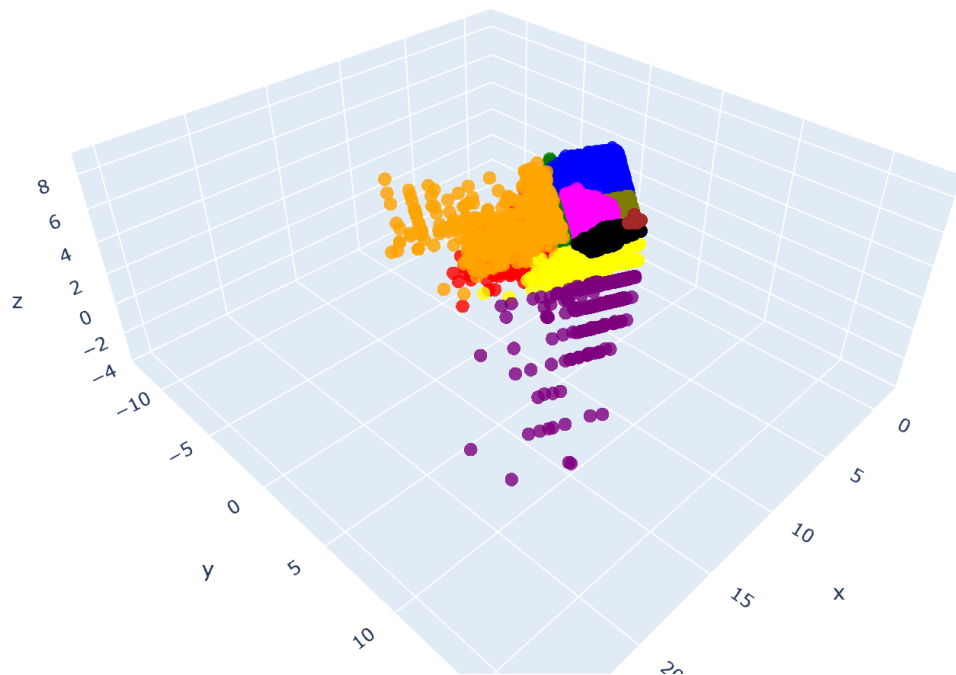
Projection PCA 3D

UMAP

PCA

Segments

- seg_10 (PCA)
- seg_1 (PCA)
- seg_8 (PCA)
- seg_4 (PCA)
- seg_6 (PCA)
- seg_9 (PCA)
- seg_2 (PCA)
- seg_7 (PCA)
- seg_5 (PCA)
- seg_3 (PCA)



Automated Customer Segmentation with K-Means and BigQuery

Here we're going to load customer data from BigQuery, applies K-Means to find the optimal clusters, assigns users to default segments (not custom ones), and produces both a segmented customer dataset and cluster centroids, saved in `semi_label_parametric` and `centroid_data` for advanced labeling purpose

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.cluster import KMeans
4 from sklearn.metrics import silhouette_score, calinski_harabasz_score, davies_bouldin_score
5 from sklearn.preprocessing import StandardScaler, MinMaxScaler
6 from google.cloud import bigquery
7
8
9 # --Charger les données depuis BigQuery
10 query = f"""
11     SELECT *
12     FROM `{PROJECT_ID}.ETL.parametric`
13 """
14
15 client = bigquery.Client(project=PROJECT_ID)
16 df = client.query(query).to_dataframe()
17
18
19 # --Features & Préparation
20 features = [
21     "total_orders",
22     "avg_basket",
23     "total_revenue",

```

```

24     "months_since_last_order",
25     "freq_orders_per_month",
26     "freq_orders_per_quarter"
27 ]
28
29 X = df[features].fillna(0)
30 X_scaled = StandardScaler().fit_transform(X)
31
32
33 # --Évaluer les k avec plusieurs métriques
34 k_values = list(range(2, 21))
35 inertias, silhouettes, calinski, davies = [], [], [], []
36
37 for k in k_values:
38     kmeans = KMeans(n_clusters=k, random_state=42, n_init=10, algorithm="elkan")
39     labels = kmeans.fit_predict(X_scaled)
40
41     inertias.append(kmeans.inertia_)
42     silhouettes.append(silhouette_score(X_scaled, labels))
43     calinski.append(calinski_harabasz_score(X_scaled, labels))
44     davies.append(davies_bouldin_score(X_scaled, labels))
45
46
47 # --Normalisation des scores & choix du meilleur k
48 scaler = MinMaxScaler()
49 inertia_norm = 1 - scaler.fit_transform(np.array(inertias).reshape(-1,1)).flatten()
50 silhouette_norm = scaler.fit_transform(np.array(silhouettes).reshape(-1,1)).flatten()
51 calinski_norm = scaler.fit_transform(np.array(calinski).reshape(-1,1)).flatten()
52 davies_norm = 1 - scaler.fit_transform(np.array(davies).reshape(-1,1)).flatten()
53
54 combined_score = (inertia_norm + silhouette_norm + calinski_norm + davies_norm) / 4
55 best_k = k_values[np.argmax(combined_score)]
56
57 print(f"✅ Best k according to combined score : {best_k}")
58
59
60 # --Clustering final avec meilleur k
61 kmeans_final = KMeans(n_clusters=best_k, random_state=42, n_init=10, algorithm="elkan")
62 df["segment_raw"] = kmeans_final.fit_predict(X_scaled)
63 df["segment"] = ["seg_" + str(i+1) for i in df["segment_raw"]]
64
65
66 # --Tableau enrichi final
67 df_result = df[[
68     "user_id", "email", "first_name", "segment", "country",
69     "total_orders", "avg_basket", "total_revenue",
70     "top_products_user", "last_event", "top_products_country",
71     "months_since_last_order", "freq_orders_per_month", "freq_orders_per_quarter"
72 ]]
73
74 print("✅ Final DataFrame with segmentation ready to use")
75 display(df_result.head())
76
77
78 # --Centroïdes par segment
79 centroids = df.groupby("segment")[features].mean().reset_index()
80 print("✅ Centroids table :")
81 display(centroids)
82
83
84 # --Sauvegarde vers BigQuery
85
86 # --Table segmentation enrichie
87 table_segmentation = f"{PROJECT_ID}.ETL.semi_label_parametric"
88 job1 = client.load_table_from_dataframe(
89     df_result,
90     table_segmentation,
91     job_config=bigquery.LoadJobConfig(
92         write_disposition="WRITE_TRUNCATE"
93     ),
94 )
95 job1.result()

```

```

96 print(f"✅ Segmentation results saved in {table_segmentation}")
97
98 # --Table centroides
99 table_centroids = f"{PROJECT_ID}.ETL.centroid_data"
100 job2 = client.load_table_from_dataframe(
101     centroids,
102     table_centroids,
103     job_config=bigquery.LoadJobConfig(
104         write_disposition="WRITE_TRUNCATE"
105     ),
106 )
107 job2.result()
108 print(f"✅ Centroids results saved in {table_centroids}")
109

```



- ✅ Best k according to combined score : 10
- ✅ Final DataFrame with segmentation ready to use

	user_id	email	first_name	segment	country	total_orders	avg_basket	total_revenue	top_pr
0	1	reneewilson@example.net	Renee	seg_10	United Kingdom	6	53.201667	319.209999	Wome Neck
1	4	erikzavala@example.net	Erik	seg_1	United States	2	96.995000	193.990000	Sa Slir
2	6	williambarber@example.net	William	seg_10	Brasil	7	77.068571	539.480000	[N2N Slee
3	11	marialloyd@example.org	Maria	seg_8	United States	1	15.670000	15.670000	[G
4	12	samueltthomas@example.org	Samuel	seg_8	Brasil	1	9.950000	9.950000	[H Cushic

✅ Centroids table :

	segment	total_orders	avg_basket	total_revenue	months_since_last_order	freq_orders_per_month	freq_orders_p
0	seg_1	1.617485	148.680380	228.435312	15.482529	0.130856	
1	seg_10	6.092105	83.192744	475.947377	9.972588	0.296769	
2	seg_2	1.878878	47.875371	88.144137	2.035542	0.767615	
3	seg_3	0.061792	1.191212	1.191212	0.033671	0.001246	
4	seg_4	1.470852	41.461919	59.944883	11.297364	0.118169	
5	seg_5	5.010499	59.029332	295.659528	0.503937	4.716535	
6	seg_6	4.005712	52.644164	204.208923	11.151397	0.188059	
7	seg_7	1.711806	411.879188	636.499688	16.946181	0.204305	
8	seg_8	1.56305	46.868280	73.797649	43.569151	0.035904	
9	seg_9	3.450895	58.016835	200.115293	1.004838	2.035889	

- ✅ Segmentation results saved in gothic-list-469006-v9.ETL.semi_label_parametric
- ✅ Centroids results saved in gothic-list-469006-v9.ETL.centroid_data

There're 2 approaches for advanced labeling purpose, especially the **HAND-ON** and **HAND-OFF** approaches, then let's dive in

✓ 1. HAND ON

Here we define a mapping data object to create a correspondence and assign the correct label to the initial standard label based on centroid data analysis. However, it requires some expertise to determine the most accurate customer designation

```

1 %%bigquery df_label_parametric_Semi_hand_off --project $PROJECT_ID
2
3 CREATE OR REPLACE TABLE `gothic-list-469006-v9.ETL.label_parametric_Sho` AS #--sho for Semi hand_on

```

```

4 SELECT
5   P.*,
6   CASE segment
7     WHEN 'seg_1' THEN 'Dormant'
8     WHEN 'seg_2' THEN 'At Risk'
9     WHEN 'seg_3' THEN 'Occasional Big Spenders'
10    WHEN 'seg_4' THEN 'New / Active Low Spenders'
11    WHEN 'seg_5' THEN 'Loyal Customers'
12    WHEN 'seg_6' THEN 'Churned / Lost'
13    WHEN 'seg_7' THEN 'One-Timers'
14    WHEN 'seg_8' THEN 'Active Mid-Value'
15    WHEN 'seg_9' THEN 'Active / Promising'
16    WHEN 'seg_10' THEN 'High-Value Dormant'
17    ELSE 'Unknown'
18  END AS segments
19 FROM `gothic-list-469006-v9.ETL.semi_label_parametric` P;
20

```

 Job ID 428fcc19-25ff-4752-b334-7b6e4307a955 successfully executed: 100%


✓ 2. HAND OFF

That's what we're going to use next, so keep your eyes open. To carry the approach forward, we're going to use `ML.GENERATE_TEXT` function, by parsing the parametric centroid data along with a specific prompt to generate high-quality designations for each row, in order to build a highly customized mapping object in a hands-off way

```

1 %%bigquery df_label_parametric_Hand_off --project $PROJECT_ID
2
3 CREATE OR REPLACE TABLE `gothic-list-469006-v9.ETL.label_parametric` AS
4 WITH AI_labels AS (
5   SELECT
6     t.segment AS segment,
7     -- Extract text from JSON output and then clean
8     TRIM(SPLIT(SAFE_CAST(JSON_VALUE(t.ml_generate_text_result, '$.candidates[0].content.parts[0].text') AS STRING)) AS STRING)
9     TRIM(SPLIT(SAFE_CAST(JSON_VALUE(t.ml_generate_text_result, '$.candidates[0].content.parts[0].text') AS STRING)) AS STRING)
10  FROM ML.GENERATE_TEXT(
11    MODEL `gothic-list-469006-v9.marketing_email_model_dataset.marketing_email_model`,
12    (
13      SELECT
14        segment,
15        CONCAT(
16          'Segment: ', segment, '. ',
17          'Total Orders: ', total_orders, '. ',
18          'Average Basket: ', avg_basket, '. ',
19          'Total Revenue: ', total_revenue, '. ',
20          'Months Since Last Order: ', months_since_last_order, '. ',
21          'Frequency per Month: ', freq_orders_per_month, '. ',
22          'Frequency per Quarter: ', freq_orders_per_quarter, '. ',
23          'Based on this data, you MUST assign a concise human-readable segment label, (without bold formatting) j
24          'First line = segment label, second line = short designation.'
25        ) AS prompt
26      FROM `gothic-list-469006-v9.ETL.centroid_data`
27    ),
28    STRUCT(0.2 AS temperature, 150 AS max_output_tokens, 0.7 AS top_p, 1 AS top_k)
29  ) AS t
30 )
31 SELECT
32   P.*,
33   A.segments,
34   A.designation
35 FROM `gothic-list-469006-v9.ETL.semi_label_parametric` P
36 LEFT JOIN AI_labels A
37 USING(segment);

```

 Job ID 1459d1f6-93d2-4c66-995f-7375bdcdbd959 successfully executed: 100%

Here's the outcome related to the `HAND OFF` approach of the labeling process

```
1 %%bigquery --project {PROJECT_ID}
2
3 SELECT
4   user_id,
5   email,
6   first_name,
7   segment,
8   segments AS type_of_customer,
9   designation,
10  country,
11  total_orders,
12  avg_basket,
13  total_revenue,
14  top_products_user,
15  last_event,
16  top_products_country,
17  months_since_last_order,
18  freq_orders_per_month,
19  freq_orders_per_quarter
20 FROM `gothic-list-469006-v9.ETL.label_parametric`
21 ORDER BY user_id;
```



Job ID 23b48763-7813-4798-be66-f777fdd5e492 successfully executed: 100%

Downloading: 100%

	user_id	email	first_name	segment	type_of_customer	designation	country	total_or
0	1	reneewilson@example.net	Renee	seg_10	Potential Churn		United Kingdom	
1	2	valerieporter@example.net	Valerie	seg_3	Infrequent Recent Buyers	Low-value, very recent customers with minimal ...	United States	
2	3	gabrielhuff@example.org	Gabriel	seg_3	Infrequent Recent Buyers	Low-value, very recent customers with minimal ...	United States	
3	4	erikzavala@example.net	Erik	seg_1	Lapsed Buyers	Infrequent purchasers with a declining engagem...	United States	
4	5	adamfrazier@example.org	Adam	seg_3	Infrequent Recent Buyers	Low-value, very recent customers with minimal ...	China	
...
99995	99996	taylorle@example.net	Taylor	seg_6	Lapsed Regulars	Infrequent purchasers with a moderate average ...	United States	
99996	99997	mackenzieramirez@example.com	Mackenzie	seg_4	Lapsed Purchasers	Infrequent buyers with a high average basket v...	United States	
99997	99998	kellygoodman@example.org	Kelly	seg_6	Lapsed Regulars	Infrequent purchasers with a moderate average ...	Brasil	
99998	99999	maryjohnson@example.net	Mary	seg_2	Recent High-Value Customer	Recently acquired customer with a high average...	Brasil	
99999	100000	amandaholloway@example.com	Amanda	seg_6	Lapsed Regulars	Infrequent purchasers with a moderate average ...	China	

100000 rows × 16 columns

Interactive 3D Visualization for now **labeled** Customer Segments using UMAP and PCA

```
1 import pandas as pd
2 import numpy as np
3 import re
4 from google.cloud import bigquery
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.impute import SimpleImputer
7 from sklearn.decomposition import PCA
8 from umap import UMAP
```




```

8 from umap import UMAP
9 import plotly.graph_objects as go
10
11
12 # --Connexion BigQuery et chargement des données
13 client = bigquery.Client(project=PROJECT_ID)
14
15 query = f"""
16 SELECT *
17 FROM `{PROJECT_ID}.ETL.label_parametric`
18 """
19 df = client.query(query).to_dataframe()
20
21
22 # --Features numériques
23 features = [
24     "total_orders", "avg_basket", "total_revenue",
25     "months_since_last_order", "freq_orders_per_month", "freq_orders_per_quarter"
26 ]
27
28 X = df[features]
29 imputer = SimpleImputer(strategy="mean")
30 X_imputed = imputer.fit_transform(X)
31
32 scaler = StandardScaler()
33 X_scaled = scaler.fit_transform(X_imputed)
34
35
36 # --PCA 3D
37 pca_3d = PCA(n_components=3, random_state=42)
38 X_pca = pca_3d.fit_transform(X_scaled)
39 df["PCA1"], df["PCA2"], df["PCA3"] = X_pca[:,0], X_pca[:,1], X_pca[:,2]
40
41
42 # --UMAP 3D
43 umap_3d = UMAP(n_components=3, random_state=42, n_jobs=-1)
44 X_umap = umap_3d.fit_transform(X_scaled)
45 df["UMAP1"], df["UMAP2"], df["UMAP3"] = X_umap[:,0], X_umap[:,1], X_umap[:,2]
46
47
48 # --Palette de couleurs et mapping segments
49 colors_palette = [
50     'red','green','blue','cyan','magenta','gold','black',
51     'orange','purple','brown','pink','lime','teal','navy',
52     'yellow','salmon','olive','maroon','aqua','violet'
53 ]
54
55 segments = df["segments"].unique()
56
57 color_map = {seg: colors_palette[i % len(colors_palette)] for i, seg in enumerate(segments)}
58
59 # Nettoyage pour enlever les * au début et à la fin
60 # clean_segments = [re.sub(r"^\*+|\*+$", "", str(seg)).strip() for seg in segments]
61
62 # Génération du color_map avec les segments nettoyés
63 # color_map = {clean_seg: colors_palette[i % len(colors_palette)] for i, clean_seg in enumerate(clean_segments)}
64
65
66 # --Graphique 3D interactif PCA et UMAP
67 fig = go.Figure()
68
69 # Ajouter PCA 3D (visible par défaut)
70 for seg in segments:
71     # clean_seg = re.sub(r"^\*+|\*+$", "", str(seg)).strip()
72     df_seg = df[df["segments"] == seg]
73     fig.add_trace(go.Scatter3d(
74         x=df_seg["PCA1"], y=df_seg["PCA2"], z=df_seg["PCA3"],
75         mode='markers',
76         marker=dict(size=5, color=color_map[seg]),
77         name=f"{seg}",
78         visible=True
79     ))
80

```

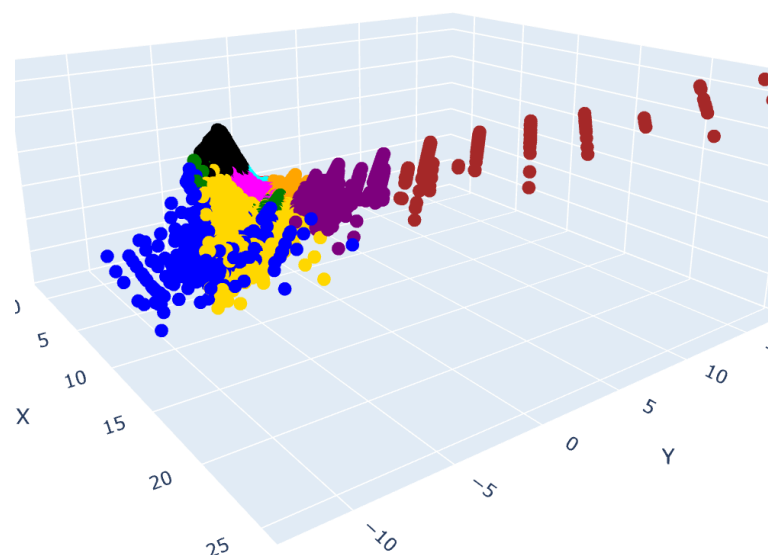
```
81 # --Ajouter UMAP 3D (invisible au départ)
82 for seg in segments:
83     # clean_seg = re.sub(r"^\*+|\*+$", "", str(seg)).strip()
84     df_seg = df[df["segments"] == seg]
85     fig.add_trace(go.Scatter3d(
86         x=df_seg["UMAP1"], y=df_seg["UMAP2"], z=df_seg["UMAP3"],
87         mode='markers',
88         marker=dict(size=5, color=color_map[seg]),
89         name=f"{seg}",
90         visible=False
91     ))
92
93 # --Dropdown pour basculer entre PCA et UMAP
94 buttons = [
95     dict(label="PCA 3D",
96         method="update",
97         args=[{"visible": [True]*len(segments) + [False]*len(segments)},
98             {"title": "PCA 3D customers per segment"}]),
99     dict(label="UMAP 3D",
100         method="update",
101         args=[{"visible": [False]*len(segments) + [True]*len(segments)},
102             {"title": "UMAP 3D customer per segment"}])
103 ]
104
105 fig.update_layout(
106     updatemenus=[dict(active=0, buttons=buttons)],
107     height=700, width=900,
108     title="PCA 3D of customer segment",
109     scene=dict(
110         xaxis_title='X',
111         yaxis_title='Y',
112         zaxis_title='Z'
113     )
114 )
115
116 fig.show()
117
```

 /usr/local/lib/python3.12/dist-packages/umap/umap_.py:1952: UserWarning:

n_jobs value 1 overridden to 1 by setting random_state. Use no seed for parallelism.

PCA 3D customers per segment

PCA 3D ▼



- Infrequent Recent Buyers
- Lapsed Buyers
- Lapsed High-Value
- Lapsed Purchasers
- Lapsed Regulars
- Potential Churn
- Dormant Low-Value
- Recent High-Value Customer
- Recent High-Frequency Spender
- Highly Active Recent Buyers

Customer Segmentation and Multilingual Marketing Email Prompt Generation in BigQuery

Now we build a final dataset: `final_data` in BigQuery (`ETL.label_parametric`) by enriching customer data with country-based language codes and detailed behavioral attributes. It assigns each customer a preferred communication language based on their country, then generates a dynamic marketing prompt tailored to their profile. The prompt incorporates purchase history, frequency, basket size, last activity, and product preferences and so on, while adapting the tone to their customer segment previously defined in the hand off labeling process. The final output provides a structured foundation for creating personalized, multilingual, and engaging marketing emails designed.

```
1 %%bigquery df_final_data --project {PROJECT_ID}
2
3 CREATE OR REPLACE TABLE `gothic-list-469006-v9.marketing_email_model_dataset.final_data` AS
4 WITH segmented AS (
5   SELECT
6     b.*,
7     CASE
8       WHEN LOWER(b.country) IN (
9         'france', 'belgium', 'switzerland', 'luxembourg',
10        'senegal', 'ivory coast', 'cameroon',
11        'madagascar', 'togo', 'mali', 'morocco', 'tunisia', 'algeria'
12      ) THEN 'fr'
13       WHEN LOWER(b.country) IN (
14        'spain', 'mexico', 'argentina', 'chile', 'colombia', 'peru',
```

```

15     'venezuela', 'ecuador', 'uruguay', 'paraguay', 'bolivia',
16     'guatemala', 'cuba', 'dominican republic'
17 ) THEN 'es'
18 WHEN LOWER(b.country) IN ('brasil', 'portugal', 'mozambique', 'angola') THEN 'pt'
19 WHEN LOWER(b.country) IN (
20     'united states', 'usa', 'united kingdom', 'canada', 'england', 'australia',
21     'ireland', 'nigeria', 'ghana', 'kenya', 'south africa', 'india',
22     'philippines'
23 ) THEN 'en'
24 WHEN LOWER(b.country) IN (
25     'egypt', 'saudi arabia', 'uae', 'united arab emirates', 'qatar',
26     'jordan', 'lebanon', 'iraq', 'syria', 'yemen'
27 ) THEN 'ar'
28 WHEN LOWER(b.country) IN ('china', 'taiwan', 'singapore') THEN 'zh'
29 WHEN LOWER(b.country) IN ('japan') THEN 'ja'
30 WHEN LOWER(b.country) IN ('germany') THEN 'de'
31 WHEN LOWER(b.country) IN ('south korea', 'korea', 'republic of korea') THEN 'ko'
32 ELSE 'en'
33 END AS lang
34 FROM `gothic-list-469006-v9.ETL.label_parametric` b
35 ),
36
37 final_data AS (
38     SELECT
39         s.user_id,
40         s.email,
41         s.first_name,
42         s.country,
43         s.lang,
44         s.segment,
45         s.segments,
46         s.designation,
47         s.total_orders,
48         s.avg_basket,
49         s.freq_orders_per_month,
50         s.freq_orders_per_quarter,
51         s.last_event,
52         s.top_products_user,
53         s.top_products_country,
54         CONCAT(
55             "Write a short, friendly, wel structured and engaging marketing email", "just the ","result"," in language o
56             "This email is intended for a ", "",s.segments,"", "without explicitelty mentioned", s.segments,"", "custo
57             " from ", s.country, ". ",
58             "They have made ", CAST(s.total_orders AS STRING), " purchases this year with an average basket of $", FORMA
59             "They buy approximately ", FORMAT('%.2f', s.freq_orders_per_month), " times per month. ",
60             "Their last activity: ", IFNULL(s.last_event, 'no recent events'), ". ",
61             "With Customer's top products: ", ARRAY_TO_STRING(s.top_products_user, ', '), ". ",
62             "Top trending products in their country: ", ARRAY_TO_STRING(s.top_products_country, ', '), ". ",
63             "Adapt tone based on segment: exclusive for VIP, motivating for Active, welcoming for Dormant, and exciting
64             "The goal is to maximize customer retention, attention and encourage engagement depending on", s.segments, "
65         ) AS prompt
66     FROM segmented s
67 )
68
69 SELECT * FROM final_data ORDER BY user_id;
70

```



Job ID 0fced6ca-6a48-4bf5-b5f1-271cac7a1fe8 successfully executed: 100%

✓ The use of ML.GENERATE_TEXT() to generate Personalized email for business activity purpose

So far, everything is working well. At this stage, our main focus is to extract first the email and prompt features along with any other optional attributes and feed them into [ML.GENERATE_TEXT\(\)](#) to generate hyper-personalized emails, and there after storing through marketing_email_P500 ou marketing_email_Pxxx

```
1 %%bigquery Hyper_personalized_email --project {PROJECT_ID}
2
3 # -- Sélectionner top 500 prompts
4 CREATE OR REPLACE TABLE `gothic-list-469006-v9.marketing_email_model_dataset.marketing_email_P500` AS
5 WITH prompts AS (
6   SELECT
7     email,
8     country,
9     lang AS language_code,
10    prompt,
11    ROW_NUMBER() OVER (ORDER BY FARM_FINGERPRINT(email)) AS prompt_id
12  FROM `gothic-list-469006-v9.marketing_email_model_dataset.final_data`
13  LIMIT 500
14  -- QUALIFY prompt_id <= 50  pour restreindre
15 ),
16
17 -- Génération du contenu marketing avec ML.GENERATE_TEXT() function.
18 marketing_mail AS (
19   SELECT
20     p.email,
21     p.language_code,
22     p.country,
23     SAFE_CAST(
24       JSON_VALUE(
25         ml_output.ml_generate_text_result,
26         '$.candidates[0].content.parts[0].text'
27       ) AS STRING
28     ) AS marketing_email
29  FROM
30    ML.GENERATE_TEXT(
31      MODEL `gothic-list-469006-v9.marketing_email_model_dataset.marketing_email_model`,
32      (
33        SELECT prompt, prompt_id FROM prompts
34      ),
35      STRUCT(
36        300 AS max_output_tokens,
37        0.7 AS temperature
38        -- ['\n\n'] AS stop_sequences
39      )
40    ) AS ml_output
41  JOIN prompts p
42    ON ml_output.prompt = p.prompt
43   AND ml_output.prompt_id = p.prompt_id
44 )
45
46 -- Résultats finaux filtrés et ordonnés
47 SELECT
48   email,
49   country,
50   language_code,
51   marketing_email
52 FROM marketing_mail
53 WHERE marketing_email IS NOT NULL
54 ORDER BY email;
```