

🌟 Draft Report: Data Cleaning Sehhaty App Reviews

1. Overview

There were **128,406** user reviews collected using the Sehhaty mobile application. Data cleaning was aimed to preprocess for providing quality by removing invalid values, eliminating short or absent reviews, deleting duplicates, and handling missing values, thereby preparing the dataset for future analysis. Reviews were collected from the Google Play Store between **2019 and 2025**, providing a longitudinal dataset in its entirety. Arabic and English reviews were cleaned and split apart individually with language-specific stopword lists and error correction rules to meet linguistic consistency and prevent cross-language contamination. This laid a strong foundation for the following analyses such as sentiment analysis and topic modeling.

(ar) تم جلب تعليق 88600	✓	✓ Retrieved 36000 reviews (en-us)...
(ar) تم جلب تعليق 88800	✓	✓ Retrieved 36200 reviews (en-us)...
(ar) تم جلب تعليق 89000	✓	✓ Retrieved 36400 reviews (en-us)...
(ar) تم جلب تعليق 89200	✓	✓ Retrieved 36600 reviews (en-us)...
(ar) تم جلب تعليق 89400	✓	✓ Retrieved 36800 reviews (en-us)...
(ar) تم جلب تعليق 89600	✓	✓ Retrieved 37000 reviews (en-us)...
(ar) تم جلب تعليق 89800	✓	✓ Retrieved 37200 reviews (en-us)...
(ar) تم جلب تعليق 90000	✓	✓ Retrieved 37400 reviews (en-us)...
(ar) تم جلب تعليق 90200	✓	✓ Retrieved 37600 reviews (en-us)...
(ar) تم جلب تعليق 90400	✓	✓ Retrieved 37800 reviews (en-us)...
(ar) تم جلب تعليق 90450	✓	✓ Retrieved 37956 reviews (en-us)...
تم حفظ 90450 تعليق عربي	✓	

2. Cleaning Steps

1 Invalid Values Removal

- ✓ Ratings were restricted to the valid range [1–5].

2 Text Cleaning

- ✓ Special symbols, repeated characters, diacritics, and URLs were removed.
- ✓ Language-specific stopword lists were applied to eliminate common phrases such as greetings and filler words, reducing linguistic noise.
- ✓ Emojis were intentionally retained, as they are important for expressing sentiment and emotional context and thus improve the accuracy of sentiment analysis.

```
# ----- Stopwords -----
STOPWORDS_AR_SENTIMENT = {
    'السلام', 'السلام عليك', 'عليكم', 'وعليكم',
    '، صباح الخير', 'مساء الخير',
    'مرحبا', 'اهلا', 'اهلين', 'هلا',
    'الله', 'اللهم', 'امين', 'آمين', 'يارب',
    'كذا', 'ترى', 'اصلا', 'اصلا', 'اصلا',
    'بس', 'يعني', 'اوكي', 'اوك', 'كمان',
    'ههه', 'ههههه', 'خخ', 'خخ', 'ممم', 'موه',
    '...', '...', '...', '...', '...' ,
    'apk', 'play', 'http', 'https', 'www', 'com', 'org',
    '500', '11', '666',
    'ذا', 'ورحمة', 'ورحمة', 'ورحمة',
    'وباركاته', 'وباركاته'
}
```

```
# ===== English Stopwords & Fixes =====
STOPWORDS_EN_SENTIMENT = {
    'the', 'a', 'an', 'and', 'or', 'but', 'to', 'of', 'in', 'on',
    'it', 'its', 'this', 'that', 'these', 'those', 'there', 'I',
    'also', 'just', 'still',
    'i', 'me', 'my', 'we', 'our', 'you', 'your', 'they', 'them',
    'be', 'am', 'is', 'are', 'was', 'were', 'been', 'being',
    'do', 'did', 'does', 'doing',
    'have', 'has', 'had', 'having',
    'will', 'would', 'should', 'could', 'can', 'may', 'might',
    'if', 'then', 'than', 'because', 'while', 'when', 'where',
    'etc', 'etc.', 'etc..',
    'hi', 'hello', 'or', 'is', 'igigi'
}
```

3. Text Normalization and Spelling Corrections

📝 Instead of discarding short reviews, a **dictionary of common spelling corrections (COMMON_FIX)** was applied for both Arabic and English text.

- This step converted frequently misspelled words into their correct forms

Arabic: رأي → رائع ، ممتازه → ممتازة

English: dont → don't

✨ By using correction rather than deletion, even short but meaningful reviews were preserved, ensuring that no potentially important words or information for the analysis were lost.

```
# ----- Common Fix -----
COMMON_FIX = {
    'ممتازة', 'ممتاز', 'ممتاز', 'ممتاز': 'ممتاز', 'ممتاز': 'ممتاز',
    'ممتاز جدا', 'ممتاز جدا': 'ممتاز جدا', 'ممتاز': 'ممتاز جدا',
    'جدا': 'رائع', 'رأيده': 'رائعة', 'رأيده': 'رائعة', 'رأي': 'رائع',
    'رووووه': 'روووهه', 'روووهه': 'روووهه', 'روووهه': 'روووهه',
    'جميل', 'جيبلل': 'جميل', 'جيبله': 'جميلة', 'جميله': 'جميلة',
    'جدا': ' جدا', 'حن': 'حن', 'حن': 'حن', 'حن': 'حن',
    'سي': 'سيي', 'سيي': 'سيي', 'سييه': 'سيي', 'سييه': 'سيي',
    'سي': 'سيي', 'سيي': 'سيي', 'سي': 'سيي',
    'بس': 'بس', 'بس': 'بس', 'بس': 'بس', 'بس': 'بس',
    ' يعني': ' يعني', ' يعني': ' يعني', ' يعني': ' يعني',
    ' يعني': ' يعني', ' يعني': ' يعني',
    'ال تمام': 'ال تمام', 'ال تمام': 'ال تمام', 'ال تمام': 'ال تمام',
    'ال شكر': 'ال شكر', 'ال شكر': 'ال شكر', 'ال شكر': 'ال شكر',
    'واش كركم': 'اشكركم', 'واشكركم': 'اشكركم', 'واش': 'اشكركم',
    'مشكله': 'مشكله', 'مشكله': 'مشكله', 'المشكله': 'المشكله',
    'التطبيق لا يصل': 'التطبيق لا يصل', 'التطبيق لا يصل': 'التطبيق لا يصل'
}
```

```
COMMON_FIX_EN = {
    "dont": "don't", "cant": "can't", "wont": "won't", "didn't": "didn't",
    "isn't": "isn't", "wasn't": "wasn't", "shouldnt": "shouldn't",
    "im": "i'm", "ive": "i've", "id": "i'd", "ill": "i'll", "u": "u",
    "plz": "please", "thx": "thanks", "thanx": "thanks",
    "definately": "definitely", "definetly": "definitely",
    "adress": "address", "addres": "address", "occured": "o",
    "enviroment": "environment", "acheive": "achieve", "wi": "wi",
    "apointment": "appointment", "appoimentement": "appoint",
    "dose": "does", "gud": "good", "bcoz": "because", "bcz": "because",
    "aap": "app", "goid": "good"
}
```

📸 Snapshot of the Preprocessing Code

📷 The images above show a snapshot of the preprocessing code, demonstrating the use of **language-specific stopword lists and spelling correction dictionaries (COMMON_FIX)** for both Arabic and English.

- These steps ensured **linguistic consistency** by removing irrelevant words and normalizing frequent misspellings before analysis.

Arabic — Stopwords & Corrections

Original	After	Status	Action
أمين	(removed)		Delete
ابل	(removed)		Delete
اصل	(removed)		Delete
اصلأ	(removed)		Delete
اصلأ	(removed)		Delete
الخ	(removed)		Delete
الخ.	(removed)		Delete
السلام	(removed)		Delete
السلام عليكم	(removed)		Delete
الصراحة	(removed)		Delete
الله	(removed)		Delete
اللهم	(removed)		Delete

English — Stopwords & Corrections

Original	After	Status	Action
a	(removed)		Delete
also	(removed)		Delete
am	(removed)		Delete
an	(removed)		Delete
and	(removed)		Delete
are	(removed)		Delete
as	(removed)		Delete
at	(removed)		Delete
be	(removed)		Delete

بطئ	بطيء		Correct
بطى	بطيء		Correct
بطيري	بطيء		Correct
جميلال	جميل		Correct
جميله	جميلة		Correct
جميل	جميل		Correct
جميله	جميلة		Correct
جميلل	جميل		Correct
حلووو	حلو		Correct
حلوووو	حلو		Correct
رائع	رائع		Correct
رائع	رائع		Correct

aap	app		Correct
accomodate	accommodate		Correct
acheive	achieve		Correct
addres	address		Correct
adress	address		Correct
apointment	appointment		Correct
appointement	appointment		Correct
apponitment	appointment		Correct
bc	because		Correct
bcoz	because		Correct
bcz	because		Correct
cant	can't		Correct

Illustration of Preprocessing Steps

- The lists above illustrate examples of text preprocessing steps:

- **Delete:** Certain words were removed as stopwords.
- **Correct:** Other words were corrected for spelling and wording.

💡 These actions helped keep the dataset **clean, normalized, and ready** for subsequent analysis.

4. Duplicates Removal

♻️ Process:

- 📁 Duplicates were removed based on ( **UserName** +  **Content**).
- 🚫 Generic placeholders such as ("A Google user" / "مستخدم Google") were ignored during duplicate detection.
- ✅ Only the first occurrence for each pair was kept.
- 📝 Only duplicate rows were removed; no other columns were modified.

📊 Results:

- Arabic: pair count **30** → **1** after deduplication.
- English: pair count **3** → **1** after deduplication.

📊 Cleaned & Deduplicated Results (Arabic & English)

- Shows record counts **before** and **after** cleaning for both languages.
- Displays how duplicate pairs were removed (keeping only the first occurrence for each pair).
- Highlights samples of cleaned Arabic and English reviews.

Before Cleaning	After Cleaning + Dedup	Removals
90,450	86,120	Dropped: 1,925 Duplicates: 2,405

Arabic — Before (User + CLEAN) | pair count: 30

User Name	Content Arabic Clean
ابو محمد	ممتاز

Arabic — After (dedup → 1)

User Name	Content Arabic Clean
ابو محمد	ممتاز

📊 Sample: Output (Before vs Cleaned)

⭐ Before_Content	🧹 Cleaned_Content
ممتاز جداً	ممتاز جداً
ممتاز جداً	ممتاز جداً
ممتاز	ممتاز
ممتاز جداً	ممتاز جداً
ممتاز ومفید	ممتاز ومفید
رائع جداً ويختصر الوقت	رائع جداً ويختصر الوقت
ممتاز	ممتاز
جيد جداً	جيد جداً
ممتاز جداً	ممتاز جداً
حلوووو	حلو
ممتاز	ممتاز
ممتاز للغاية	ممتاز للغاية

Before Cleaning

37,956

After Cleaning + Dedup

34,441

Removals

Dropped:
2,985

Duplicates:
310

English — Before (User + CLEAN) | pair count: 3

English — After (dedup → 1)

User Name	Clean Content	User Name	Clean Content
Shahid Ali	good	Shahid Ali	good
Shahid Ali	good		
Shahid Ali	good		

Sample: Output (Before vs Cleaned)

Before_Content	Cleaned_Content
good	good
thank you so much for helping us inshallah 🙏	thank so much helping us inshallah 🙏
good	good
So good	so good
Excellent	excellent
good	good
google play store very good app internet work very fast	google play store very good app internet work very fast
most beautiful aap	most beautiful app
nice	nice
good	good
Very good	very good
good luck	good luck

5. Handling Missing Values

- The blank text columns were replaced with NA.
- Date columns containing blank cells were converted to text format and filled with NA.
- The dataset is now completely free of empty or missing values.

6. Manual Checking and Cleaning

- In addition to the automated treatment, a **manual verification** was performed to detect and delete or replace any remaining non-Arabic or non-English text.
- This step ensured that the dataset was **standardized, consistent, and fully ready** for further analyses.

মনা ও তশ্বিন যাই আল জামদাল হাল কেবি কেবি

3. Results

-  **Before cleaning:** 128,406 rows
-  **After cleaning:** 120,561 rows
-  **Removed:** 7,845 rows (6.1%)

Language Distribution (After):

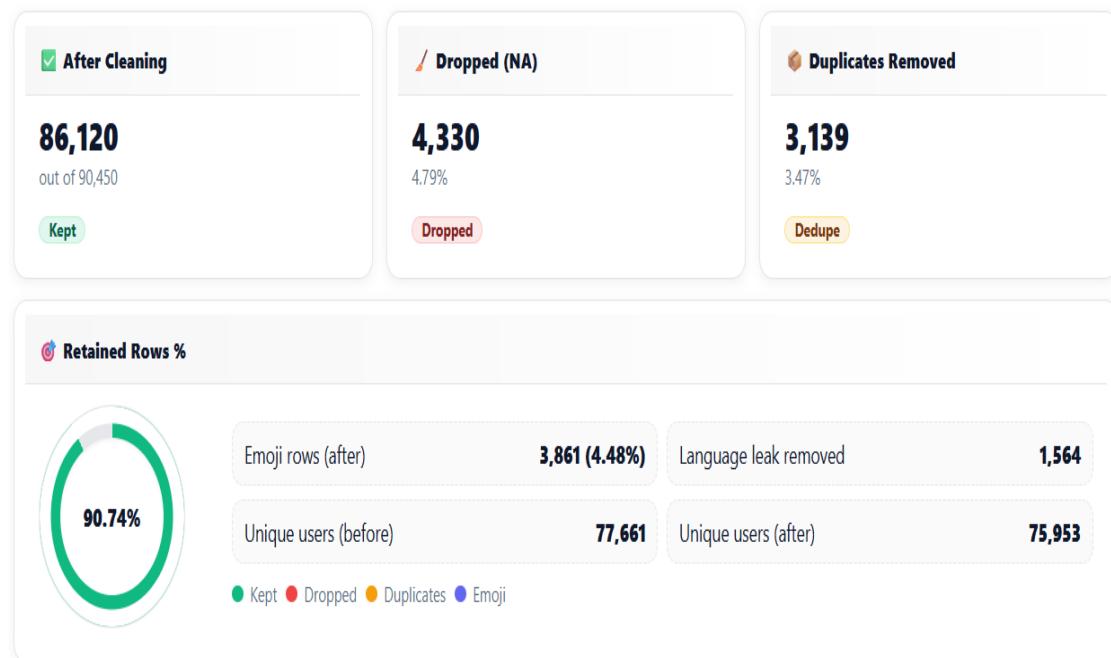
- **SA Arabic:** 90,450 → 86,120 (removed 4,330 | 4.8%)
- **GB English:** 37,956 → 34,441 (removed 3,515 | 9.3%)

Data Cleaning Dashboard

An interactive dashboard was created to summarize the results of the data cleaning process for both Arabic and English reviews:

-  **Key metrics displayed:** number of rows before and after cleaning, dropped rows (NA), duplicates removed, and percentage of retained rows.
-  **Additional insights:** unique users before and after cleaning, percentage of rows containing emojis, and language-leak removal counts.
-  **Outcome:** These dashboards provide a clear visual representation of how the dataset was standardized and prepared for analysis.

Arabic — Dashboard



Summary

Metric	Value	Extra
Total Rows (Before)	90,450	
Valid Rows (After)	86,120	95.21%
Dropped (NA)	4,330	
Duplicates Removed	3,139	
Language leak removed	1,564	
Unique Users (Before)	77,661	
Unique Users (After)	75,953	
Emoji Rows (After)	3,861	4.48%
Avg Length (Before)	21	
Avg Length (After)	21	
Median Length (After)	7	
URLs (Before)	0	
URLs (After)	0	
URLs Removed	0	

English — Dashboard

✓ After Cleaning

34,441

out of 37,956

Kept

✗ Dropped (NA)

3,515

9.26%

Dropped

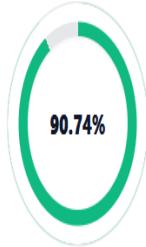
DUP Duplicates Removed

395

1.04%

Dedupe

⌚ Retained Rows %



Emoji rows (after)

2,036 (5.91%)

Language leak removed

2,458

Unique users (before)

35,399

Unique users (after)

32,509

● Kept ● Dropped ● Duplicates ● Emoji

Summary		
Metric	Value	Extra
Total Rows (Before)	37,956	
Valid Rows (After)	34,441	90.74%
Dropped (NA)	3,515	
Duplicates Removed	395	
Language leak removed	2,458	
Unique Users (Before)	35,399	
Unique Users (After)	32,509	
Emoji Rows (After)	2,036	5.91%
Avg Length (Before)	26	
Avg Length (After)	20	
Median Length (After)	9	
URLs (Before)	0	
URLs (After)	0	
URLs Removed	0	

4. Quality Assurance

- ✓ **No missing values:** empty/blank entries were standardized to NA.
- ✓ **Ratings check:** all ratings are restricted to the valid range [1–5].
- ✓ **No duplicates remain:** (*UserName* + *Clean_Content*) pairs were deduplicated — first occurrence kept; generic placeholders such as "A Google user / مستخدم Google" were ignored during duplicate detection.

🔍 Language Leakage Check:

- **SA Arabic:** 0 English-only rows (NA-like blanks excluded).
- **GB English:** 0 rows containing Arabic letters.

🔍 Arabizi Removal Check:

The dataset was additionally scanned for **Arabizi words** (transliterated Arabic written in Latin letters, e.g., *alhamdulillah*, *mashallah*, *tmam*).

📊 Sample Before/After Cleaning:

(Table below shows the first 15 affected rows before and after cleaning — demonstrating that Arabizi words were removed while keeping meaningful English words intact).

🔍 Number of rows containing Arabizi: 223

Row_Number		Before	After
0	3	thank so much helping us inshallah 🍏	thank so much helping us 🍏
1	213	alhamdulillah	NA
2	291	tmam	NA
3	324	acha	NA
4	475	mashallah good excellent 🍏	good excellent 🍏
5	558	alhamdulillah	NA
6	585	masha allah very good app	very good app
7	678	alhamdulillah 's good	s good
8	832	mashallah very very good	very very good
9	1025	very good mashallah	very good
10	1240	alhamdulillah	NA
11	1345	alhamdulillah	NA
12	1370	alhamdulillah allahuakbar	NA
13	1408	masha allah very nice	very nice
14	1453	mashallah	NA

✓ Verification Result:

The scan confirmed that **0 rows containing Arabizi words remain** — indicating that the dataset is fully clean and safe for downstream analysis.

✅ The file is clean — no Arabizi words detected.

✓ Status: CLEAN 100% for both Arabic and English.

🔍 Final Quality Check

The following tables present the final quality assurance results after cleaning and deduplication:

- **SA Arabic Data:** 86,120 rows retained, 0 duplicates, 0 English-only rows, all ratings within [1–5].
- **GB English Data:** 34,441 rows retained, 0 duplicates, 0 rows with Arabic letters, all ratings within [1–5].

✓ Final Status: Both Arabic and English datasets are confirmed **CLEAN 100%** and fully ready for analysis.

Arabic — After Cleaning										
Arabic — After Cleaning										
Rows (After)	Dup (User+Content) — After	Dup (Content) — After	Dup (Clean Col) — After	Leak Rows (after)	Missing Rating %	Rating Dist (top)	User Col	Content Col	Clean Col	Rating Col
0	86,120	0	54,910	56,449	0	0.00% {"1": 12085, "2": 1822, "3": 2837, "4": 4139, "5": 65237}	UserName	Content	Content_Arabic_Clean	Rating

English — After Cleaning										
English — After Cleaning										
Rows (After)	Dup (User+Content) — After	Dup (Content) — After	Dup (Clean Col) — After	Leak Rows (after)	Missing Rating %	Rating Dist (top)	User Col	Content Col	Clean Col	Rating Col
1	34,441	0	21,030	22,923	0	0.00% {"1": 5558, "2": 1023, "3": 1561, "4": 2858, "5": 23441}	UserName	Content	Content_English_Clean	Rating

☒ Quality Check (READ ONLY — No NA Reclassification)

Rows_Total	NonEmpty_Content	English-only (in Arabic)	Arabic-letters (in English)	User+Content Dups (rows)	UC Check	Status
------------	------------------	--------------------------	-----------------------------	--------------------------	----------	--------

Language

Arabic	86,120	86,120	0	-	0	OK ✓	CLEAN 100%
English	34,441	34,441	-	0	0	OK ✓	CLEAN 100%



5. Attachments

⌚ Sample of Cleaned Data: ⏪ (Screenshot)

📁 Final Cleaned File: Sehhaty_Database_Clean — ready for analysis ✅

💾 Available in two formats: Excel (.xlsx) and CSV (.csv)

📋 Arabic Sample — Important Columns

User Name	Rating	Content	Developer Reply	Review Year	Review Month	Review Day	Review Hour 12	Content Arabic Clean
حبيبة عبد الغفار بین الزمن	1	ممتاز جداً	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	7	ممتاز جداً
Emad ahmed	5	ممتاز جداً	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	7	ممتاز جداً
samira posaily	5	ممتاز	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	5	ممتاز
الجروح أرواح	5	ممتاز جداً جداً	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	3	ممتاز جداً جداً
زيد الكندي	5	ممتاز ومبين	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	1	ممتاز ومبين
Mohamed Ahmad Elbealy	5	رائع جداً ويختصر الوقت	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	12	رائع جداً ويختصر الوقت
الحاج جيلاوي	5	ممتاز	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	11	ممتاز
Mahmoud Saeed	4	جيد جداً	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	9	جيد جداً
علي	5	ممتاز جداً	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	6	ممتاز جداً
Zz Bb	5	طورو	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	5	طورو
ابو ياسر	5	ممتاز	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	4	ممتاز
علاء هلي	5	ممتاز للغاية	شكراً لمشاركة رأيك بتنبئي لكل يوم الصحة والعافية	2025	8	7	2	ممتاز للغاية

English Sample — Important Columns

UserName	Rating	Content	DeveloperReply	Review_Year	Review_Month	Review_Day	Review_Hour_12	Content_English_Clean
Khan Sohrab	3	good	we would like to thank for your feedback on the app and sharing your experience with us.	2025	9	3	12	good
Mine 14	5	thank you so much for helping us inshallah 🍀	we would like to thank for your feedback on the app and sharing your experience with us.	2025	9	2	9	thank so much helping us inshallah 🍀
Saleem Sulemani	5	good	we would like to thank for your feedback on the app and sharing your experience with us.	2025	9	2	9	good
Aftab Choudhary	5	So good	we would like to thank for your feedback on the app and sharing your experience with us.	2025	9	2	9	so good
Fahad Alswailim	5	Excellent	we would like to thank for your feedback on the app and sharing your experience with us.	2025	9	2	8	excellent
Ansari Aslam	5	good	we would like to thank for your feedback on the app and sharing your experience with us.	2025	9	2	4	good
Aaftab Rizvi	5	google play store very good app internet work very fast	we would like to thank for your feedback on the app and sharing your experience with us.	2025	9	2	3	google play store very good app internet work very fast

In []:

```
Total number of columns: 17
UserName, ⭐ Rating, 📄 Content, 💬 DeveloperReply,
Review_Year, 📆 Review_Month, 📅 Review_Day,
⌚ Review_Hour_12, ⚡ Review_AM_PM, 🕒 Review_Period,
🕒 Comment_Length, 📅 Reply_Year, 📆 Reply_Month, 📅 Reply_Day,
⌚ Reply_Hour_12, ⚡ Reply_AM_PM,
📝 Content_Arabic_Clean / 📝 Content_English_Clean
```

📌 The tables shown here represent **only a sample subset** of the columns, focusing on displaying the cleaned text column (Content_Arabic_Clean or Content_English_Clean). In the full dataset, there are **17 columns** to provide comprehensive coverage and enable broader analysis.

Figure 1. Summary of data cleaning results — showing **before and after cleaning**, **language distribution**, and **kept vs. removed reviews**.

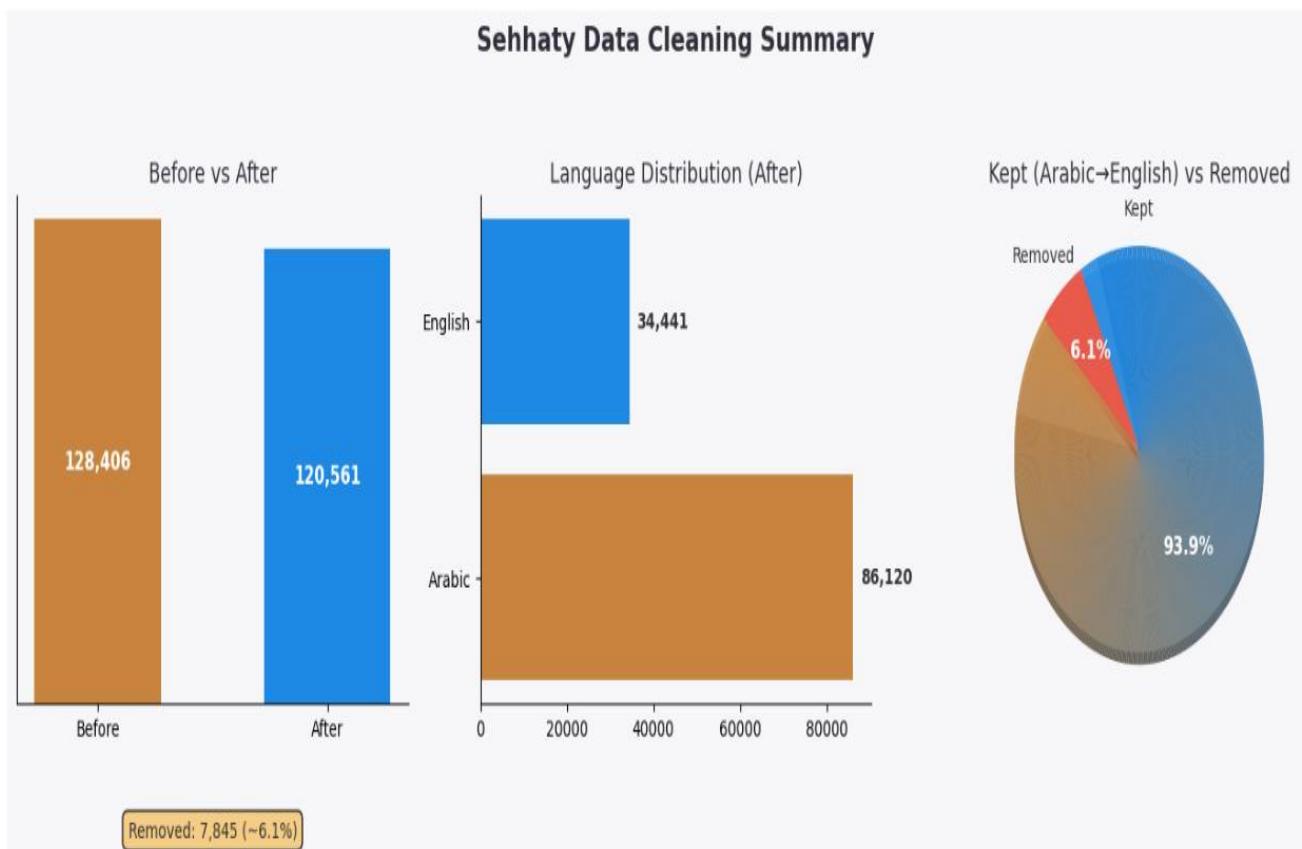
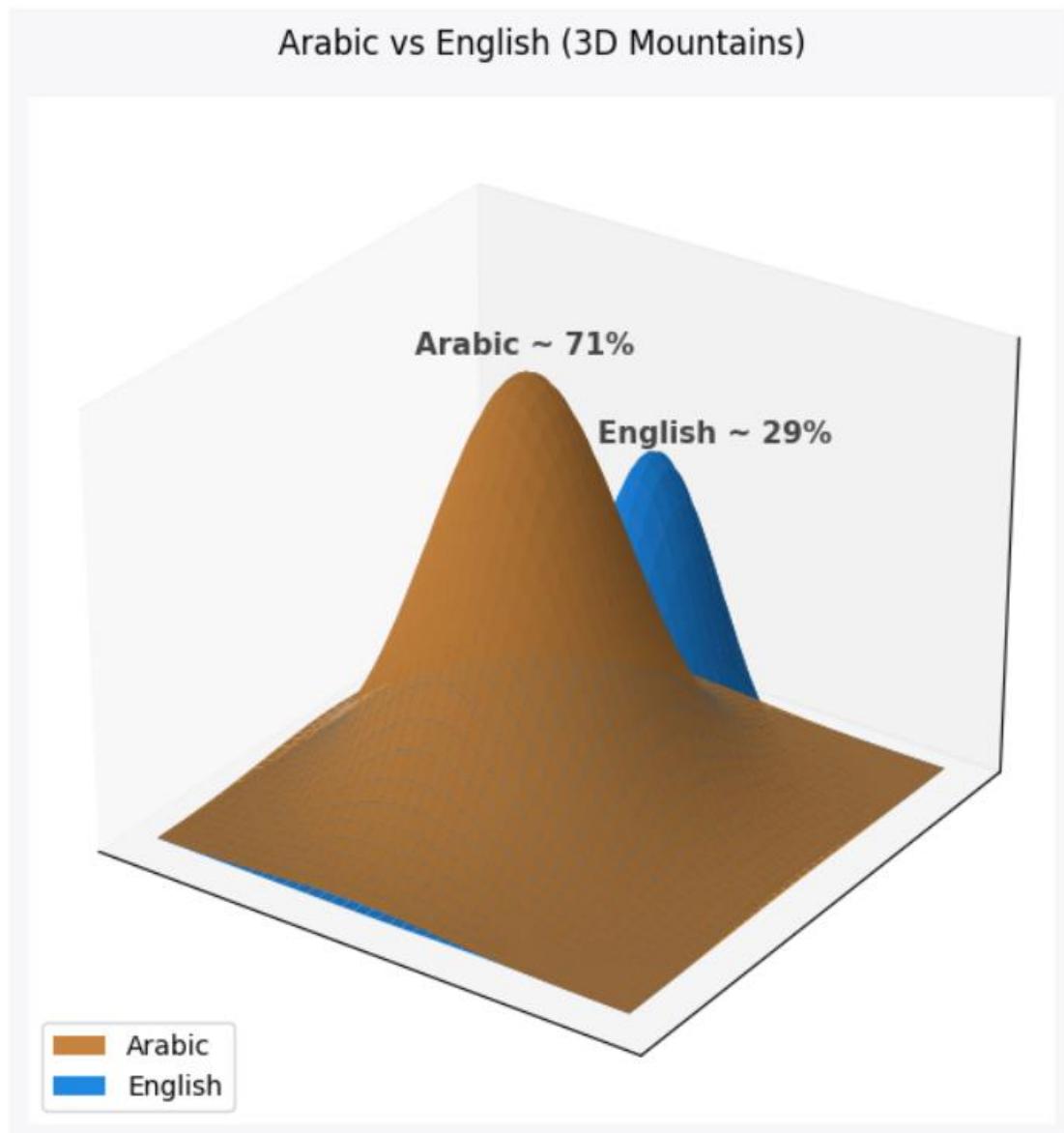


Figure 2. Distribution of the dataset by language using a 3D mountain plot. The figure shows that **Arabic reviews represent the majority (~71%)**, while **English reviews account for a smaller portion (~29%)**.



📊 === Data Cleaning Summary ===

📁 Before	┃	👑 Arabic: 90,450	┃ GB English: 37,956	┃	📦 Total: 128,406
🧹 After	┃	👑 Arabic: 86,120	┃ GB English: 34,441	┃	📦 Total: 120,561
✗ Removed	┃	▼ 7,845 (6.1%)			
✓ Kept	┃	☑ 120,561 (93.9%)			

Figure 3. Word Clouds (Before vs. After)

The word clouds highlight the most frequent terms in the reviews **before and after cleaning**. After cleaning, noise (e.g., greetings, links, repetitions) was removed, making key terms such as *التطبيق / التطبيق / good / update* more prominent and easier to analyze.



👉 This visualization clearly demonstrates the effectiveness of the cleaning process in enhancing the quality and focus of the data.

📊 The cleaning process refined the *Sehhaty* datasets, ensuring completeness, consistency, and readiness for future analysis.