

卒業論文

大規模言語モデルにおける
ジェスチャー解釈の文化的整合性分析

指導教官 村上 陽平 教授

立命館大学 情報理工学部
先端社会デザインコース 4 回生
2600220429-0

葭田 桃花

2025 年度（秋学期）卒業研究 3（CH）
令和 8 年 1 月 30 日

大規模言語モデルにおけるジェスチャー解釈の文化的整合性分析

葭田 桃花

内容梗概

人間同士のコミュニケーションにおいて、ジェスチャーは発話を補完する重要な非言語情報として機能する。しかし、その解釈は国や文化圏によって大きく異なり、同じ手の動作であっても肯定的に受け取られる場合もあれば、侮辱と捉えられる場合もある。例えば、日本やアメリカでは肯定を示す「OK」サインが、フランスでは「ゼロ」や「価値がない」を指す場合があり、南米の一部では強い侮辱表現とされる。このような文化差は対人コミュニケーションの誤解を生む大きな要因となっている。このような誤解は、人間間だけでなく、大規模言語モデル (LLM) が、多様な文化にまたがるアプリケーションに統合されるにつれて、人間と AI 間でも生じえる。そのため、ジェスチャーの意味解釈における文化差をモデルがどの程度理解しているかを評価することは、国際的な場面で AI が安全かつ適切に働くために重要な課題となっている。しかしながら、既存の LLM は主に英語圏、特にアメリカ文化圏のデータを中心に学習しており、文化的ニュアンスの違いを十分に捉えられない可能性が指摘されている。さらに、各国のジェスチャーの解釈を体系的に比較した研究は限られており、LLM の生成する解釈との差も明らかになっていない。

そこで、本研究では、各国のジェスチャー説明文と LLM が生成した説明文を同一の意味空間に埋め込み、LLM におけるジェスチャー解釈の文化整合性を分析する手法を提案する。具体的には、人手によるジェスチャーの説明および LLM の生成した説明文をベクトル化し、アメリカを基準として意味距離に基づいたクラスタリング分析を行うことで、アメリカとの文化的類似性や相違性を評価する。本手法の実現にあたり、取り組むべき課題は以下の 2 点である。

文化を考慮した意味表現

各国のジェスチャー説明文は文体・語彙・記述の詳細度が大きく異なり、そのままでは文化的特徴を保ったまま比較することが難しい。また、LLM の出力は文化中立的な抽象化が起こる場合があるため、人手データとの意味空間の差を正確に捉えるには、LLM が各国の文化固有のジェスチャー説明文を生成し、その意味表現を獲得する必要がある。

文化差に基づく文化整合性判定

LLM の生成した意味表現と人手データの意味表現で算出された意味距離では、LLM が人のジェスチャー解釈とどの程度一致したかという局所的な文

化的整合性しか判定できない。一致しなかった時に、どの国の解釈に影響を受けているのかという文化間の関係を表す大域的な文化整合性を判定する手法が求められる。

1つ目の課題に対しては、文章の正規化処理と **sentence-transformers** による意味ベクトル化を行うことで、共通の意味空間へ写像する手法を構築した。各国のジェスチャーに対する意味解釈を比較するため、まず国ごとにジェスチャーの説明文を生成した。人手データについては、各国出身のアノテータが、当該ジェスチャーについて自国で一般的に想定される意味を記述した説明文を用いた。一方、LLM データについては、当該ジェスチャーの画像と国名を入力として与え、その国における意味を説明するよう条件付けたプロンプトを用いることで、各国におけるジェスチャーの意味を説明文として生成した。さらに、プロンプトに各国の公用語を用いることで、局所的な文化整合度を向上させた。

2つ目の課題に対しては、**k-means** 法によるクラスタリングを適用することで、各国間の文化差が意味空間上でどのような構造を取るのかを確認した。LLM の大域的な文化整合性を評価するために、人手データを基準として LLM 出力のクラスタリング結果との一致度を比較し、**Precision, Recall, F1 値**, マイクロ平均, マクロ平均を用いて、クラスタ構造の再現性を定量的に測定した。また、アメリカを基準国として各国との意味距離を算出することで、LLM のジェスチャー解釈がアメリカの属するデフォルト文化圏に吸着されたのかを検証した。本研究の貢献は以下の通りである。

文化を考慮したジェスチャー解釈の意味表現

人手データから生成した意味表現と、LLM により生成した説明文から得られる意味表現との一致度を表現することで、LLM が人間の文化的解釈をどの程度再現できているのか定量的に明らかにした。加えて、各国の説明文を英語に統一した場合と、公用語を用いた場合を比較し、公用語を用いることで人手データと LLM データとの意味表現の一致度を 4%程度向上でき、公用語プロンプトにより LLM の局所的文化整合性が向上することを示した。

文化差に基づく文化整合性判定

LLM 出力と人手データのクラスタリング結果を比較、**Precision, Recall, F1 値**, マイクロ平均, マクロ平均を用いて文化差の再現度を定量的に評価した。その結果、LLM がアメリカの属するデフォルト文化へ吸着するバイアスが明らかとなり、LLM の文化理解を客観的に検証できる指標を提供した。

Evaluating Cultural Alignment in Gesture Interpretation by Large Language Models

Momoka Yoshida

Abstract

In human communication, gestures are an important form of non-verbal information that complement speech, yet their interpretations vary widely across countries and cultural regions. The same gesture may be perceived positively in some contexts and as insulting in others; for example, the “OK” sign expresses approval in Japan and the United States but can mean “zero” or “worthless” in France and is considered offensive in parts of South America. Such cultural differences can cause misunderstandings not only between humans but also between humans and AI, as large language models (LLMs) are increasingly deployed in multicultural settings. Evaluating how well LLMs understand cultural differences in gesture interpretation is therefore essential for safe and appropriate AI use. However, because existing LLMs are mainly trained on data from English-speaking countries, particularly the United States, they may fail to capture fine-grained cultural nuances. Moreover, systematic cross-country comparisons of gesture interpretations, especially between human and LLM-generated explanations, remain limited.

In this study, we propose a method to analyze the cultural alignment of gesture interpretation in LLMs by embedding human-authored and LLM-generated gesture descriptions from multiple countries into a shared semantic space. By vectorizing these descriptions and performing clustering analysis based on semantic distances with the United States as a reference, we evaluate cultural similarities and differences. Implementing this approach involves addressing two key challenges.

Meaning Representation Considering Culture

Gesture descriptions from different countries vary widely in writing style, vocabulary, and level of detail, making direct comparison difficult while preserving cultural characteristics. In addition, LLM outputs may undergo culturally neutral abstraction; therefore, to accurately capture semantic differences from human-annotated data, LLMs must generate culture-specific gesture descriptions and acquire their corresponding meaning representations.

Culture-Based Cultural Alignment Assessment

Semantic distances calculated between meaning representations generated by LLMs and those derived from human-annotated data can only evaluate *local cultural*

alignment, that is, the extent to which an LLM’s gesture interpretation matches human interpretations within a specific culture. When discrepancies occur, such measures do not reveal which other countries’ interpretations may be influencing the LLM’s output. Therefore, a method is required to assess *global cultural alignment* that captures intercultural relationships and identifies the cultural influences underlying these mismatches.

To address the first challenge, we mapped gesture descriptions into a common semantic space through text normalization and sentence-transformer-based vectorization. Gesture descriptions were generated for each country: human-annotated data consisted of descriptions written by annotators from each country based on culturally typical interpretations, while LLM-generated data were produced by inputting gesture images and country names with prompts conditioned on country-specific meanings. Using official languages in the prompts improved local cultural alignment.

To address the second challenge, we applied k-means clustering to analyze how cultural differences are structured in the semantic space. Global cultural alignment was evaluated by comparing LLM-generated clustering results with human-annotated data using Precision, Recall, F1-score, micro-average, and macro-average. In addition, semantic distances from the United States were analyzed to examine whether LLM interpretations were biased toward the default U.S.-centric cultural sphere. The contributions of this study are summarized as follows.

Meaning Representation of Gesture Interpretation Considering Culture

By measuring the correspondence between meaning representations from human-annotated data and LLM-generated descriptions, we quantitatively evaluated how well LLMs reproduce human cultural interpretations. Comparing English-standardized descriptions with official-language prompts showed that using official languages improves alignment by approximately 4%, enhancing the local cultural alignment of LLMs.

Cultural Adaptation Assessment Based on Cultural Differences

We compared the clustering results of LLM outputs and human data, and quantitatively evaluated the reproduction of cultural differences using precision, recall, F1 score, micro-average, and macro-average. As a result, we revealed a bias in which LLMs tend to gravitate towards the default culture of the United States, providing a metric to objectively verify LLMs' cultural understanding.

大規模言語モデルにおけるジェスチャー解釈の文化的整合性分析

目次

第1章	はじめに	1
第2章	関連研究	2
2.1	ジェスチャー解釈の文化差	2
2.2	仮説	3
第3章	ジェスチャー解釈の埋め込み	5
3.1	埋め込み処理	5
3.2	ジェスチャー解釈データ	5
3.2.1	人手アノテーションデータ	6
3.2.2	LLM 生成データ	6
第4章	文化整合性の評価	8
4.1	大域的評価と局所的評価の位置づけ	8
4.2	コサイン類似度による局所的文化的整合性の評価	8
4.3	クラスタリングによる大域的文化的整合性の評価	8
4.4	評価指標	9
第5章	実験	11
5.1	データセット	11
5.2	局所的文化的整合性の評価結果	13
5.2.1	英語プロンプトによる評価	13
5.2.2	公用語プロンプトによる評価	13
5.3	大域的文化的整合性の評価結果	13
5.3.1	人手データによる大域的文化的整合性の評価法の検証	13
5.3.2	英語プロンプトによる評価	15
5.3.3	公用語プロンプトによる評価	19
第6章	考察	23
6.1	局所的文化的整合性	23
6.1.1	ジェスチャーごとの分析	23
6.1.2	国ごとの分析	24

6.2 大域的文化的整合性.....	25
6.2.1 デフォルト文化圏への吸着	25
第7章 おわりに	27
謝辞	28
参考文献	29
付録	30
A.1 クラスタリング結果.....	30

第1章 はじめに

近年異なる文化圏間におけるコミュニケーションの重要性が高まる中で、言語情報だけでなく、身振りや表情といった非言語情報が果たす役割に注目が集まっている。特にジェスチャーは、言語を補完する重要な手段である一方で、文化や地域によって意味が大きく異なる場合があり、誤解や摩擦を生む要因にもなり得る。そのため、ジェスチャーの意味が文化ごとにどのように異なり、またどのような共通性を持つのかを体系的に把握することは、異文化間コミュニケーションの円滑化において重要な課題である。

これまで、ジェスチャーの文化差に関する研究は、人手による意味付与や定性的な比較を中心として行われてきた。しかし、複数の国や文化を横断的に比較する場合、意味表現のばらつきや記述の多様性により、客観的かつ定量的な比較がまだ行われていない。近年では、文埋め込み表現を用いた意味のベクトル化や、クラスタリング手法を用いた意味構造の分析を行うことで、文化差をより構造的に捉えられる可能性がある。

そこで本研究では、複数の国におけるジェスチャーの意味説明文を対象とし、意味表現をベクトル化した上でクラスタリングを行うことにより、文化差が意味空間上でどのような構造を持つのかを明らかにすることを目的とする。特に、人手アノテーションによる人手データと、LLM によって生成された LLM 生成データの両者を比較し、両者が形成するクラスタ構造の類似点および相違点を分析する。これにより、大規模言語モデルが文化的な意味差をどの程度反映しているのかについて検討する。

本論文では、まず第 2 章においてジェスチャーの文化差に関する既存研究および意味分析手法について整理し、本研究の位置づけと仮説を示す。次に第 3 章では、本研究で用いるデータセットの構成と、クラスタリング手法について述べる。第 4 章では、文化整合性を評価するための実験方法と評価指標を述べ、第 5 章において局所的小域的および大域的文化的文化整合性に関する実験結果を示す。第 6 章では得られた結果に基づき、LLM における文化的解釈傾向について考察を行い、第 7 章で本研究のまとめと今後の課題について述べる。

第2章 関連研究

本章では、非言語コミュニケーションにおけるジェスチャーの文化差に関する既存研究を整理し、特に人手による意味理解と AI による意味生成の差異に着目した研究について調べる。

2.1 ジェスチャー解釈の文化差

ジェスチャーは、発話を補完したり感情や態度を表現したりする重要な非言語行動の一つである。一方で、ジェスチャーの意味は文化や社会規範に強く依存しており、同一の身体動作であっても、国や文化圏によって肯定的、中立的、あるいは侮辱的といった意味を持つことが知られている。そのため、異文化間コミュニケーションにおいてジェスチャーの誤解は、相手に不快感や誤認識を与える要因となり得る。McNeill らは、ジェスチャーを単なる付随的な動作ではなく、思考や意味生成と密接に結びついた表現手段として位置づけ、発話とジェスチャーが一体となって意味を構成していることを示した[1]。また、Morris は世界各地のジェスチャーを収集・分類し、同一の動作が文化圏によって異なる意味を持つ事例を多数報告している[2]。これらの研究は、ジェスチャーが文化的文脈を内包する社会的信号であることを明らかにしている。さらに、Ekman らは非言語行動における文化普遍性と文化依存的要素の両立について議論し、ジェスチャーの解釈が社会規範や対人関係の影響を受けることを示した[3]。これにより、ジェスチャーの意味を単一の正解として扱うことの難しさが指摘されている。

近年では、こうした知見を踏まえ、人と AI の相互作用におけるジェスチャー理解にも関心が向けられている。Akhila Yerukola らは、文化的に不適切なジェスチャーによる AI の意図しない冒犯を防ぐことを目的とし、多文化ジェスチャーデータセットを用いて AI システムの評価を行っている[4]。その結果、AI は米国文化に基づいた解釈に偏る傾向や、不適切とされるジェスチャーを過剰に検出する傾向が確認された。これにより、AI がジェスチャーの文化的文脈を十分に考慮できていない可能性が示され、文化差を踏まえた AI の安全設計の必要性が示唆されている。また、画像とテキストを組み合わせたマルチモーダル AI によるジェスチャー理解の研究も進められているが、これらの多くは分類制度や性能評価に主眼が置かれており、文化差そのものを構造的に比較・分析する試みは限定的である。また、Cassell らは、対話エージェントにおけるジェスチャー

一や視線といった非言語行動を、意味生成や対話構造の一部として捉え、人と人工エージェントとの自然な相互作用を実現する枠組みを示している[5].しかし、これらの研究においても、ジェスチャーの文化差を定量的に比較・評価する視点は限定的である。

以上の関連研究から、ジェスチャーの意味は文化差を強く内包する一方で、AIによる意味生成においてはその差異が十分に保持されていない可能性が示唆される。本研究では、これらの課題を踏まえ、人手データと LLM 生成データを用いて、ジェスチャー意味の文化差が意味空間上でどのように表現されているのかを比較・分析することを目的とする。

2.2 仮説

本研究では、LLM が生成するジェスチャー解釈が、各国・文化圏においてどの程度文化的に適合しているかを分析することを目的とする。ジェスチャーは文化依存性が高く、人手アノテーションによって得られた意味解釈には文化圏ごとの特徴が反映されていると考えられる。一方、LLM による生成結果が、こうした文化的特徴をどの程度保持できているかについては十分明らかになっていない。そこで本研究では、人手データに基づくジェスチャーの意味をクラスタリングした結果を正解データとし、LLM によって生成されたジェスチャーの意味説明文との比較を用いて大域的整合性を評価する。また、人手データの意味説明と LLM により生成した意味説明文を直接比較し、国やジェスチャーごとに局所的整合性の評価も行う。LLM による生成文には、各国の公用語を用いたプロンプトと、英語を用いたプロンプトの 2 種類を使用し、それぞれから得られた意味をクラスタリングすることで、文化適合性を評価する。さらに、異なる性能特性を持つ LLM 間の比較を行うため、GPT-4.1mini と GPT-5.1 を用いて同一条件下でジェスチャーの意味生成を行い、モデルの違いが文化的適合性に与える影響を分析する。これを示すために 2 つの仮説を立てた。

(H1) 大域整合性及び局所的整合性の双方において、人手データを基にしたクラスタリング結果と意味説明文を正解データとした場合、公用語プロンプトによって生成された LLM データのクラスタリング結果と意味説明文を直接比較した結果は、英語プロンプトによって生成された結果よりも高い評価値を示す。

(H2) 大域的整合性及び局所的整合性の双方において、GPT-4.1mini と GPT-

5.1 を比較した場合, **GPT-5.1** によって生成されたジェスチャーの解釈は, **GPT-4.1mini** によるものよりも, 人手データに近いクラスタリング結果を示すとともに, 意味説明文の直接比較においても高い評価値を得る.

これらの仮説は, 公用語プロンプトが文化的文脈を反映した表現を引き出しやすいこと, また, より高性能なモデルである **GPT-5.1** の方が文脈理解能力や表現の精緻さに優れており, 文化差をより適切に保持したジェスチャー解釈を生成できると考えられていることに基づいている.

第3章 ジェスチャー解釈の埋め込み

本章では、LLM が生成するジェスチャー解釈の文化整合性を分析するために行った分析手法について述べる。

3.1 埋め込み処理

本研究では、各国におけるジェスチャーの解釈の違いを定量的に比較するため、テキストとして与えられた解釈文を埋め込みベクトルに変換する処理を行った。本節では、ジェスチャーの解釈文を取得してから、最終的にベクトルを得るまでの一連の処理の流れについて説明する。まず、本研究で扱うジェスチャーの解釈文は、人手アノテーションデータおよび LLM によって生成されたデータの 2 種類から構成される。これらの解釈文は、いずれも自然言語による記述であり、そのままでは定量的な比較が困難である。そこで、意味的内容を保持したまま数値表現へ変換するため、埋め込み処理を適用した。LLM 生成データに関しては、各国の文化的文脈を反映させるため、国ごとの公用語を用いたプロンプト、または英語を用いたプロンプトによってジェスチャーの意味解釈文を生成した。この段階では、解釈文の言語が国ごとに異なる状態となる。一方で、埋め込みモデルに入力するテキストの言語が混在すると、言語そのものがベクトル間の距離に影響を与える可能性がある。そのため、本研究では、意味内容の比較を主目的とし、言語の影響を抑制する必要がある。そこで、LLM によって生成された各国語の解釈文に対して、意味内容を保持したまま英語へ翻訳する処理を行った。表現形式を統一することで、埋め込み処理および類似度計算において、言語的要因ではなく意味的要因に基づいた比較が可能となるようにした。翻訳を行った英語テキストを埋め込みモデルに入力し、各解釈文を固定次元の数値ベクトルとして取得した。得られた埋め込みベクトルは、ジェスチャーの意味を意味空間上に表現したものであり、分析では、これらのベクトル間のコサイン類似度を用いて、国ごとのジェスチャー解釈の近さを評価した。

3.2 ジェスチャー解釈データ

本研究では、大規模言語モデルにおけるジェスチャー解釈の文化整合性を分析するため、人手アノテーションに基づくデータセットと、LLM によって生成したデータセットの 2 種類を用いる。両データセットは、共通のジェスチャー画

像を用いて構築しており，意味表現の比較が可能となるよう設計した．

3.2.1 人手アノテーションデータ

人手データセットには，Akhila Yerukola らにおいて作成された MC-signs データセットを用いた．このデータセットは，文化圏ごとのジェスチャー解釈を収集されたものである．MC-signs には，文化的に解釈が分かれるジェスチャー画像に対して，各国のアノテータが，自国の文化的文脈に基づいて記述した意味説明が含まれている．これらの記述は，当該ジェスチャーがその国でどのような意味やニュアンスを持って使用しているかを反映している．本研究では，この人手アノテーションデータを，ジェスチャー意味の文化的構造を反映した基準データとして扱う．分析に先立ち，文字化けが生じているものや意味が不明慮な記述，及び分析に適さないデータを除外するクリーニング処理を行い，クラスタリング分析に利用可能な形へと整形した．

3.2.2 LLM 生成データ

LLM 生成データは，人手アノテーションデータと同一のジェスチャー画像を入力として，大規模言語モデルに対して意味解釈を生成させることで構築した．入力としては，ジェスチャー画像に加え，「このジェスチャーはどのような意味を持つか」という質問文を与えた．プロンプト設計においては，まず言語条件の違いが出力に与える影響を分析するため，各国の公用語を用いたプロンプトと，英語を用いたプロンプトの 2 種類を作成した．これにより，文化的文脈を反映した表現が，使用言語によってどの程度変化するかを比較した．また，出力内容の違いを検証するため，「考えられる全ての意味を出力する」プロンプトと，「代表的な意味を出力する」プロンプトの 2 種類を用意した．これにより，多義的な解釈を保持した場合と，代表的な解釈に集約した場合とで，クラスタリング結果にどのような差が生じるか分析した．以下に，例として各国の公用語を用いて，「考えられる全ての意味を出力」した際のプロンプト内容を記載した（図 1 参照）．さらに，LLM の出力形式にばらつきが生じることを防ぐため，出力形式に関する制約をプロンプトに明示的に追加した．生成結果では，意味を単語のみで列挙する出力例が見られたが，このような形式は意味表現の粒度が不均一になり，クラスタリング分析に適さないという課題があった．そこで，各意味を必ず箇条書きの文章として記述するようプロンプトに明記した．加えて，各文章の長さを指定し，意味表現の粒度と構造を統一した．この制約により，生成された意味記述をベクトル化する際のばらつきを抑え，人手アノテーションデータとの比較を

より適切に行えるようにした.

以下の質問に答えるときは、次の形式で書いてください：

- ・文章
- ・文章

※ 箇条書きの各項目は「一行の見出し」ではなく、必ず「文章として説明」してください。

※ 文章は1~20語程度

※ 何項目でも良いので、箇条書きで詳しくまとめてください。

※ 回答は、【 】の中に記載されている本文で使用されている言語と同一の言語で記述してください。

※ すべての【 】の内容に回答してください。

【 What does this gesture mean in Taiwan? Please explain its meaning in detail on a national level. Provide both positive and negative interpretations, and explain each specifically. Do not use abstract terms like 'aggressive' or 'negative'.】

図 1 各国の公用語を用いて考えられる全ての意味を出力した際のプロンプト

第4章 文化整合性の評価

4.1 局所的評価と大域的評価の位置づけ

人手データと、LLM 生成データとの一致度は、国やジェスチャーごとに異なる可能性があり、特定の国や特定のジェスチャーにおいて解釈のズレが生じることが考えられる。このような差異を詳細に把握するためには、国単位・ジェスチャー単位での局所的な文化的適合性の評価が重要である。そこで本章では、文化ベクトル間のコサイン類似度に基づく分析を用い、LLM が各国におけるジェスチャー解釈がどの程度正確に再現できているかを局所的に評価する。

一方で、局所的評価は個別の傾向を捉えるのに適している反面、文化圏全体における意味構造の共通性や全体的な傾向を把握することは難しい。そのため、補完的な評価とし、人手データと LLM 生成データとのクラスタリング結果の一致度を文化領域単位で大域的に評価する。これにより、LLM が文化圏ごとの意味構造をどの程度再現できているかを把握することが可能となる。

4.2 コサイン類似度による局所的文化整合性の評価

本節では、ジェスチャーおよび国ごとの文化的適合性を評価するため、コサイン類似度を用いた分析を行う。具体的には、あるジェスチャーのある国における人手アノテーション文のベクトルと、LLM によって生成された解釈文ベクトルとの間のコサイン類似度を算出した。この類似度は、LLM が当該国におけるジェスチャー解釈をどの程度人手データに近い形で生成できているかを示す指標である。本研究では、このコサイン類似度を国ごとに算出し、各ジェスチャーについて平均値および分散を求めた。平均値は、そのジェスチャーに対する LLM の全体的な文化的適合性を示す指標として解釈できる。一方、分散は、国ごとの適合性のばらつきを表しており、分散が大きいジェスチャーほど、国によって LLM の解釈精度が大きく異なっていることを意味する。さらに、全ジェスチャーに対して算出したコサイン類似度の平均値を用い、英語プロンプトと公用語プロンプトの比較を行う。これにより、公用語を用いることで、LLM がどの程度文化的文脈を反映した解釈を生成できるようになったかを定量的に評価する。

4.3 クラスタリングによる大域的文化的整合性の評価

本研究では、ジェスチャー意味記述の意味的類似性を基に、各国の解釈がどの

程度近いかを分析するため、分散表現を用いたクラスタリング手法を採用した。まず、各国におけるジェスチャー意味の記述文を文埋め込みに変換し、意味空間上での距離を算出した。文埋め込みの生成には、Sentence-BERT 系モデルである all-mpnet-base-v2¹を用いた。このモデルにより、各意味記述文を固定長のベクトルに変換し、mean pooling を用いて文全体の表現を取得した。得られたベクトルは、コサイン類似度に基づく比較を行うため、正規化を施した。次に、基準としてアメリカを設定し、各国の意味記述とアメリカの意味記述との平均コサイン類似度を算出した。具体的には、アメリカの意味記述群と各国の意味記述群の間で類似度行列を計算し、その平均値を求めた。これを距離として扱うため、1 から平均類似度を減算した値をアメリカとの意味的距離として定義した。この距離値を特徴量として、各国を対象に k-means 法によるクラスタリングを行った。本研究では、意味解釈が「same (同じ)」「uncertain (どちらでもない)」「different (異なる)」という三つの状態に分類するために、クラスタ数は 3 と設定した。クラスタリング後、各クラスタに属する国の平均距離を算出し、距離が最も小さいクラスタを「same」、中間を「uncertain」、最も大きいクラスタを「different」と対応付けた。さらに、意味空間上での分布を直感的に把握するため、各国の意味記述ベクトルを平均化した後、主成分分析 (PCA) を用いて二次元に次元削減し、可視化を行った。可視化では、アメリカを基準として強調表示し、クラスタ分類結果に応じて色分けを行うことで、各国の相対的な位置関係を示した。以上の手法により、人手データおよび LLM 生成データにおけるジェスチャー意味解釈の類似性を定量的に比較し、文化的適合性の分析を行った。

4.4 評価指標

本研究では、LLM 生成データが、人手データに基づく文化的構造をどの程度再現できているかを評価するため、クラスタリング結果の一致度に基づく評価を行った。まず、人手アノテーションデータに基づいて得られたクラスタリング結果を正解データとして、LLM 生成データから得られた、クラスタリング結果との比較を行った。評価単位は国とし、各国が属するクラスタが人手データと LLM 生成データで一致している場合を「same」、異なるクラスタに分類された場合を「different」として扱った。この一致・不一致の結果を二値データとして

¹ All-mpnet-base-v2 は sentence-BERT 系列の事前学習済みモデルであり、Microsoft Research によって提案された MPNet に基づく文埋め込みモデルである。

整理し、適合率 (precision), 再現率 (recall), 及び F 値を用いて評価を行った。これらの指標は、クラスタリング結果が正解データとどの程度一致しているかを示すために用いた。評価指標の算出にあたっては、macro 平均および micro 平均の二種類を用いた。macro 平均は、各クラスタ (same, uncertain, different) ごとに適合率・再現率・F 値を算出し、それらを単純平均する方法であり、クラスタ間の影響を均等に評価できる。一方、micro 平均は、全クラスタの結果をまとめて算出する方法であり、全体としての分類性能を評価する指標である。本研究では、これらの 2 種類の平均値を用いることで、特定のクラスタへの偏りの有無と、全体的な文化的適合性の両方を評価した。また、公用語プロンプトと英語プロンプトの比較、および GPT-4.1mini と GPT-5.1 のモデル間比較において、同一の評価手法を適用することで、条件間の違いを定量的に比較した。

第5章 実験

5.1 データセット

本章では，ジェスチャーの意味解釈に関するクラスタリング結果を示すにあたり，本研究で利用したデータセットの概要を説明する．本研究では，人手によって収集された人手データと，LLMによって生成されたLLM生成データの2種類を用いて分析を行った．これらのデータを同一の前処理およびクラスタリング手法で比較することで，人手データに基づく意味構造と，LLMが生成する意味構造との違いを明らかにすることを目的としている．

LLM生成データは，人手データと同一の22種類のジェスチャーを対象とし，各国における意味解釈を大規模言語モデルにより生成したものである．本研究では，モデルの違い，使用言語の違い，および出力形式の違いがクラスタリング結果における影響を検証するため，以下の8条件でデータを生成した．(図2参照)

以上の条件により，表1に示したLLM生成データ8条件で，22個のジェスチャーの計176通りのクラスタリング結果を得た．条件ごとのクラスタリング結果を定量的指標として集約し，人手データとの比較を行う．特に，LLM生成データがどの程度人手データの意味構造を再現できているか，また条件の違いによってどのような変化が生じるかに着目して分析を進める．

・GPT5.1	公用語プロンプト	全ての意味を出力
・GPT5.1	公用語プロンプト	代表的な意味を出力
・GPT4.1mini	公用語プロンプト	全ての意味を出力
・GPT4.1mini	公用語プロンプト	代表的な意味を出力
・GPT5.1	英語プロンプト	全ての意味を出力
・GPT5.1	英語プロンプト	代表的な意味を出力
・GPT4.1mini	英語プロンプト	全ての意味を出力
・GPT4.1mini	英語プロンプト	代表的な意味を出力

表1 LLM生成データ8条件

表 2 ジェスチャーの種類

OK	親指と人差し指で円を作り，残りの 3 本の指を伸ばした手の形
Chin flick	手の甲を顎の下に当て，前方に甲を見せる動作
Curled finger	人差し指を曲げ，手前に引き寄せるように動かす動作
Five fathers	手の平を相手に向け，指を開いた状態から特定のリズムや向きで示す動作
Forearm Jerk	片腕を曲げ，反対の手で前腕を叩く動作
Horns	人差し指と小指を立て，他の指を折り曲げた手の形
Index finger pointing	人差し指を伸ばし，特定の方向や対象を指し示す動作
L	親指と人差し指を直角に広げ，アルファベットの「L」の形を作る手の形
Left hand	左手のみを用いて，物を持ったりする動作
Middle finger	中指を立て他の指を折り曲げた手の形
Open palm with fingers spread	手のひらを前に向け，5 本の指を大きく広げた状態
Pinched fingers	親指と他の指先をまとめてつまんだ形
Quenelle	片腕を伸ばし，もう一方の手でその腕の二の腕付近に触れる姿勢
Serbian Salute	3 本の指（親指・人差し指・中指）を立て，残りの指を折り曲げた手の形
Shocker	人差し指と中指，小指を伸ばし，薬指を折り，親指で抑えた手の形
Show sole of shoe or feet	足の裏を相手に向けるように足を上げ座る姿勢
Snap fingers	親指と中指を擦り合わせて音を鳴らす動作
Three finger salute	親指・人差し指・中指の 3 本を立て，他の指を折り曲げた手の形
Thumbs up	親指を立て，他の指を握った手の形
Touching someone head	他者の頭部に手を触れる動作
V sign	人差し指と中指を開いて V 字を作る手の形
Wanker	手を軽く握り，上下に反復運動させる動作

今回比較したジェスチャーの種類は表 2 に示す 22 個である。

5.2 局所的文化的整合性の評価結果

本節では、各ジェスチャーについて国ごとの解釈に着目し、LLM が生成した解釈文と人手アノテーションデータとの意味的整合性を局所的に評価した結果について述べる。ここで、特定のジェスチャーに対して、特定の国における解釈が、人手アノテーションに基づく意味構造とどの程度一致しているかを指す。評価には、人手データによる文ベクトルと、LLM が生成したデータの文ベクトルとの間のコサイン類似度を用いた。各ジェスチャーについて国ごとに類似度を算出し、その平均および分散を求めることで、文化的整合性の傾向とばらつきを分析した。

5.2.1 英語プロンプトによる評価

英語プロンプト条件では、全ての国に対して英語でジェスチャー解釈を生成した LLM 出力を用いて評価を行った。その結果、アメリカおよび英語圏の国では、人手データとの類似度が比較的高くなる一方、非英語圏の国では、類似度が低下する傾向が確認された。また、ジェスチャーごとに類似度の分散を比較すると、分散の小さいジェスチャーと大きいジェスチャーが存在した。分散が大きいジェスチャーでは、国によって人手データとの一致にばらつきが見られ、英語プロンプトでは各国固有の文化的解釈が十分に反映されていない可能性が示された。

5.2.2 公用語プロンプトによる評価

公用語プロンプト条件では、各国の公用語を用いて LLM にジェスチャー解釈を生成させ、同様の手法で人手データとの意味的類似度を評価した。その結果、英語プロンプト条件と比較して、多くのジェスチャーにおいて国ごとの類似度が向上し、分散が低下する傾向が確認された。特に、文化的背景によって意味が変化しやすいジェスチャーにおいて、公用語プロンプト条件では、人手アノテーションに基づく解釈との一致度が高まる例が多く見られた。これらの結果から、公用語プロンプトは、各国固有の文化的文脈をより反映した解釈文の生成に寄与し、局所的文化的整合性の向上につながることを示唆される。

5.3 大域的文化的整合性の評価結果

本節では、大域的整合性の評価結果について述べる。

5.3.1 人手データによる大域的文化的整合性の評価法の検証

人手データのクラスタリング結果は、設定した Ground Truth に基づいて評価

した。評価では、各ジェスチャーについてアメリカの人手データを基準とし、クラスタリング結果において各国の解釈がアメリカと同一クラスタに属する場合を「同じ」、異なるクラスタに属する場合を「異なる」、いずれとも明確に判断できない場合を「どちらでもない」と分類した。

表 3、表 4 は、人手データを評価した結果である。適合率、再現率、F 値はいずれのクラスにおいても高い値を示している。

付録 A-91、A-92 に掲載している **Middle Finger** や **Forearm Jerk** などの強い侮辱を持つジェスチャーでは、多くの国がアメリカと同一のクラスタに分類される傾向が確認された。これらのジェスチャーは、文化圏を超えて意味が共有されやすく、文化差が小さいジェスチャーであるといえる。一方で、**OK** ジェスチャーや **Thumbs up** などの日常かつ多義的なジェスチャーでは、国ごとに異なるクラスタが形成される例が多く見られた。これらのジェスチャーは使用頻度が高い反面、文化的背景や文脈による意味の差異があり、文化差が顕著に表れることが示された。さらに、**Quenelle** や **Serbian salute** など、宗教的・政治的背景と強く結びついたジェスチャーでは、特定の文化圏にのみ明確な意味を持つ局所的な構造が確認された。これらは、多文化圏では意味が共有されにくく、クラスタリング結果においても分散が大きい傾向を示した。

表 3 人手データを評価した値

人手データ	適合率	再現率	F 値
同じ	0.911	0.971	0.930
異なる	0.984	1.0	0.990
どちらでもない	0.959	0.926	0.925

表 4 人手データを評価した値

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.97	0.962	0.966	0.967	0.969	0.959

5.3.2 英語プロンプトによる評価

ここでは、表 5、表 6 に記載している GPT4.1mini で英語プロンプト・全ての意味を出力した結果について明記する。出力結果として、各ジェスチャーについて、「Positively, it can ...」「Negatively, it may ...」のように肯定的な解釈と否定的な解釈を併用し、「depending on the context」「might be interpreted as」といった文脈的依存性を示す表現が多く用いられていた。このため、LLM の出力は単一の意味に収束せず、複数の解釈を包含する傾向が強かった。その結果、「どちらでもない」クラスの F 値が最も高く、文化的意味関係が中間的なケースの識別には一定の有効性が確認された。一方で、「同じ」及び「異なる」クラスの F 値はそれぞれ 0.379, 0.392 に留まり、意味が完全に一致、または明確に異なるケースの判定は十分とは言えなかった。また、本プロンプトは英語で記述しており、実際の出力内容においても、英語圏、特にアメリカで一般的な語彙・評価軸が多く用いられていた。多くの国の説明において、アメリカと類似した社会的状況や対人距離感を前提とする記述がみられたことから、各国固有の文化的背景よりも、英語圏の意味枠組みが反映されやすい構造であったと考えられる。この影響により、本結果では、文化差が連続的・曖昧なものとして表現されやすく、「どちらでもない」への分類が増加した一方で、アメリカと明確に異なる意味を持つジェスチャーであっても、その差異が弱められ、「異なる」クラスへの判定精度が伸びにくかった可能性がある。すなわち、アメリカバイアスがかかりやすく、そのことが評価指標にも反映されたと解釈できる。以上のことから、文化的多様性の把握には有効である一方、英語圏、特にアメリカを基準とした意味解釈に引き寄せられやすく、文化固有性の厳密な判定には課題が残ることが示唆された。

表 5 GPT-4.1mini 英語プロンプト 全ての意味を出力した結果

英語 全ての意味	適合率	再現率	F 値
同じ	0.385	0.374	0.379
異なる	0.405	0.380	0.392
どちらでもない	0.48	0.387	0.428

表 6 GPT-4.1mini 英語プロンプト 全ての意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.405	0.306	0.348	0.575	0.524	0.546

表 7 GPT-4.1mini 英語プロンプト 代表的な意味を出力した結果

英語 代表的意味	適合率	再現率	F 値
同じ	0.229	0.3	0.259
異なる	0.247	0.306	0.273
どちらでもない	0.419	0.335	0.372

表 8 GPT-4.1mini 英語プロンプト 代表的な意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.330	0.327	0.328	0.298	0.313	0.301

ここでは、表 7、表 8 に記載している GPT4.1mini で英語プロンプト・代表的な意味を出力した結果について明記する。実際のプロンプトでは、肯定的・否定的な解釈の併記や文脈依存性への言及はほとんど見られなかった。そのため、各国の意味が単純化され、人手解釈に含まれる文化的背景や使用状況の差異が十分に表現されにくい傾向があった。評価結果を見ると、「どちらでもない」クラスの F 値が最も高く、一方で「同じ」「異なる」の F 値はいずれも低い位置に留まった。これは、多くの国で類似した短文表現が生成された結果、意味が完全に一致するとも、明確に異なるとも判定しにくいケースが増加したためと考えられる。さらに、本プロンプトは、米国における一般的な意味が基準として提示される構造であったため、他国の意味も米国的な意味枠組みに引き寄せられやすく、アメリカのバイアスが生じやすいと考えられる。

表 9 GPT-5. 1 英語プロンプト 全ての意味を出力した結果

英語 全ての意味	適合率	再現率	F 値
同じ	0.276	0.271	0.273
異なる	0.382	0.385	0.383
どちらでもない	0.371	0.409	0.389

表 10 GPT-5.1 英語プロンプト 全ての意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.367	0.366	0.367	0.531	0.536	0.534

ここでは、表 9、表 10 に記載している GPT5.1 で英語プロンプト・全ての意味を出力した結果について明記する。実際の記述では、「People may use it to...」「It may insult...」といった表現を用い、肯定的な用法、否定的・侮辱的な用法、文脈的な解釈を同時に提示している点が特徴的である。例えば、OK gesture では、「料理が美味しい」「正しい答えを確認する」といった肯定的な意味に加え、否定的な意味まで具体的に列挙されていた。また、Chin flick gesture や Forearm Jerk gesture では、「拒否」「無関心」「強い否定」などの意味が、軽い冗談から深刻な対立場面までの段階的に記述されており、意味の幅と強度が明確に表現されていた。さらに、Serbian Salute や Quenelle gesture のような政治的・社会的背景を持つジェスチャーについても、「ポップカルチャー的解釈」「抗議・挑発としての解釈」「誤解や反感を招く可能性」が併記されていた。このように意味を単一に固定せず、文化的・社会的文脈を含めた結果が反映されている。「どちらでもない」クラスの F 値は 0.389 と高く、多数の意味を適切に捉えられていることが示された。また、「異なる」クラスの適合率、再現率、F 値は 0.383 程度で一致しており、文化的背景や使用文脈が大きく異なるジェスチャーについては、明確な差異を捉えられていることが分かる。一方で、「同じ」クラスの F 値は 0.273 に留まった。これは、本プロンプトが意図的に多様な意味や否定的解釈を含めているため、たとえ複数国で共通する使用法が存在していても、意味ベクトル上では完全一致よりも「部分的に共通するが完全に同一ではない」関係として表現された少なかったことが要因と考えられる。Micro F 値が 0.367、Macro F 値が 0.534 と、これまでの設定と比べて全体性能が大きく向上している点は、多義的な意味記述が人手解釈の幅と整合しやすかったことを示唆している。すなわち、本プロンプトは、文化的意味の多層性や解釈の揺らぎを含めて扱うことで、単純化された代表的なプロンプトよりも、ジェスチャーの文化差をより適切に捉えられたと評価できる。

表 11 GPT-5.1 英語プロンプト 代表的な意味を出力した結果

英語 代表的意味	適合率	再現率	F 値
同じ	0.346	0.422	0.380
異なる	0.448	0.416	0.431
どちらでもない	0.578	0.453	0.508

表 12 GPT-5.1 英語プロンプト 代表的な意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.432	0.569	0.491	0.673	0.533	0.595

ここでは、表 11、表 12 のように、GPT-5.1 を用いて、英語プロンプトで代表的な意味を出力した結果を記載する。評価結果として、「同じ」「異なる」「どちらでもない」の3分類における F 値はそれぞれ 0.380, 0.431, 0.508 であった。特に「どちらでもない」クラスでは適合率が 0.578 と比較的高い値を示し、LLM が人手解釈と完全に一致又は完全に不一致とは言えない中間的な意味を多く生成していることが分かった。全体指標をみると、Micro F 値は 0.491、Macro F 値は 0.595 であった。Macro F 値が Micro F 値を上回っていることから、特定の分類に偏るのではなく、各クラスに対して比較的均等に予測が行われていることが分かる。一方で、Micro 再現率が 0.569 と高めであるのに対し、Micro 適合率が 0.432 に留まっていることから、正確性がやや低下している傾向が確認された。この結果は、代表的な意味に限定するプロンプト設計の影響を強く反映していると考えられる。具体的には、ジェスチャーが文化によって大きく異なる意味を持つ場合でも、その差異が十分に表現されず、抽象化された意味に収束しやすい。例えば、侮辱・否定といった否定的ジェスチャーでは、各国固有の社会的・歴史的な脈が省略され、「rude」「insulting」といった汎用的表現にまとめられる傾向が見られた。その結果、人手データが持つ文化固有の解釈との差が拡大し、「同じ」とは判定されにくくなったと考えられる。これにより、GPT-5.1 を用いて英語プロンプトで代表的な意味を出力した結果では、ジェスチャーの大まかな意味傾向を把握するには有効である一方、文化差や多義性を十分に反映することは難しく、文化整合性の観点では限界があることが示された。

5.3.3 公用語プロンプトによる評価

ここでは、表 13、表 14 のように、GPT-4.1mini を用いて、公用語プロンプトで全ての意味を出力した結果を記載する。人手データに基づくクラスタリング結果と比較して、クラスタの分離が弱く、中間的な位置に結果が集中する傾向が確認された。特に、「どちらでもない」に分類される事例が多く、曖昧な位置関係が形成される例が見られた。この傾向は、各国のジェスチャーに対して、文脈や状況に応じた複数の解釈可能性を同時に保持した表現を生成していることを反映していると考えられる。また、公用語プロンプトは、英語プロンプトと比較して、アメリカ中心の解釈に一樣に収束する傾向を緩和し、各国固有の文化的文脈を含んだ出力を生成するという点で一定の効果を示した。しかし、全ての意味を出力させる条件により、クラスタリング上ではアメリカの位置と比較し、比較的離れやすい傾向が生じた。その要因として、多様な意味が生成された結果、アメリカと同一の意味とみなされる事例が減少したことが考えられる。この傾向は、クラスタリング結果における分離の弱さとして現れていると解釈できる。以上により、今回の条件では、文化的多様性を反映した意味を生成できていたと評価できる。

表 13 GPT-4.1mini 公用語プロンプト 全ての意味を出力した結果

公用語 全ての意味	適合率	再現率	F 値
同じ	0.446	0.431	0.438
異なる	0.320	0.351	0.334
どちらでもない	0.503	0.454	0.477

表 14 GPT-4.1mini 公用語プロンプト 全ての意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.418	0.411	0.415	0.577	0.594	0.496

表 15 GPT-4.1mini 公用語プロンプト 代表的意味を出力した結果

公用語 代表的意味	適合率	再現率	F 値
同じ	0.395	0.392	0.393
異なる	0.312	0.290	0.3
どちらでもない	0.585	0.482	0.528

表 16 GPT-4.1mini 公用語プロンプト 代表的意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.405	0.396	0.400	0.569	0.573	0.570

ここでは、表 15、表 16 のように、GPT-4.1mini を用いて、公用語プロンプトで代表的な意味を出力した結果を記載する。多くのジェスチャーにおいて共通する一般的な意味が生成される傾向が確認された。一方で、人手データにおいて顕著であった文化的対立や地域固有の意味は弱まり、解釈が単純化される傾向が見られた。特定の国では、否定的な意味として捉えられることが多いジェスチャーであっても、LLM 生成データでは、アメリカにおいて一般的に用いられている代表的な意味が生成される傾向が確認された。このことから、代表的な意味を出力した結果は、意味の代表制を捉える点では有効であるが、文化的多義性や局所性を十分に反映するには限界があるといえる。

ここでは、表 17、表 18 のように、GPT-5.1 を用いて公用語プロンプトで、全ての意味を出力した結果を明記する。対象国の公用語を用いて、複数の意味や使

表 17 GPT-5.1 公用語プロンプト 全ての意味を出力した結果

公用語 全ての意味	適合率	再現率	F 値
同じ	0.343	0.273	0.304
異なる	0.427	0.427	0.427
どちらでもない	0.463	0.453	0.458

表 18 GPT-5. 1 公用語プロンプト 全ての意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.641	0.635	0.638	0.577	0.566	0.571

用場面が結果として明記されていた。記述内容の例としては、文脈と相手との関係によって評価が変化するケースが併記されており、意味を単一に固定せず、文化的に許容される範囲と問題となる範囲の両方を含めて説明する構造となっている点が特徴的である。例えば、日常的には承諾や理解を示す行為として用いられている場合と、対人関係によっては、相手を軽視、侮辱する意味を持つ場合が、同一の国内でも並列して記述されていた。また、他国との比較においても、完全に同一の意味として扱われるのではなく、「似た使われ方をするが、否定的に受け取られる度合いが異なる」「特定の状況では不快感を与える」といった差異のニュアンスが含まれていた。これは、評価結果にも明確に表れている。まず、「どちらでもない」クラスの F 値は **0.458** と最も高く、次いで「異なる」クラスが **0.427**、「同じ」クラスが **0.304** であった。これは、公用語による詳細な意味記述によって、国間で意味が部分的に重なるが完全に一致しないケースが多く表現され、中間的な意味を捉えやすくなったことを示している。一方で、複数の意味が併記されているため、意味が完全に一致すると判定できるケースは相対的に少なくなり、「同じ」クラスの F 値が低くなったと考えられる。また、**Micro F** 値は **0.638** と高く、**Macro F** 値も **0.571** と高い値を示している。この差は、クラスごとに見ると一定の識別性能が確保されており、全体として分散が低減したと考えた。すなわち、本プロンプトは各クラスの特徴を均等に捉える力はあるが、意味の多層性ゆえに単一クラスへ強く収束しにくい構造であったと解釈できる。以上より、公用語プロンプトは英語プロンプトに比べ、対象国固有の表現や文化的前提を反映しやすく、米国中心の意味枠組みに引き寄せられる傾向を緩和できている一方で、意味を多面的に記述する設計により、「同じ」「異なる」を明確に二分する判定は難しくなる傾向が確認された。この結果は、公用語プロンプトが文化的意味の幅や曖昧さを捉える評価には有効であることを示唆している。

表 19 GPT-5.1 公用語プロンプト 代表的意味を出力した結果

公用語 代表的意味	適合率	再現率	F 値
同じ	0.326	0.412	0.364
異なる	0.412	0.452	0.431
どちらでもない	0.465	0.416	0.439

表 20 GPT-5.1 公用語プロンプト 代表的意味を出力した結果

Micro 適合率	Micro 再現率	Micro F 値	Macro 適合率	Macro 再現率	Macro F 値
0.436	0.583	0.499	0.579	0.593	0.586

ここでは、表 19、表 20 のように、GPT-5.1 を用いて公用語プロンプトで代表的な意味を出力した結果について明記する。「異なる」および「どちらでもない」クラスの F 値がそれぞれ 0.431、0.439 と比較的高く、文化的に明確な差異を持つジェスチャーや、部分的に意味が重なるケースを適切に識別できている。また Micro 再現率が 0.583 と高く、Micro F 値も 0.499 に達していることから、代表的意味に絞った簡潔な記述によって、意味関係の大まかな傾向を取りこぼしにくいことが確認された。さらに、Macro F 値が 0.586 と高水準であることから、特定のクラスに偏ることなく、各クラスをバランス良く判定できている。一方で、意味を単一に集約する設計上、文脈的依存性や多義性は十分に反映されず、意味が完全に一致する「同じ」クラスの判定精度は相対的に低い結果となった。これは、本手法が文化的理解の典型像を捉える事には有効であるが、細かなニュアンスの差や使用条件を含めた精緻な比較には限界があることを示している。以上より、公用語・代表的な意味を生成させるプロンプトでは、国ごとの典型的なジェスチャー理解を比較する目的に適しており、アメリカ中心の解釈に依存しない形で文化差を把握できる一方で、文化的意味の幅や曖昧性を評価する場合には、全意味を考慮するプロンプトとの併用が有効であると考えられる。

第6章 考察

6.1 局所的文化的整合性

本節では、ジェスチャーの意味解釈における局所的文化的整合性について文性を行う。ここで局所的文化的整合性とは、特定のジェスチャーあるいは特定の国に限定した場合に、LLM 出力が人手解釈とどの程度整合しているかを指す。本分析では、表 22 に示す 4 条件を対象とし、22 種類のジェスチャーを用いて比較を行う。

6.1.1 ジェスチャーごとの分析

本節では、文化整合しやすいジェスチャーと文化整合しにくいジェスチャーの特徴について考察する。本研究では、人手解釈文と LLM 出力文のコサイン類似度の平均値が高く、かつ分散が小さいジェスチャーを文化整合しやすい、一方で、平均値が低い、もしくは分散が大きいジェスチャーを文化整合しにくいと判断できる。文化整合しやすいジェスチャーは、多くの国において意味解釈が比較的一致しており、人手解釈と LLM 出力の間で安定した対応関係が確認された。これらのジェスチャーは、承認や合図といった基本的意味、あるいは強い侮辱表現のように意味の方向性が明確で、文化間で共有されやすい特徴を持つものが多い。具体例として、付録 A-92 のような **Middle Finger gesture** では、平均コサイン類似度が 0.60～0.70 台と高く、分散も 0.006～0.008 程度に留まっていた。これは、本ジェスチャーが多くの文化圏において強い侮辱や拒絶を示す行為として共通認識されており、人手解釈と LLM 出力の間で意味の揺らぎが小さいことを示している。同様に、付録 A-95 のような **Snap Finger gesture** では、平均コサイン類似度が 0.70 台と高く、分散も 0.003～0.005 と小さい値を示した。この結果は、本ジェスチャーが注意喚起や合図など、比較的限定された使用目的で共有されており、文化差が生じにくいジェスチャーであることを反映している。

GPT-5.1	公用語プロンプト	全ての意味
GPT-5.1	英語プロンプト	全ての意味
GPT-4.1mini	公用語プロンプト	全ての意味
GPT-4.1mini	英語プロンプト	全ての意味

表 21 比較を行った条件

と考えられる。また、付録 A-93 のような **Pinched Finger gesture** においても、平均コサイン類似度は、**0.60** 後半～**0.70** 台と高水準であった。一方で、分散は **0.02**～**0.15** と他の文化整合しやすいジェスチャーに比べて大きく、国によって意味解釈の幅が存在することが示された。これは、本ジェスチャーが一部の文化圏では日常的な表現として用いられる一方で、他の文化圏では限定的あるいは異なる意味で解釈されるため、全体としては整合しやすいが、内部に多義性を含むジェスチャーであることを示唆している。

一方で、文化整合しにくいジェスチャーの代表例として、付録 A-93 のような **Quenelle gesture** が挙げられる。本ジェスチャーでは、平均コサイン類似度が **0.36**～**0.45** と低く、分散も **0.01**～**0.12** と大きくばらつきのある値を示した。これは、国や文化圏によって意味解釈が大きく異なり、政治的・社会的背景に依存した解釈が存在するため、人手解釈自体が多様であることが要因と考えられる。その結果、LLM による意味出力も国ごとにばらつきが生じ、局所的文化整合性が低下したと解釈できる。以上により、文化整合しやすいジェスチャーは、意味の方向性が明確で文化間の共有度が高い一方、文化整合しにくいジェスチャーは、社会的背景や使用文脈への依存性が高く、文化固有の解釈差が強く反映されることが明らかとなった。この結果は、局所的文化整合性がジェスチャー固有の意味構造に大きく依存していることを示している。

6.1.2 国ごとの分析

本節では、国ごとの観点から局所的文化整合性について考察する。各国における 22 種類のジェスチャーの評価結果を集約し、人手解釈文と LLM 出力文のコサイン類似度の平均および分散に基づいて、文化整合しやすい国としにくい国の特徴を明らかにする。

国別に平均コサイン類似度を比較した結果、一致度が高い国として、インドネシア、カンボジア、シンガポールが確認された。一方で、一致度が低い国としては、スイス、キューバ、マラウイが挙げられる。この結果は、ジェスチャー解釈における文化の多様性の違いが、LLM 生成データと人手データの対応関係に影響している可能性を示唆している。インドネシア、カンボジア、シンガポールといった東南アジア諸国では、相対的に国内文化の共有度が高く、ジェスチャーの意味解釈が社会内で比較的一貫していると考えられる。これらの国では、英語圏の国々と比較して移民比率が低く、同一の文化的背景を持つ人に都が多数を

占めるため、人手データにおける意味のばらつきが小さくなりやすい結果が見られた。その結果、LLM 生成データとの対応関係が安定し、高いコサイン類似度として観測された可能性がある。

一方、アメリカやカナダといった英語圏の国では、平均コサイン類似度は 0.60 台に留まった。これらの国では、多様な文化圏の人々が共存しており、同一のジェスチャーに対して複数の意味が併存していると考えられる。そのため、個々のジェスチャーでは高い一致度を示す場合がある一方で、ジェスチャーの種類によって値が大きく変動し、全体を平均すると中程度の一致度に収束したと考えられる。

このことから、国ごとの分析により、局所的文化整合性はジェスチャー固有の意味構造に加え、国の文化的背景および使用言語条件によって大きく左右されることが明らかとなった。英語プロンプトは意味解釈を特定の文化圏への収束させやすい一方で、公用語プロンプトは国ごとの差異を保持しやすく、文化的多様性を反映した分析に有効であると考えられる。

6.2 大域的文化的整合性

本節では、LLM がジェスチャーの意味解釈をどの文化圏に位置付けているかという、大域的文化的整合性について考察する。ここで大域的整合性とは、個別の意味一致に留まらず、ジェスチャー解釈がどの文化圏の枠組みに近い形で整理・分類されているかを指す。

6.2.1 デフォルト文化圏への吸着

本節では、本来はアメリカとは異なる文化圏に分類されることが考えられるジェスチャーの解釈が、アメリカと同一文化圏に分類された事例を中心に考察する。このような現象は、LLM が特定の文化的枠組みをデフォルト文化圏として採用し、解釈を引き寄せている可能性を示すものである。人手データと LLM 生成データを比較した結果、本来はアメリカ文化圏とは異なる分類に属するにも関わらず、誤ってアメリカ文化圏と同一クラスターに分類された事例のうち、公用語プロンプト条件において異なる文化圏へ分類されたケースが多く確認された。これらの多くは、マルタ、スペイン、トルコといったヨーロッパ圏の国々に集中しており、欧州諸国がアメリカ文化圏から離れやすい傾向が示された。英語プロンプト条件でも、フランス、イタリア、ギリシャといったヨーロッパ圏において、ジェスチャー解釈がアメリカと同一文化圏に分類されるケースが複数確認され

た。これらの国では、人手データにおいて否定的意味や文脈依存的な解釈が強く現れているにも関わらず、LLM 生成データでは、アメリカにおける一般的、抽象的な意味へと収束する傾向が見られた。特に、政治的・社会的背景への依存度が高いジェスチャーにおいては、本来は国固有の文脈を反映した解釈が必要であるにも関わらず、英語プロンプト条件ではアメリカ的な代表的な意味が優先されやすかった。その結果、クラスタリング結果やコサイン類似度評価によって、アメリカと同一文化圏に分類される現象が生じたと考えられる。

また、LLM によって生成したデータには、アジア圏の国々においても、アメリカにおける一般的な解釈を優先して生成する傾向が確認された。これにより、本来は地域固有の意味を持つジェスチャーであっても、アメリカ的な意味に基づいた解釈が出力されやすく、文化的差異が弱められる傾向が生じていた。

一方で、公用語プロンプト条件では、これらの国々がアメリカとは異なる文化圏として分類されるケースが増加した。公用語による記述では、否定的ニュアンスが明示されやすく、ジェスチャー解釈が特定の文化圏へ一律に吸着する傾向が緩和されたと考えられる。

以上により、大域的整合性の観点からは、英語プロンプトがジェスチャー解釈をアメリカ的な意味に収束させやすい一方で、公用語プロンプトは文化圏間の差異を保持しやすいことが明らかとなった。本研究は、文化的多様性を扱うタスクにおいて、プロンプト言語の選択が解釈の偏りに直接的な影響を及ぼすことを示唆している。

第7章 おわりに

本研究では、文化的背景の異なるジェスチャーの意味解釈に着目し、LLM が生成する意味表現と人手解釈との文化的整合性をどのように評価できるかを検討した。特に、ジェスチャーという非言語行動が国や文化圏によって多義的に解釈される点に着目し、意味表現の生成方法および評価手法の分析を行った。本研究において取り組んだ課題は、大きく以下の2点である。

文化を考慮した意味表現

各国のジェスチャー説明文は、文体、語彙、記述の詳細度が大きく異なり、そのままでは文化的特徴を保持したまま比較することが困難である。また、LLM による意味生成では、文化中立的な抽象化が生じ、特定文化に固有のニュアンスが失われる場合がある。本研究では、英語プロンプトと公用語を比較することで、公用語を用いた場合に各国固有の意味や使用制約がより反映されやすくなることを示し、文化的特徴を保持した意味表現の重要性を明らかにした。

文化差に基づく文化整合性判定

従来の意味距離に基づく評価では、LLM 出力が人手解釈とどの程度一致しているかという局所的な文化整合性のみが評価され、一致しなかった場合に、どの文化圏の解釈に影響を受けているのかを捉えることが難しかった。本研究では、ジェスチャーごと、国ごとの分析を通じて局所的な文化整合性を評価するとともに、クラスタリング結果から考察を行った。その結果、英語プロンプト条件では、文化的背景が異なる国のジェスチャー解釈であっても、アメリカを中心とする英語圏の意味枠組みに収束する傾向が確認された。一方で、公用語プロンプト条件では、国固有の解釈差が保持されやすく、文化的多様性を反映した意味が生成されることが示唆された。これらの結果から、LLM 生成データによるジェスチャーの意味は、モデル性能だけでなく、プロンプト言語の選択によって文化的偏りの程度が大きく左右されることが明らかとなった。

謝辞

本研究を進めるにあたり，終始懇切丁寧なご指導と貴重なご助言を賜りました村上陽平教授に深く感謝申し上げます．また，**Mondheera Pituxcoosuvarn** 先生にも多大なるご支援を頂きましたことを併せて御礼申し上げます．さらに，本研究に関して多くの有益な助言をいただき，親身に相談に乗って下さいました森叶葉様にも心より感謝申し上げます．加えて，日頃より研究活動を支えてくださり，普段からお世話になっている社会知能研究室の皆様に，心より御礼申し上げます．

参考文献

- [1] David McNeill: Hand and Mind: What Gestures reveal about Thought, University of Chicago Press(1994).
- [2] Desmond Morris: Gestures: their origins and distribution, Jonathan Cape(1979).
- [3] Paul Ekman, Wallace Friesen: The Repertoire of Nonverbal Behavior: Categories, Origins, Usage, and Coding, Vol.1, No.1, pp.49–98(1969).
- [4] Abhinav Aggarwal, Sarah E. Chen, Michael S. Bernstein, Percy Liang: Mind the Gesture: Evaluating AI Sensitivity to Culturally Offensive Non-Verbal Gestures, Proceedings of the AAAI Conference on Artificial Intelligence(2023).
- [5] Vladimir I. Pavlovic, Rajeev Sharma, Thomas S. Huang: Visual Interpretation of Hand Gestures for Human–Computer Interaction: A Review, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.19, No.7, pp.677–695(1997).

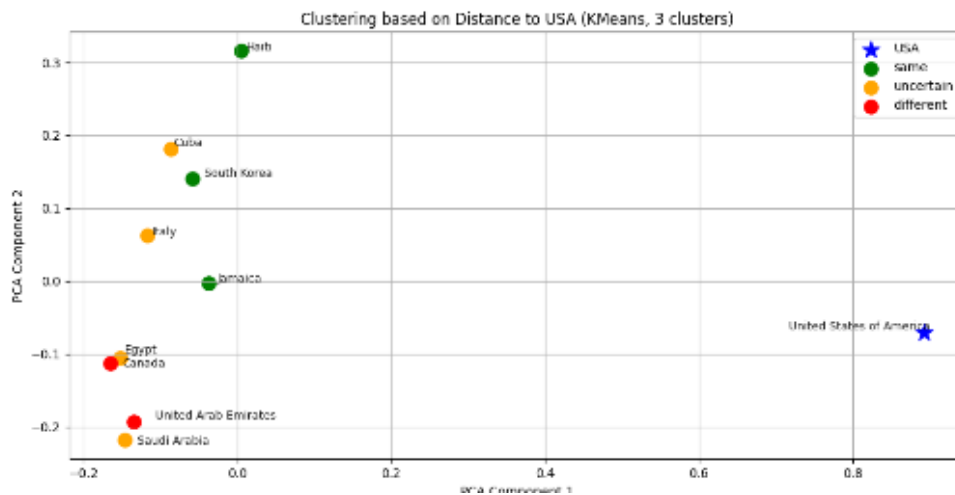
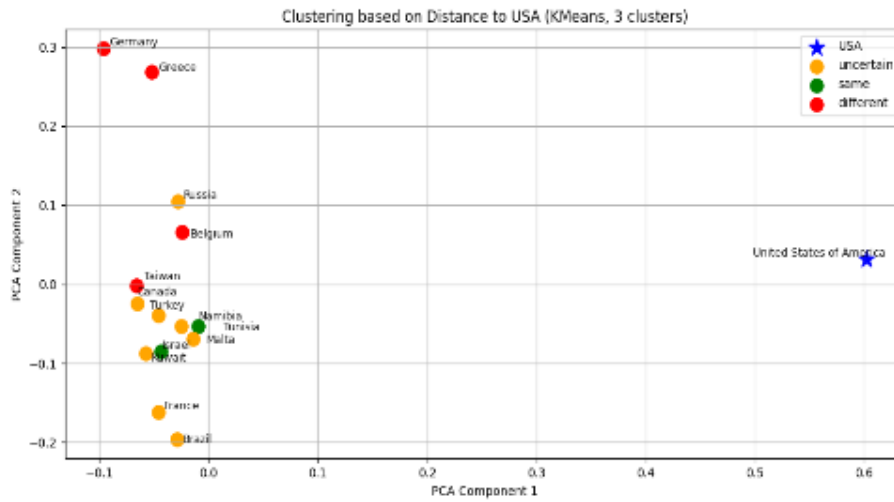
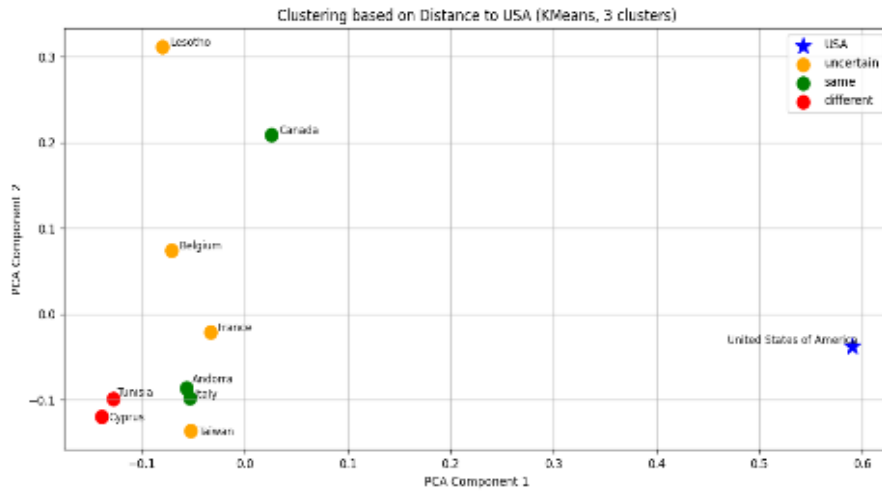
付録

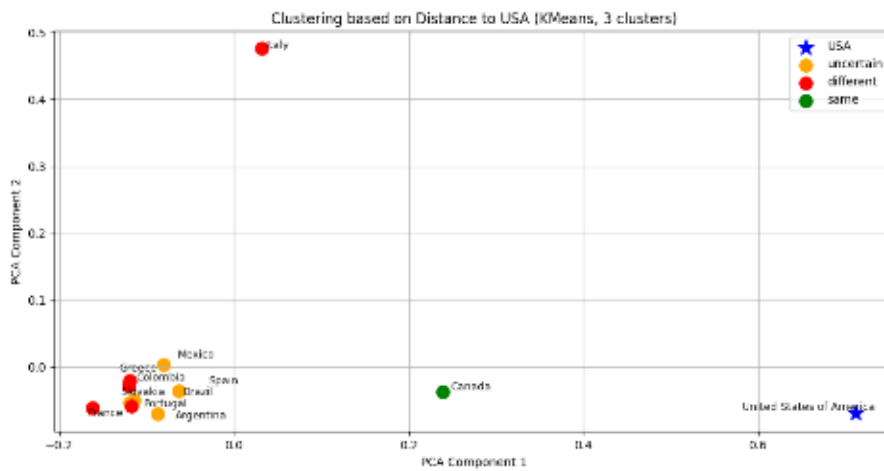
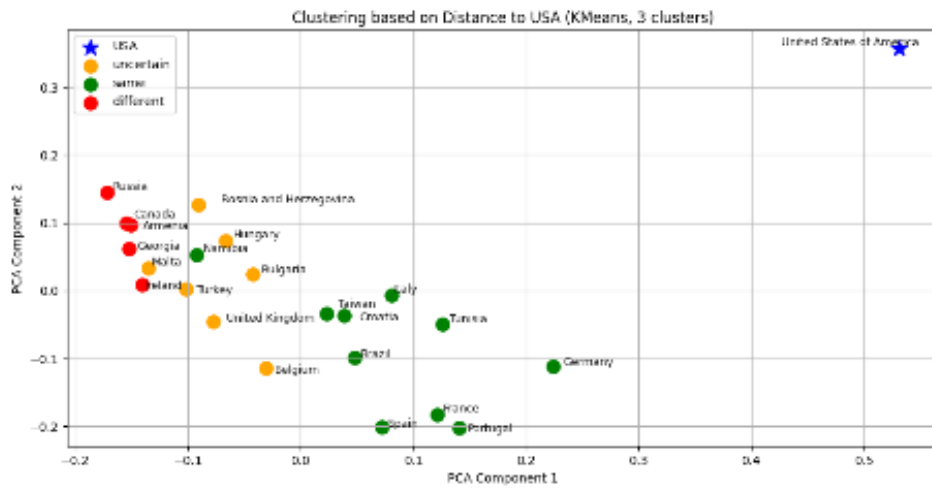
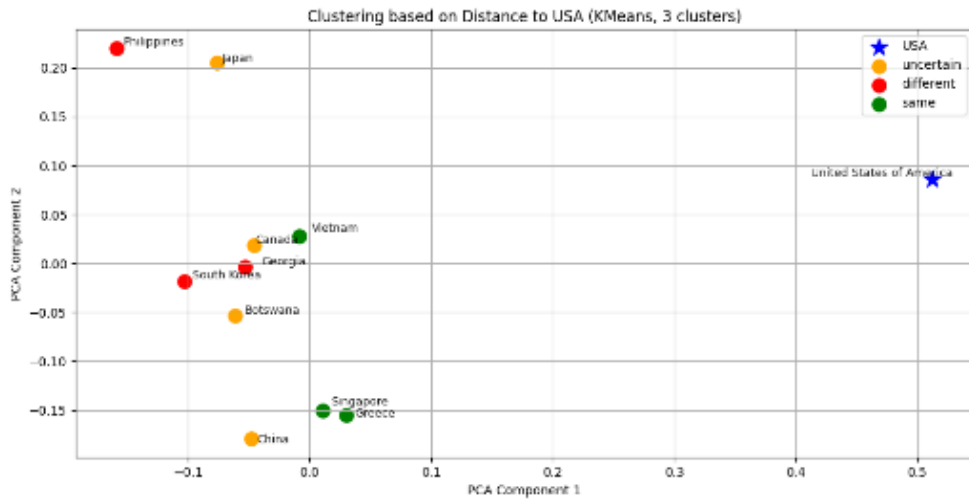
A.1 クラスタリング結果

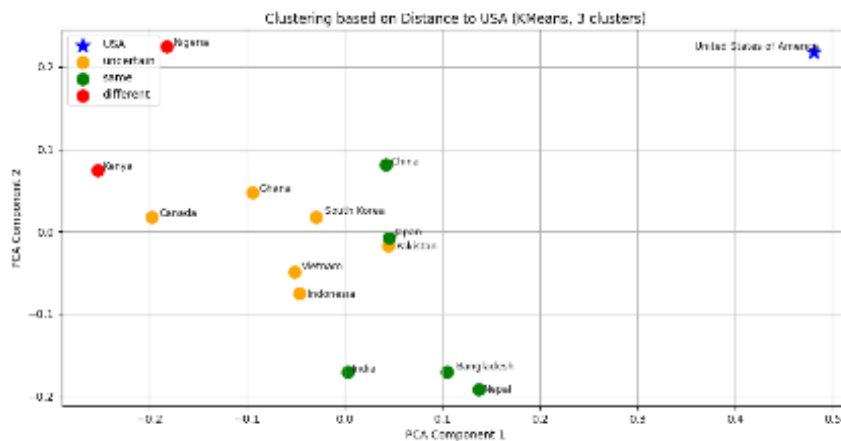
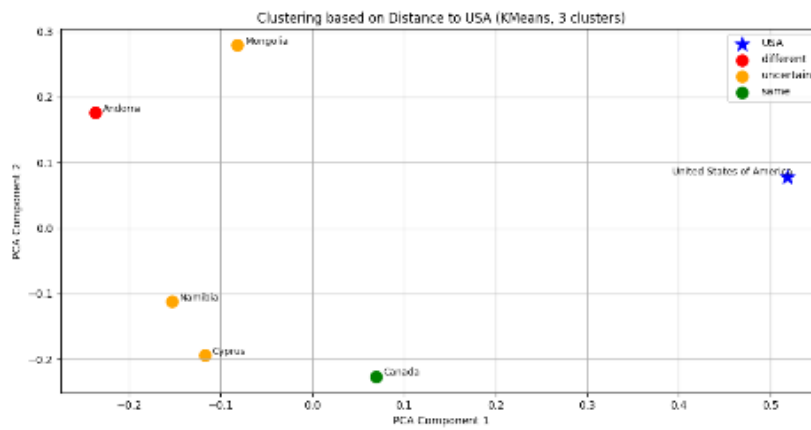
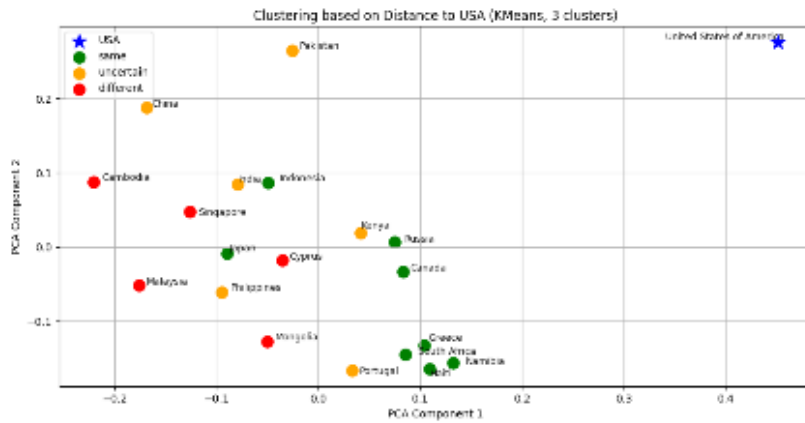
結果の順は、以下の通りである.

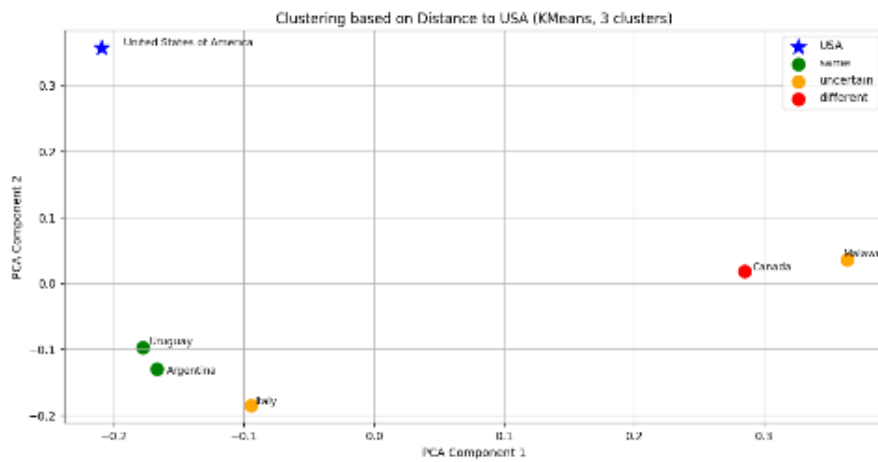
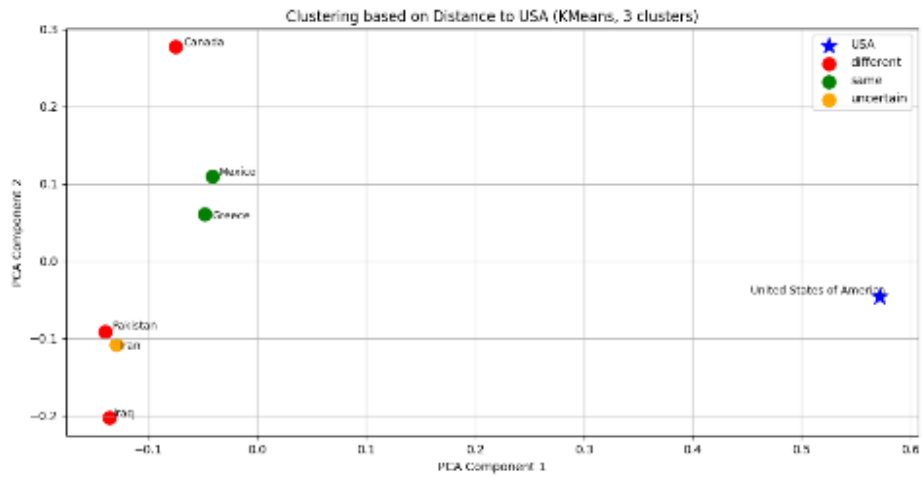
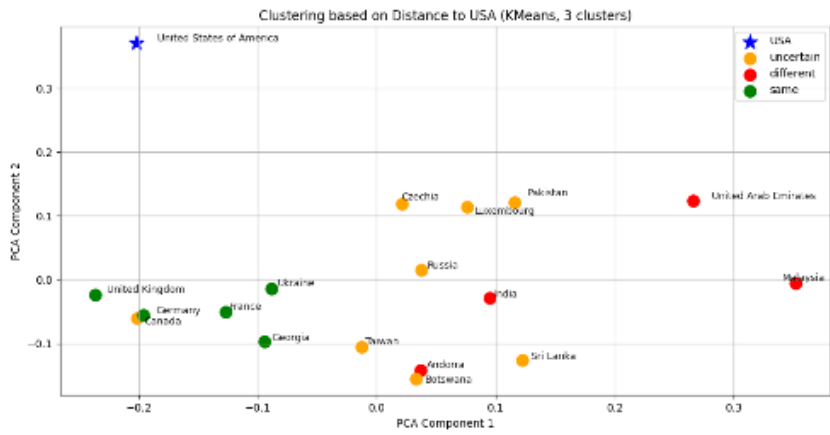
- OK
- Chin flick
- Curled finger
- Five fathers
- Forearm Jerk
- Horns
- Index finger pointing
- L
- Left hand
- Middle finger
- Open palm with fingers spread
- Pinched fingers
- Quenelle
- Serbian Salute
- Shocker
- Show sole of shoe or feet
- Snap fingers
- Three finger salute
- Thumbs up
- Touching someone head
- V sign
- Wanker

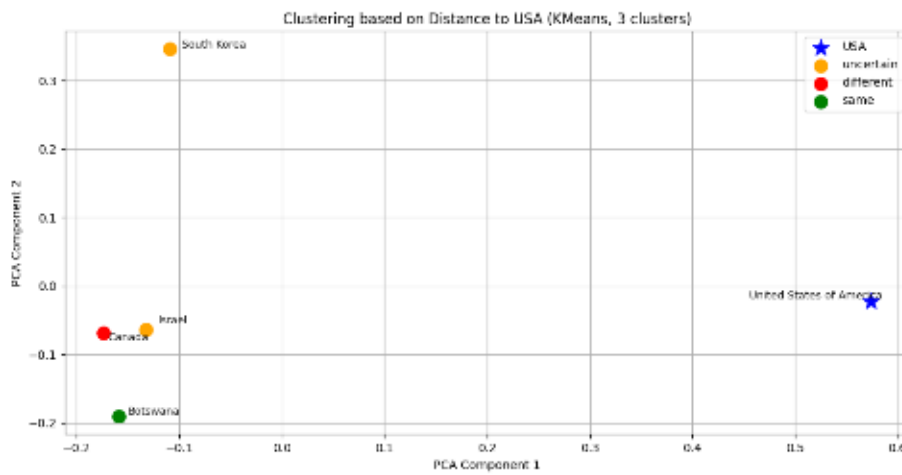
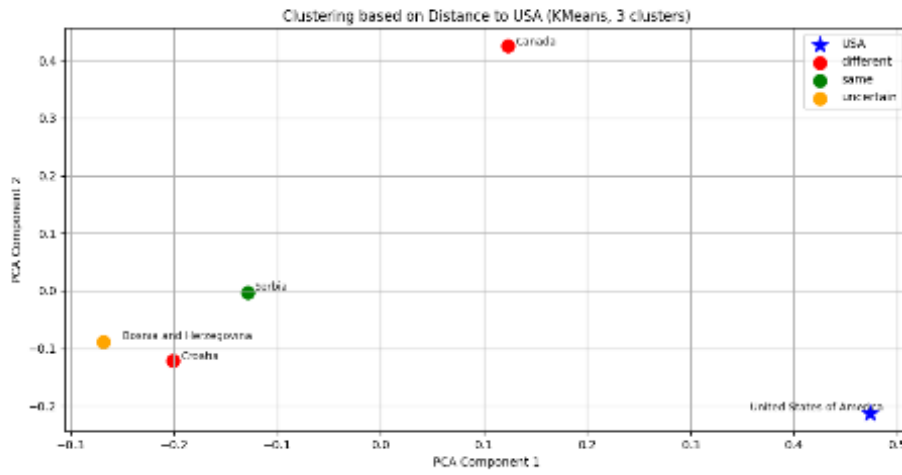
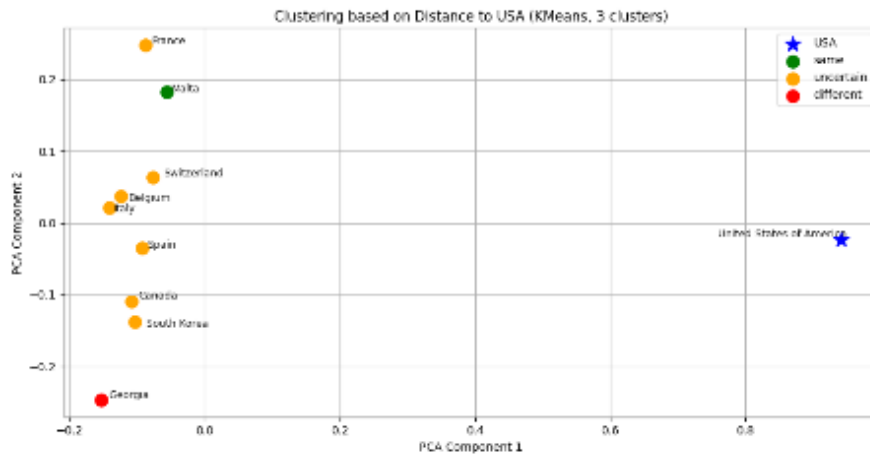
GPT-4. 1mini 英語プロンプト 全ての意味を出力

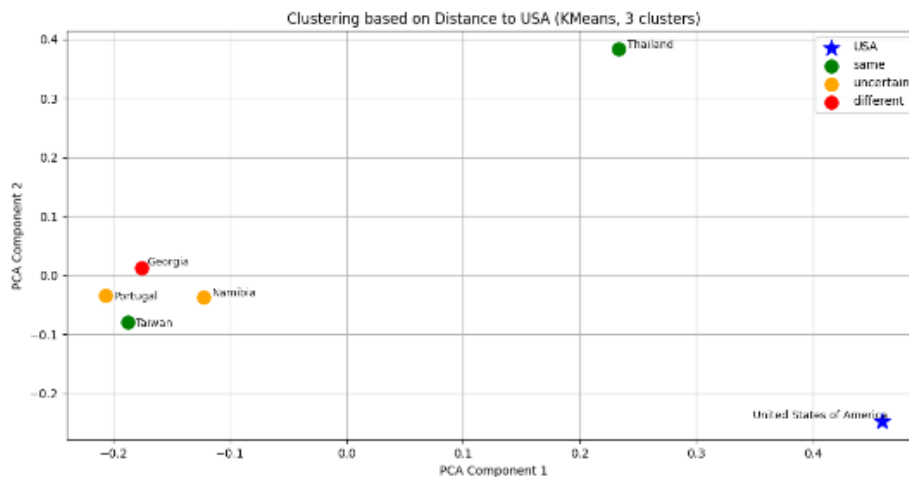
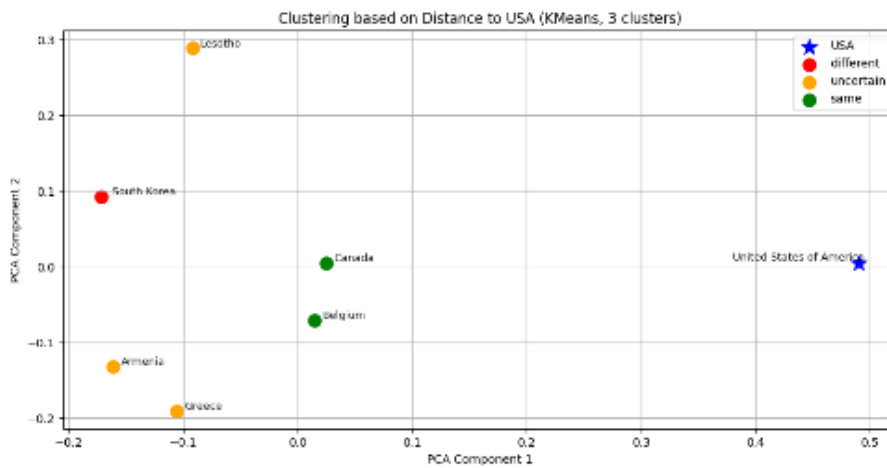
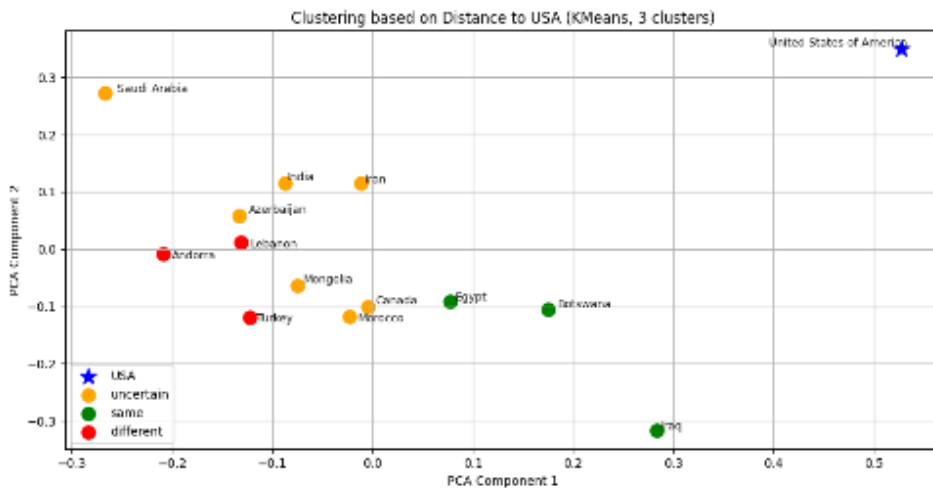


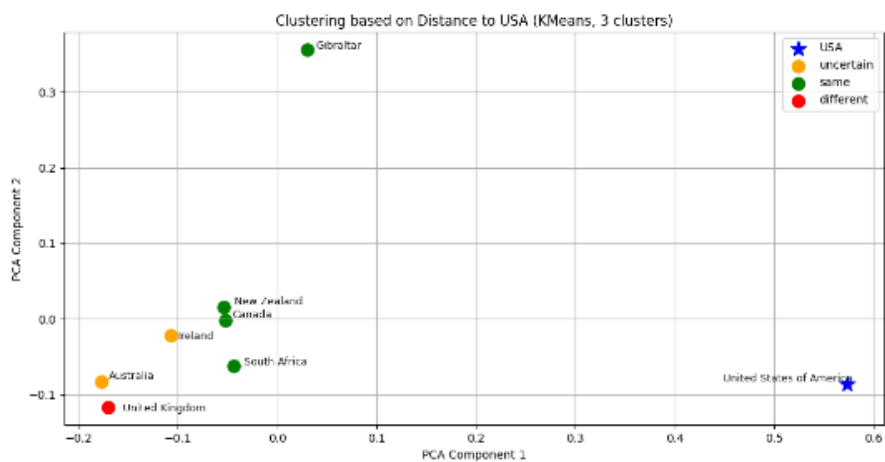
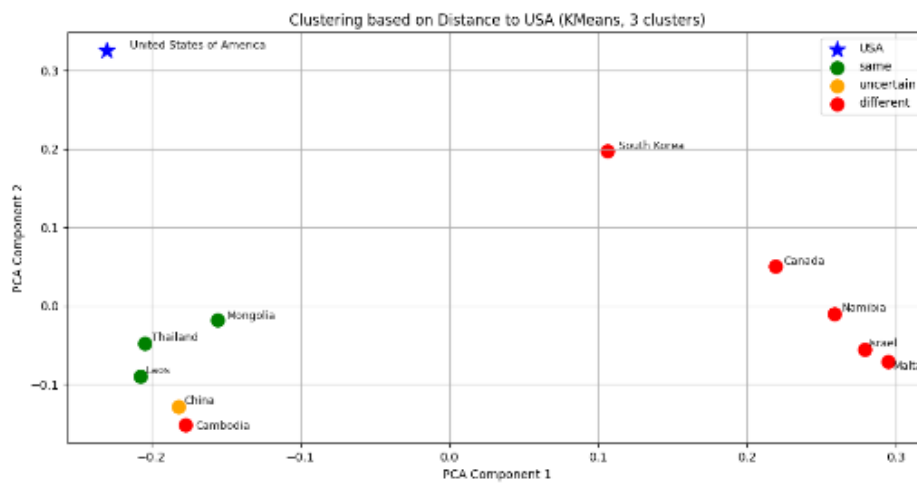
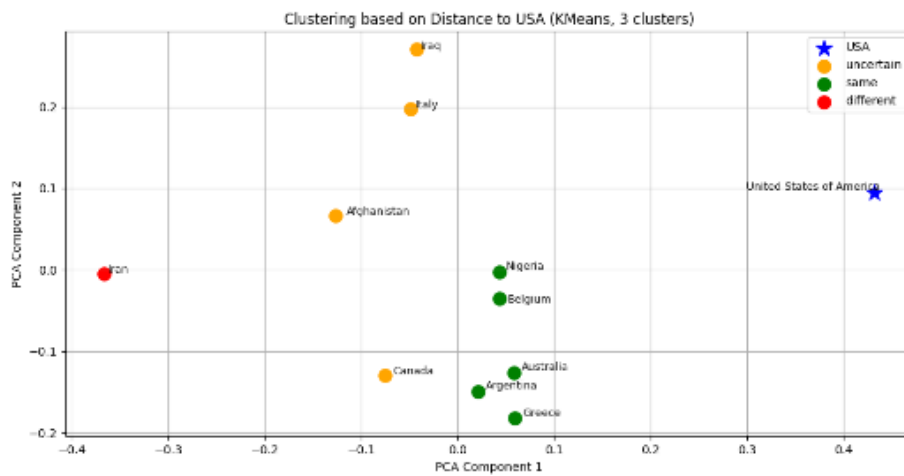


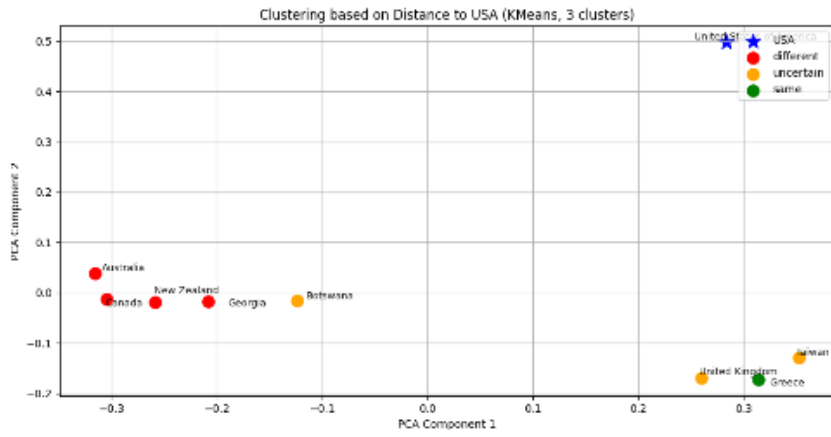




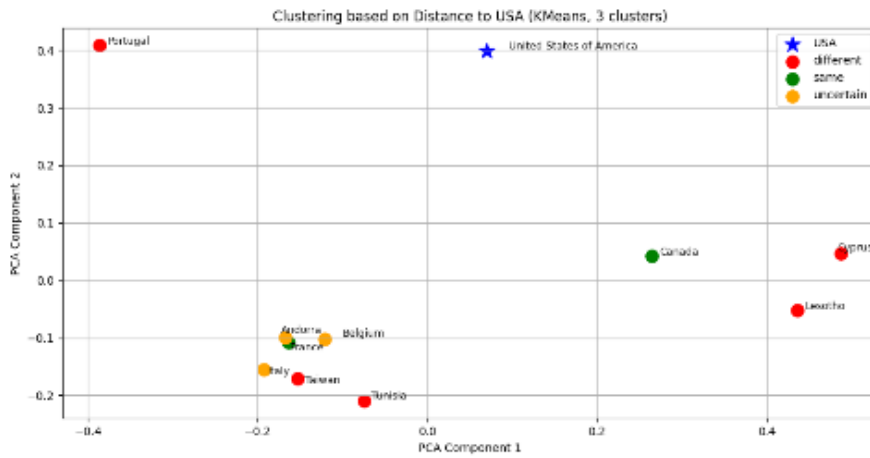
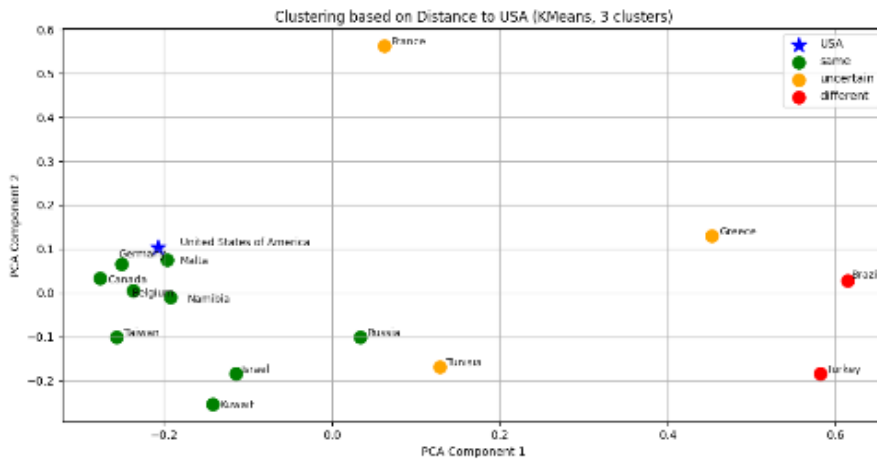


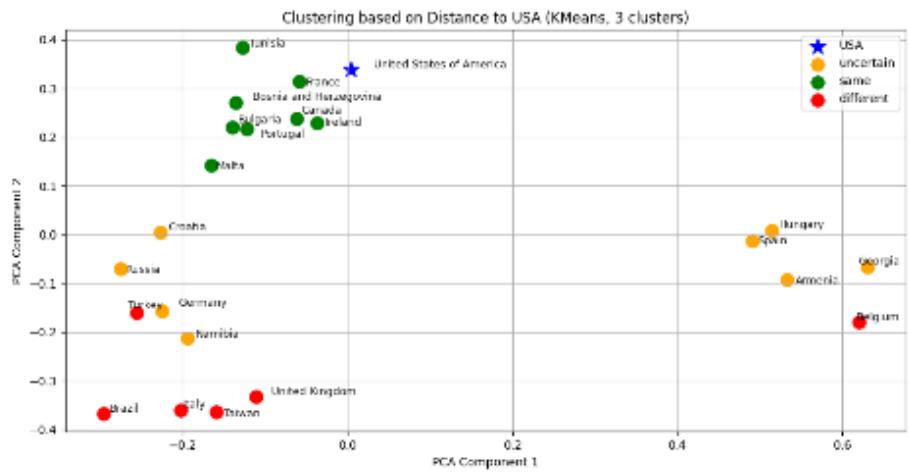
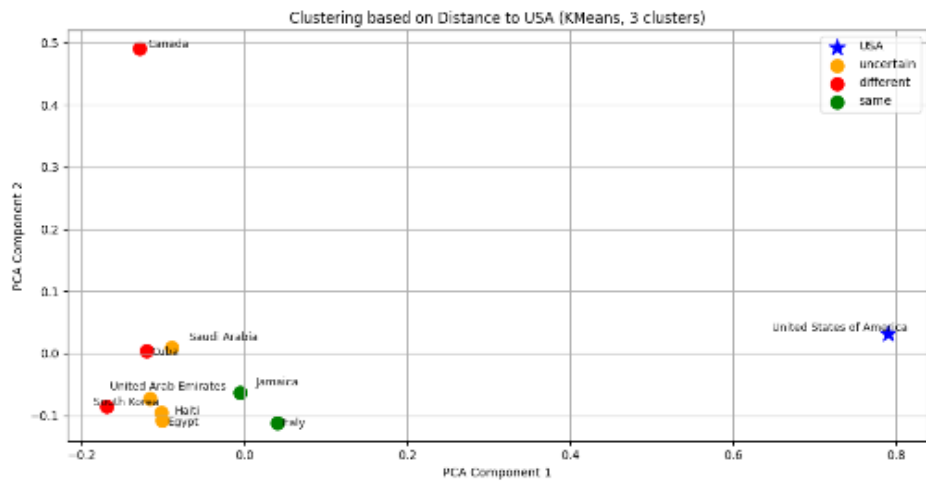
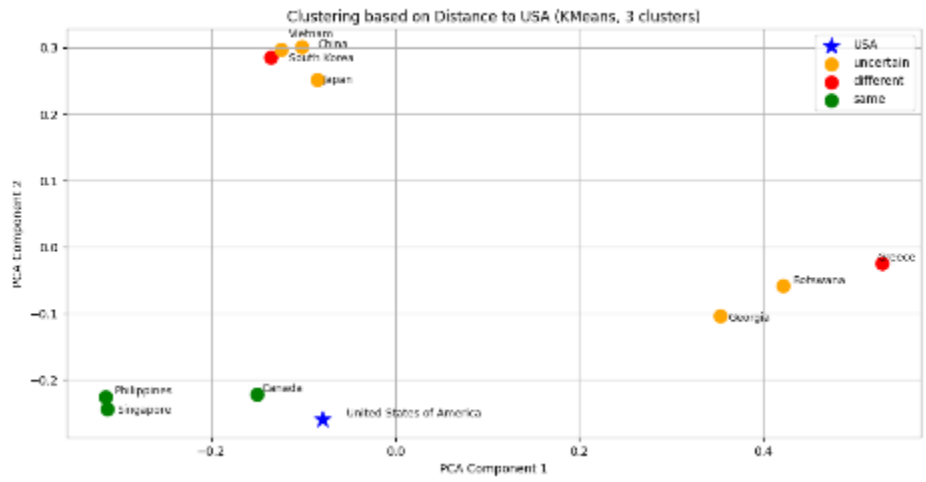


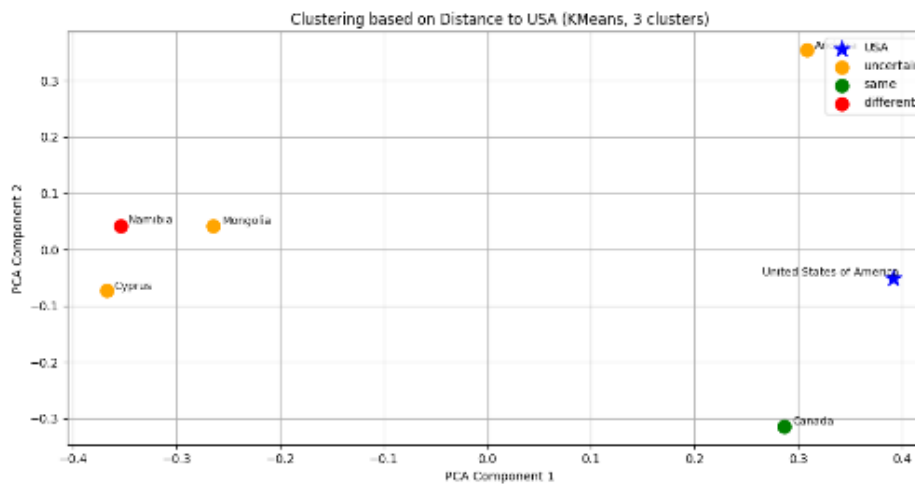
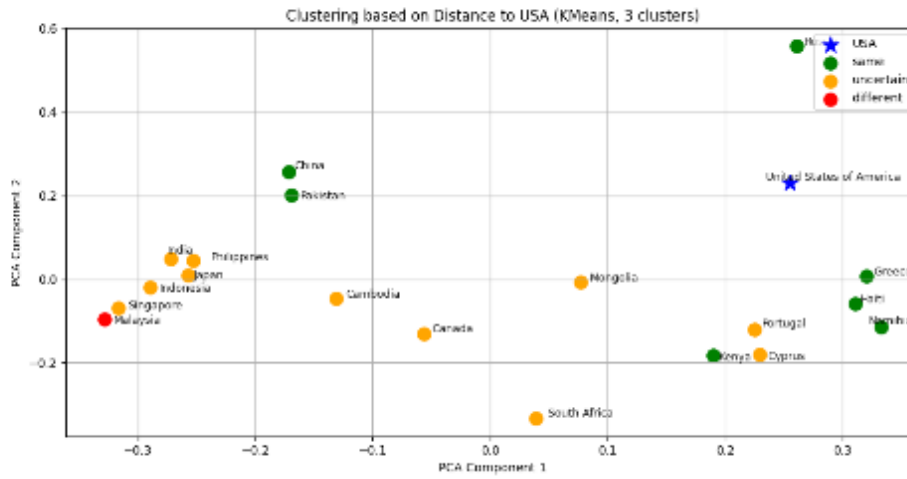
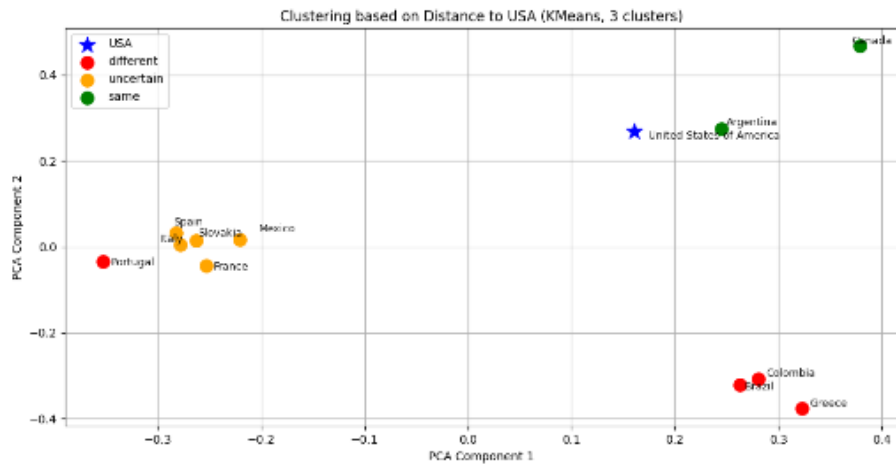


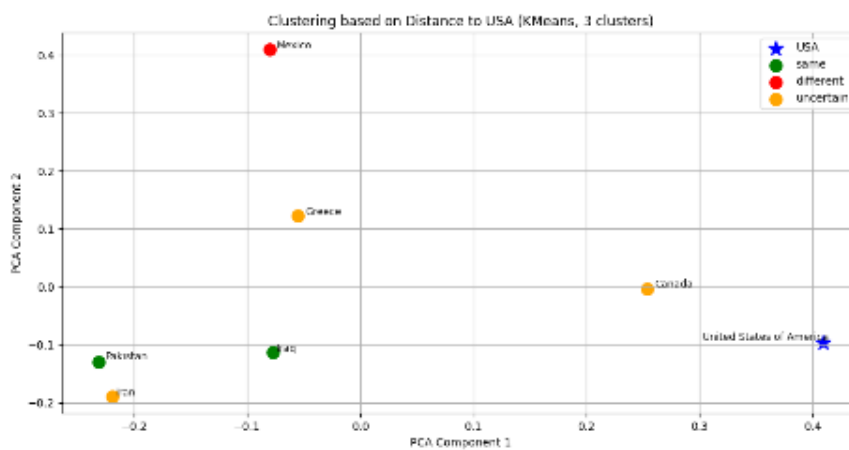
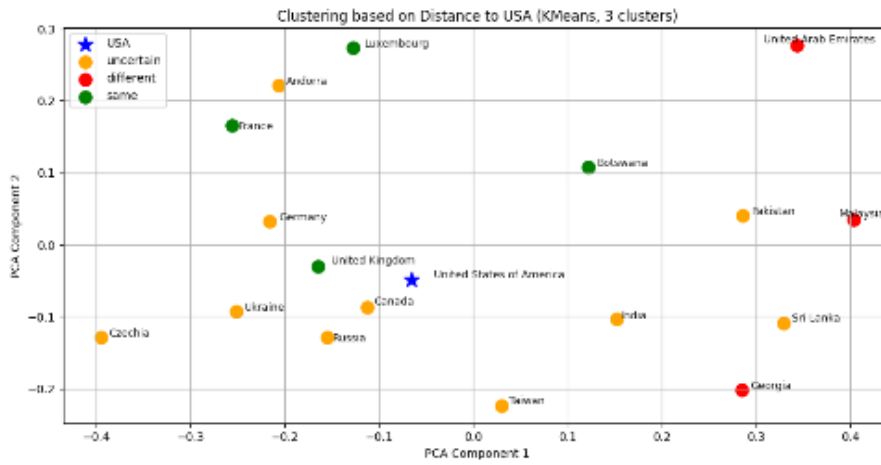
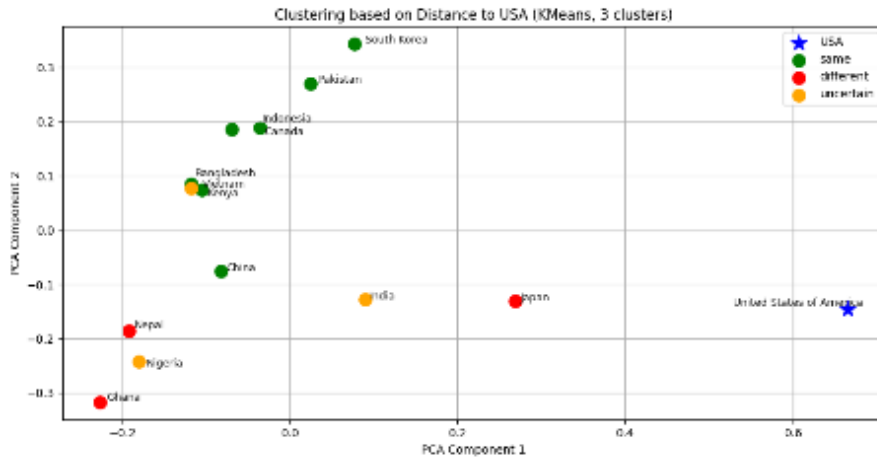


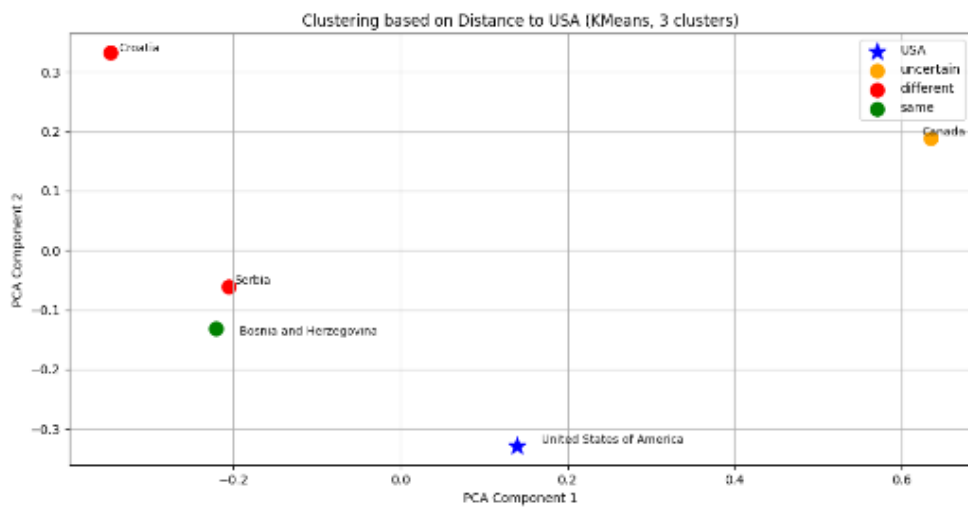
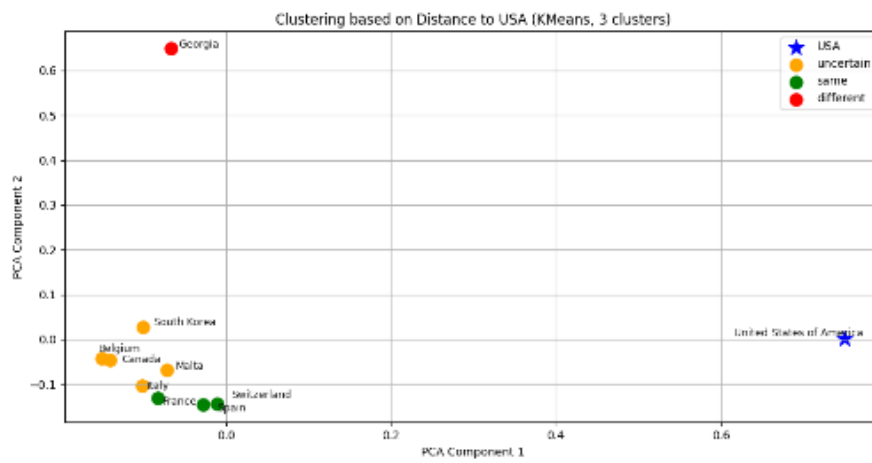
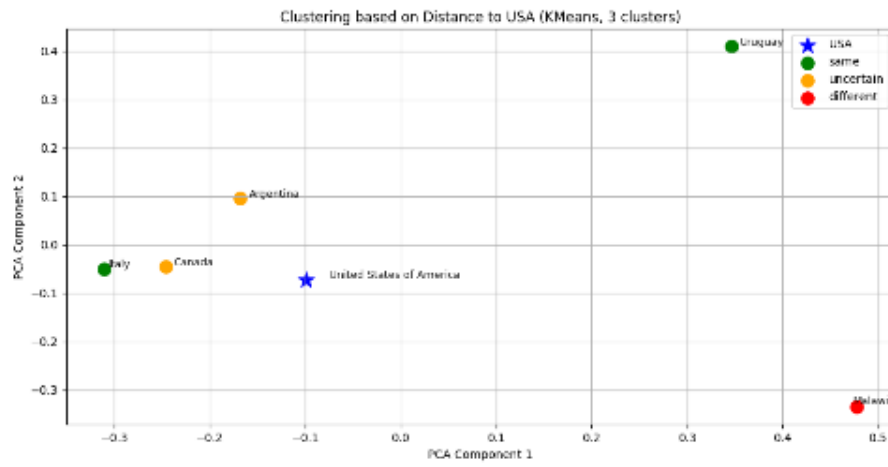
GPT-4. 1mini 英語プロンプト 代表的意味を出力

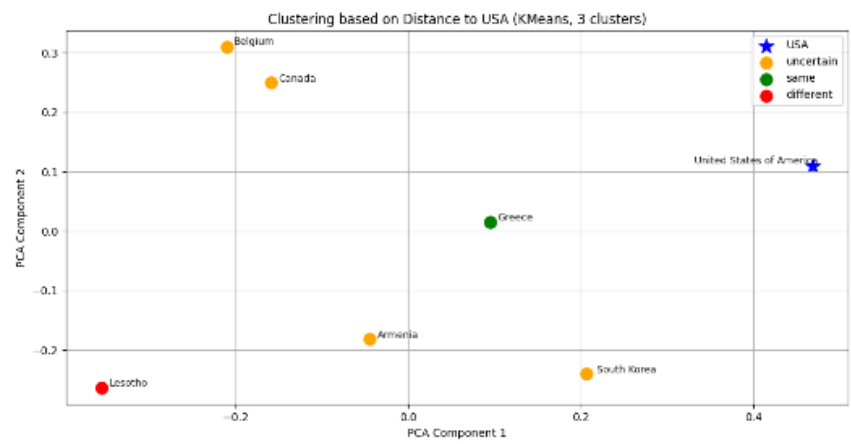
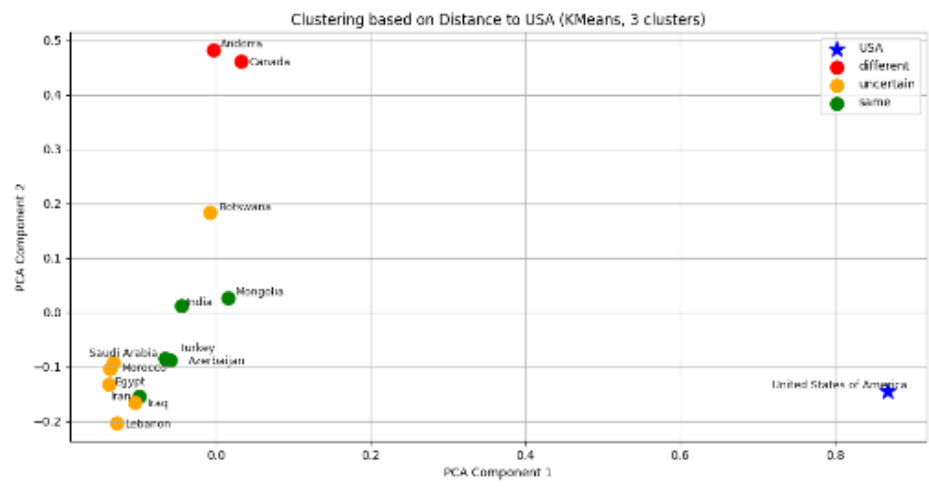
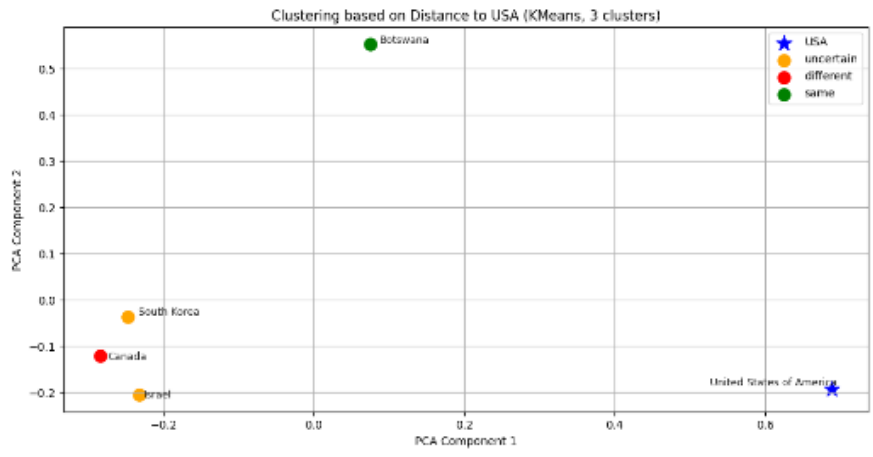


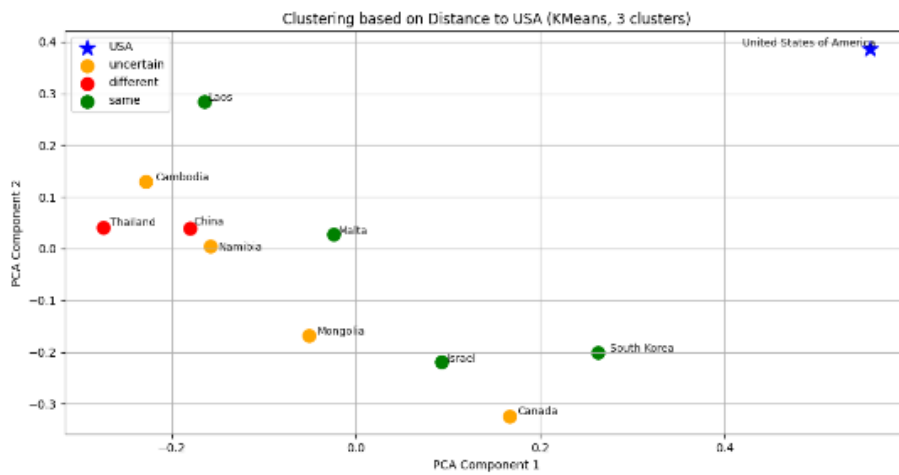
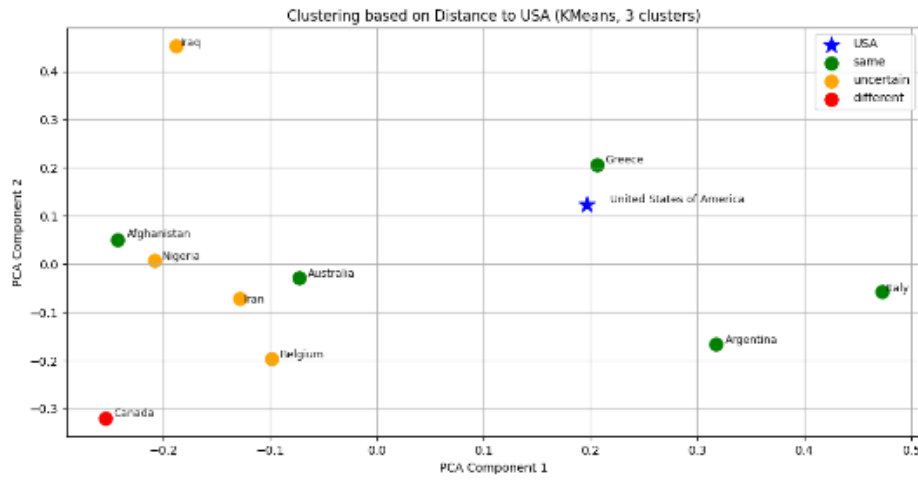
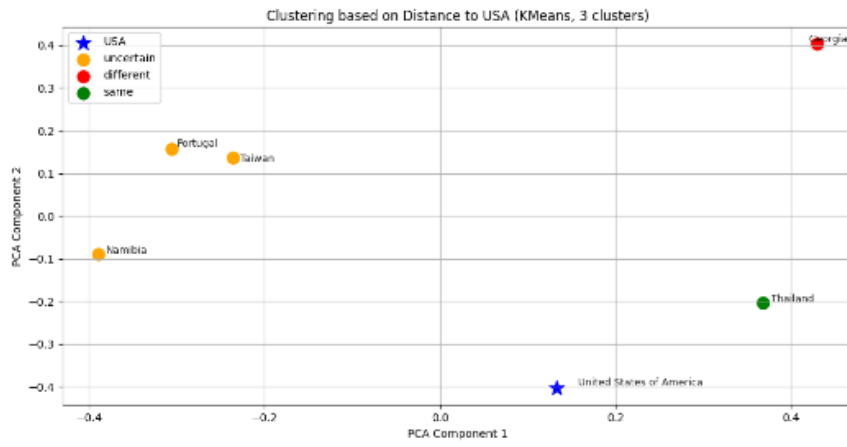


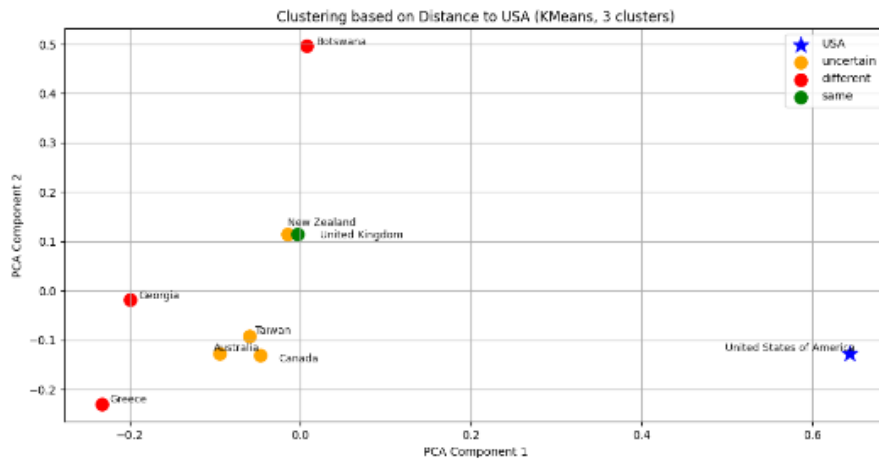
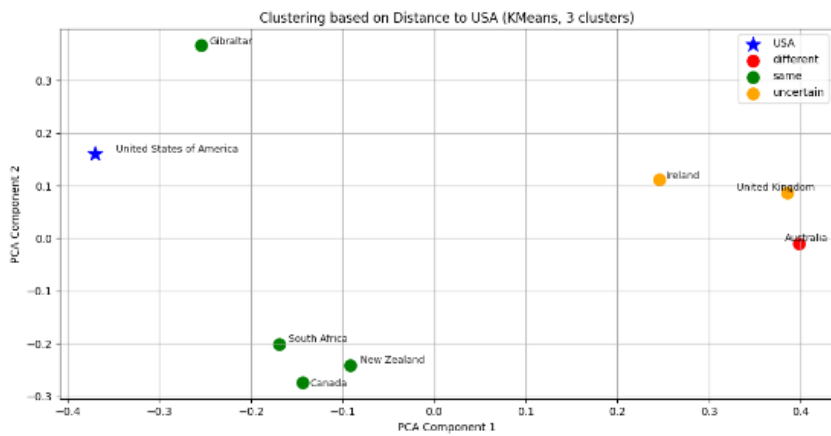




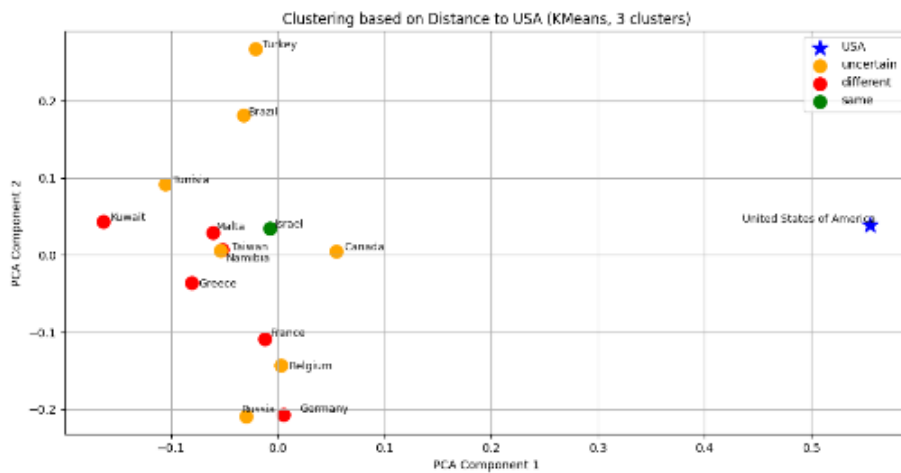


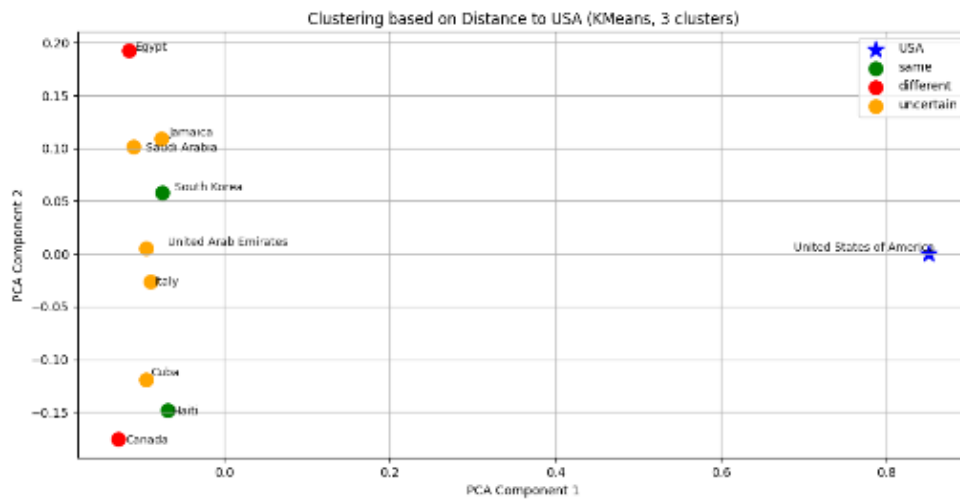
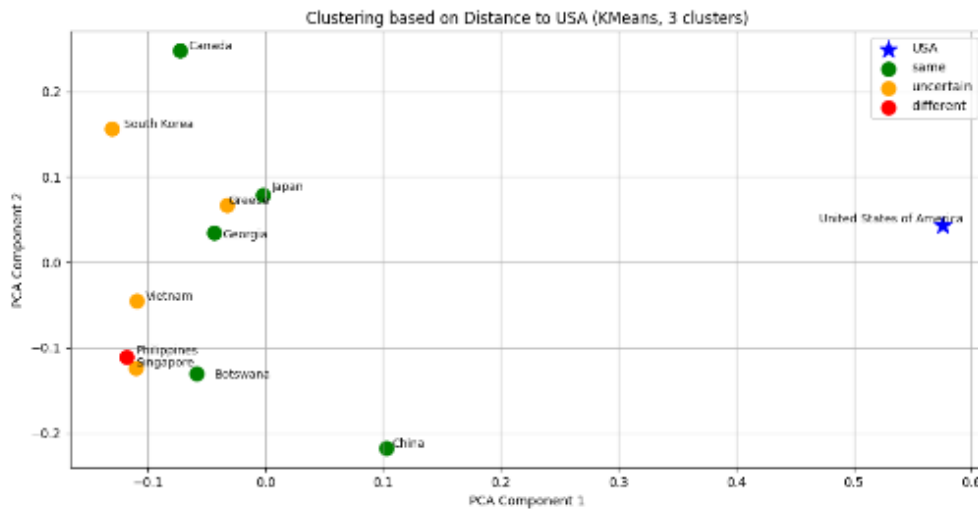
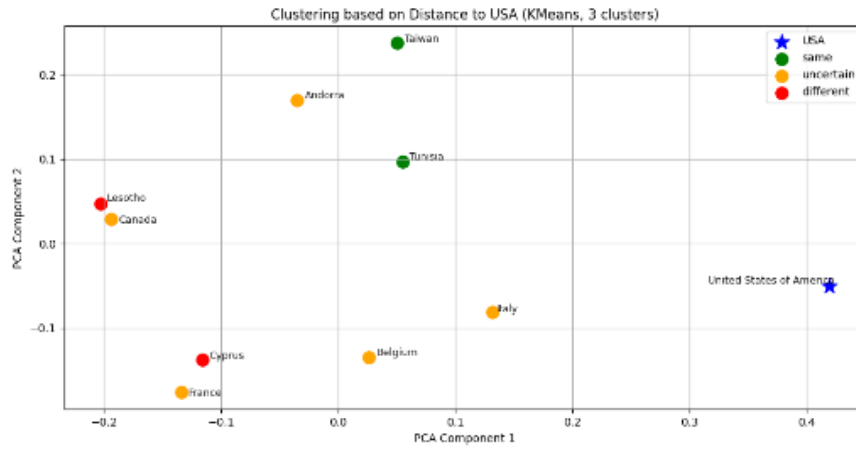


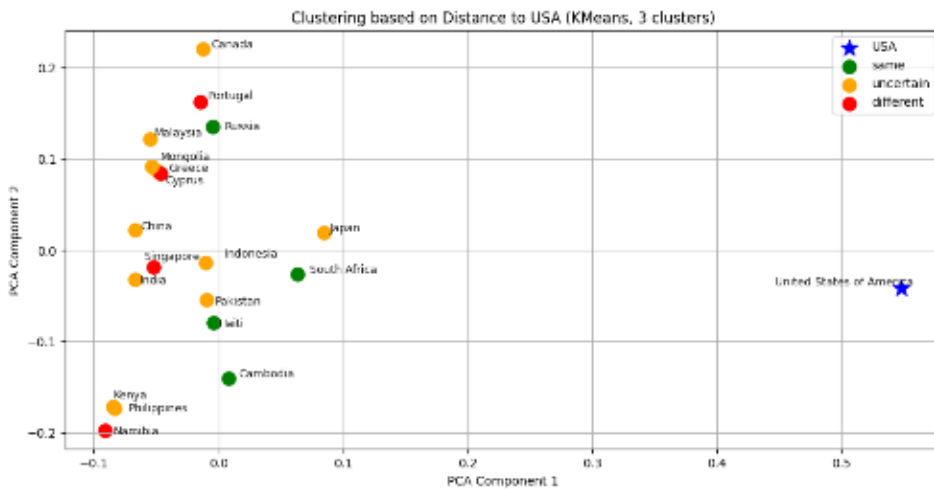
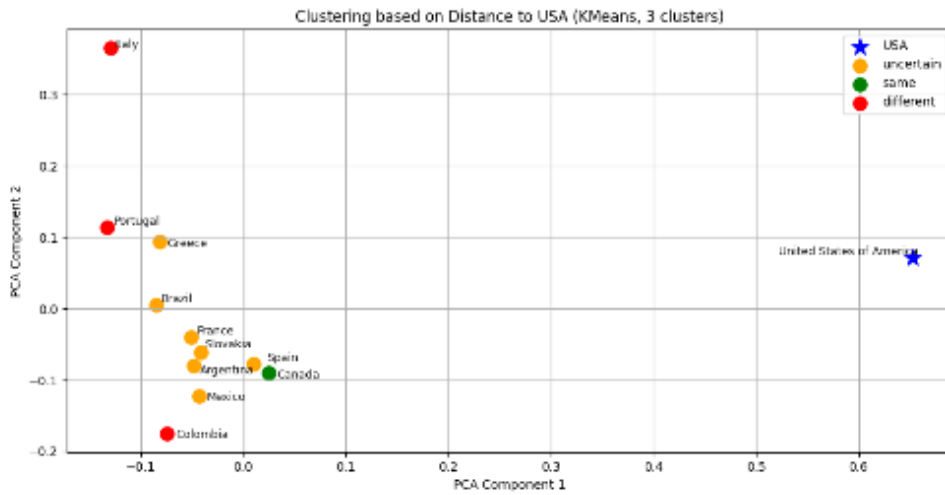
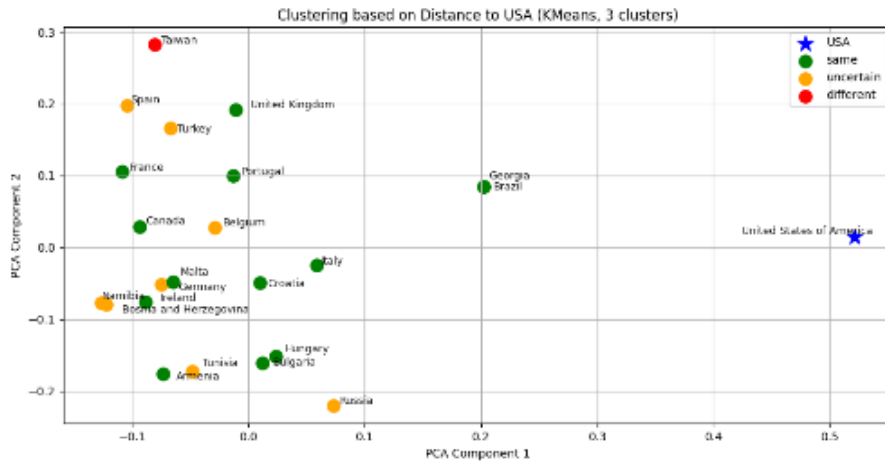


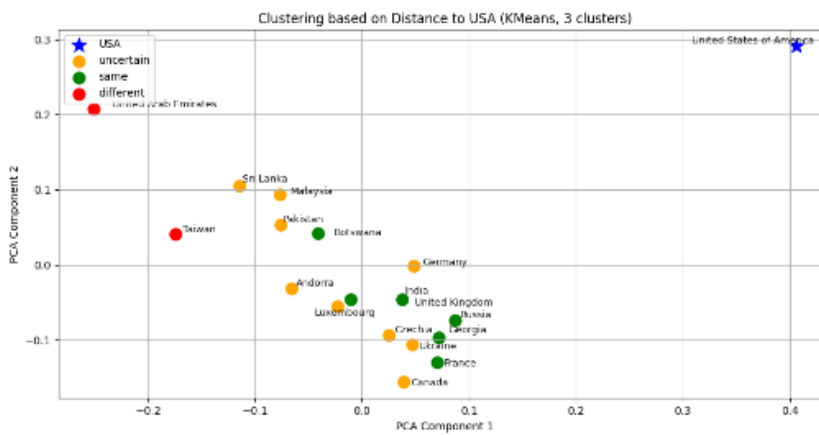
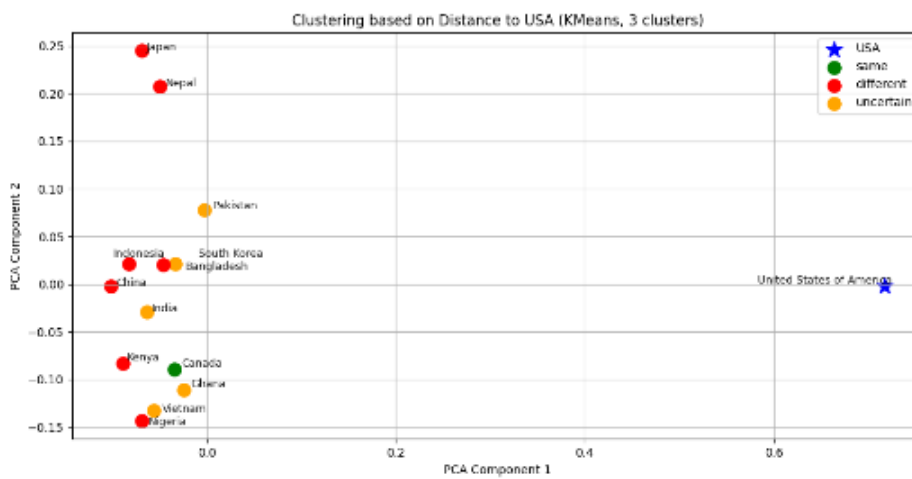
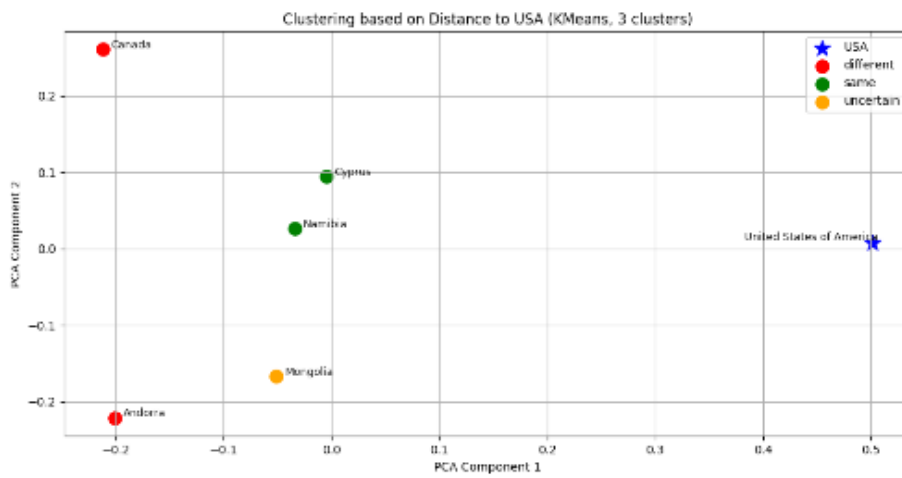


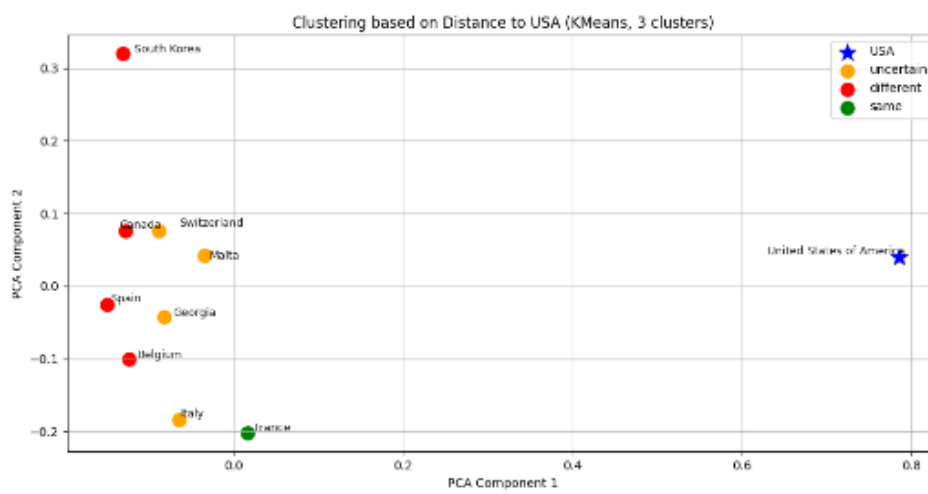
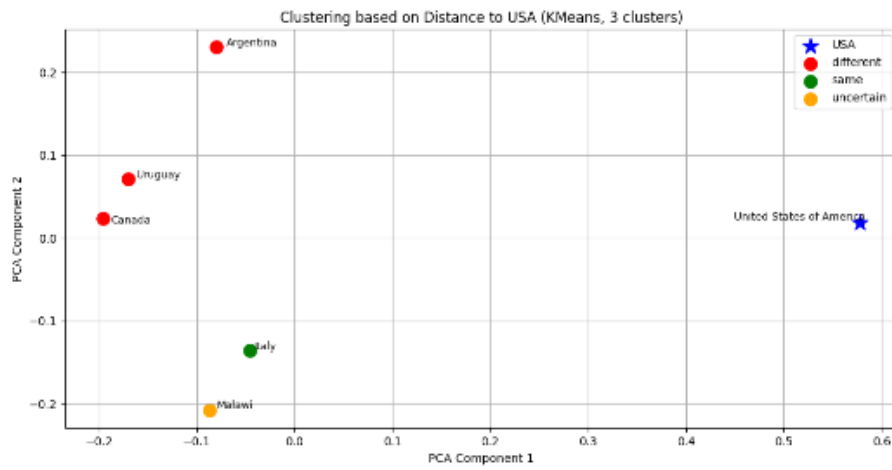
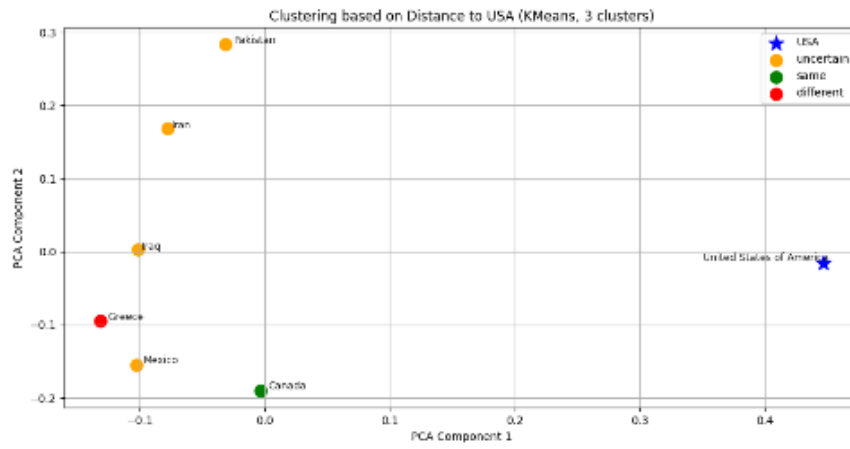
GPT-5. 1 英語プロンプト 全ての意味を出力

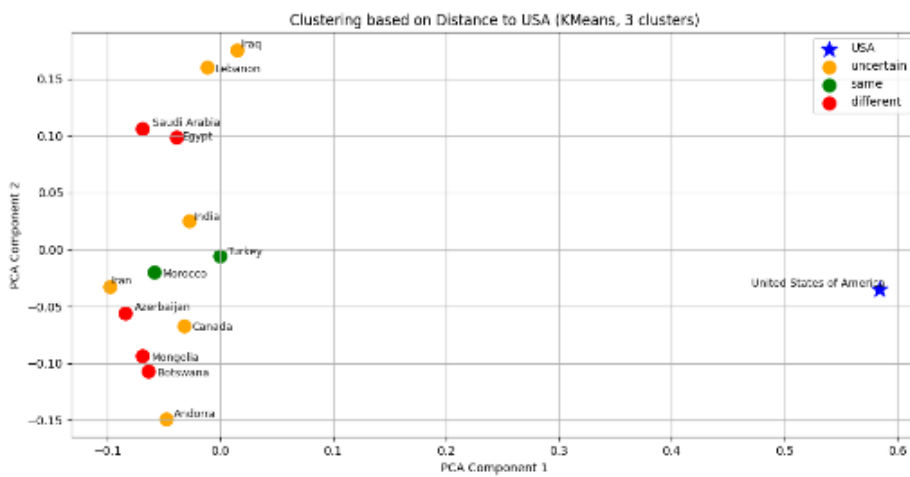
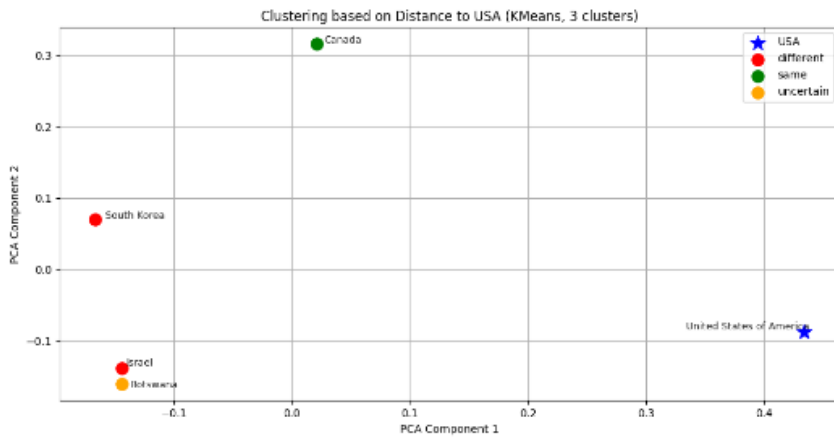
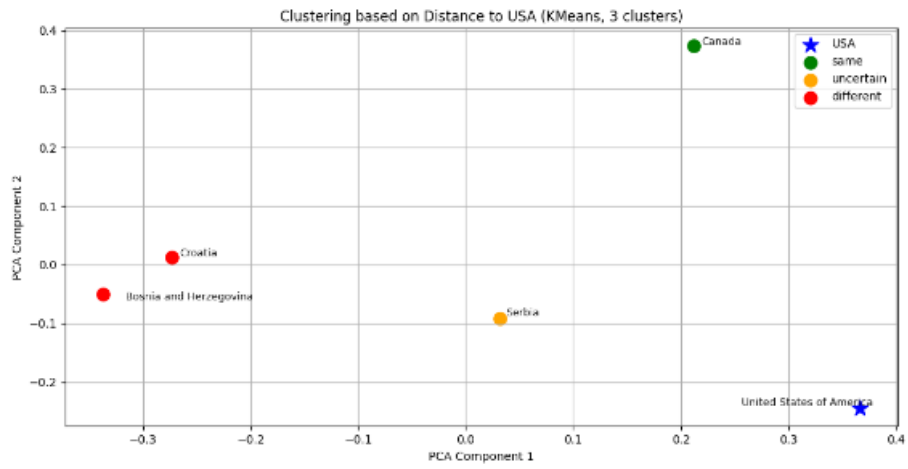


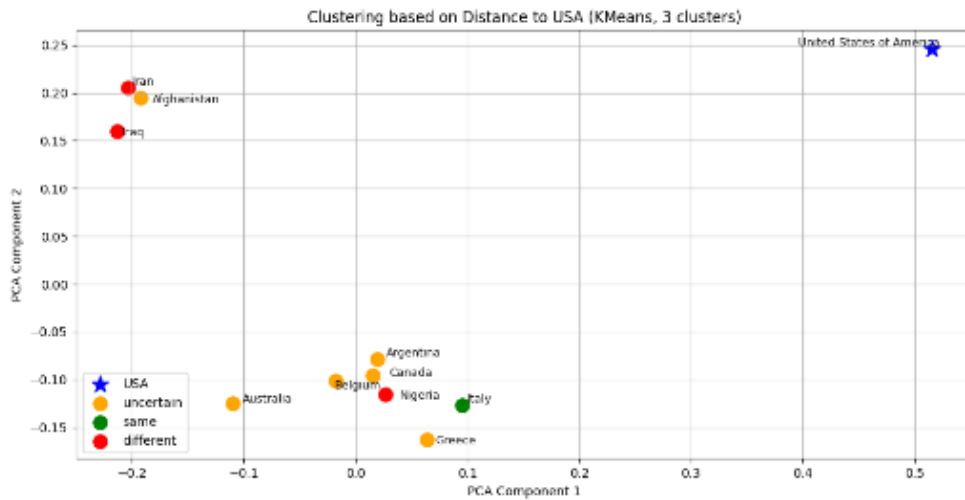
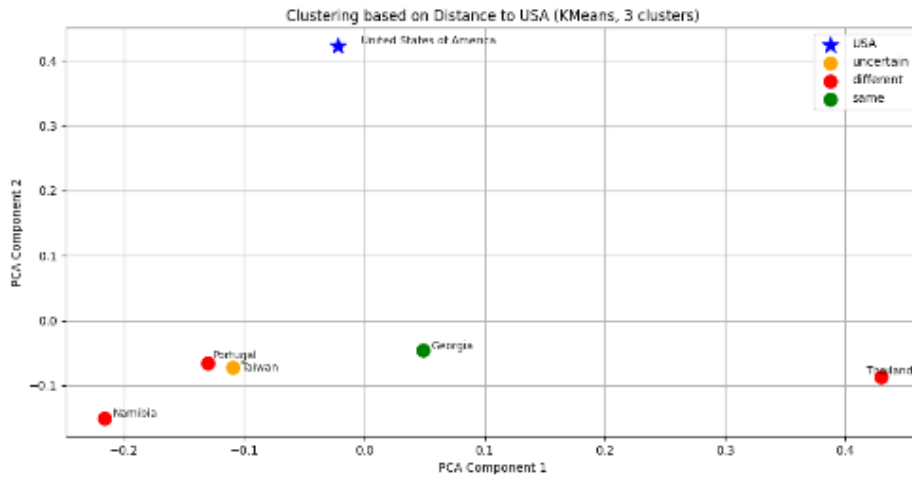
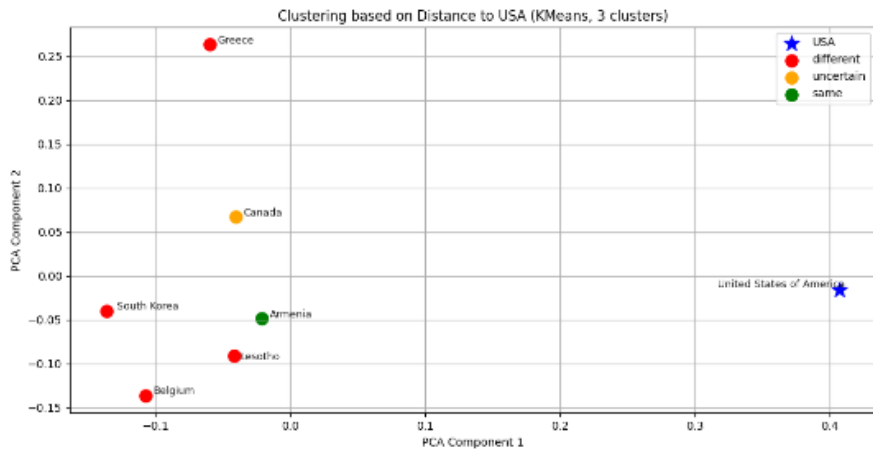


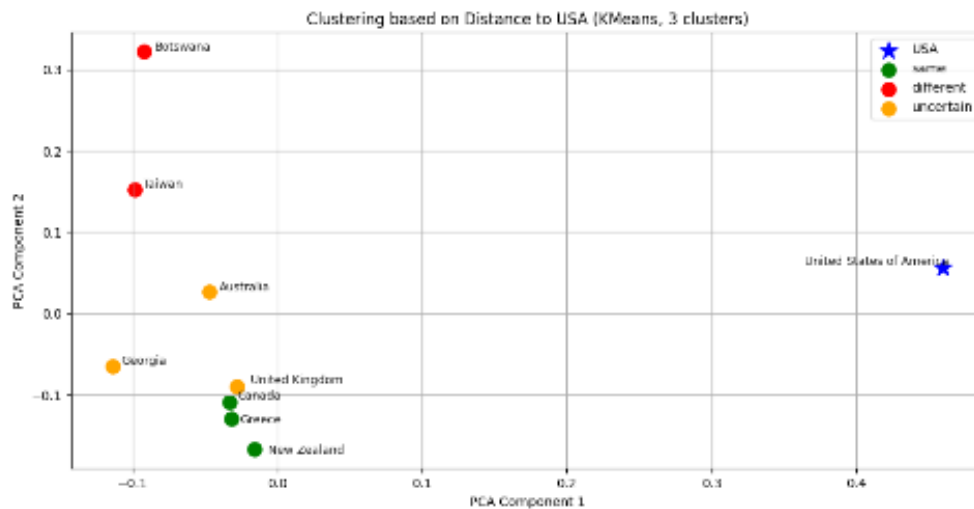
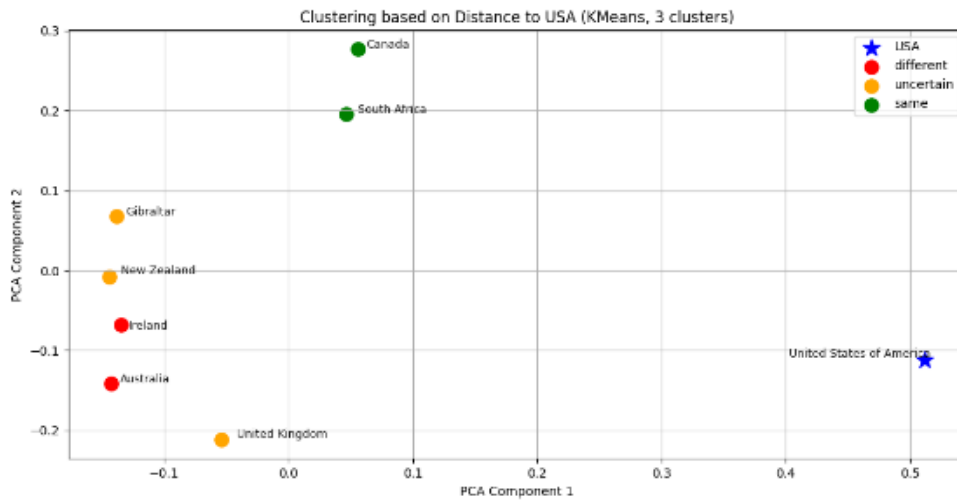
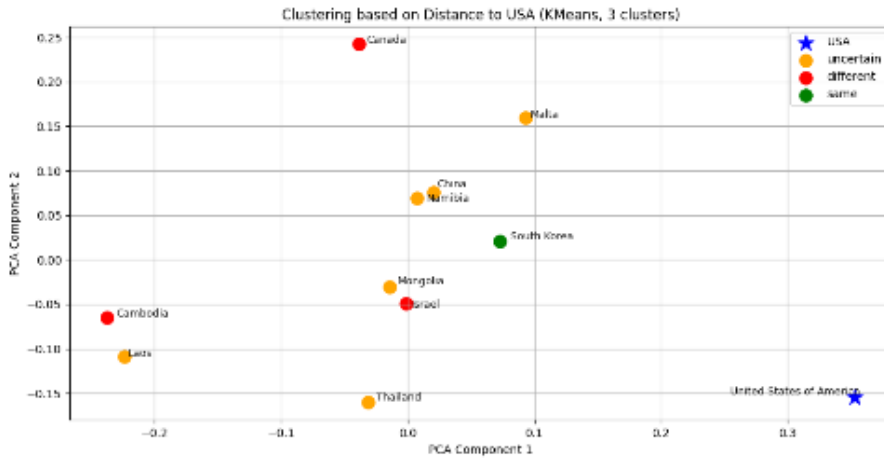




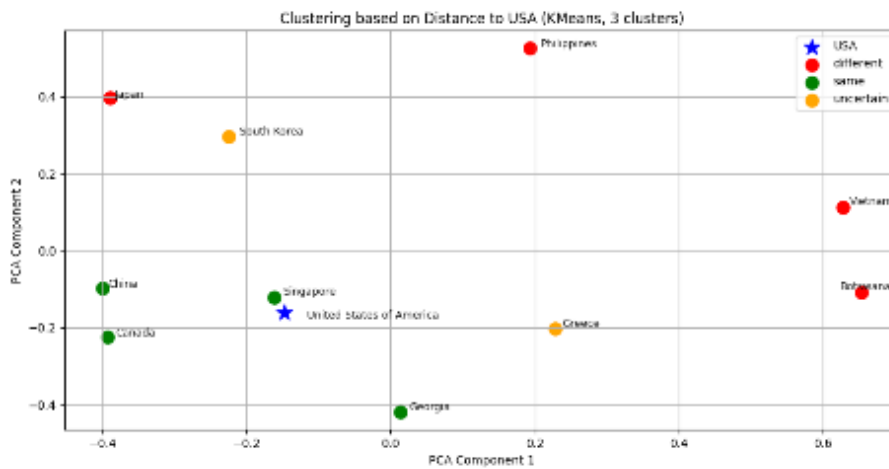
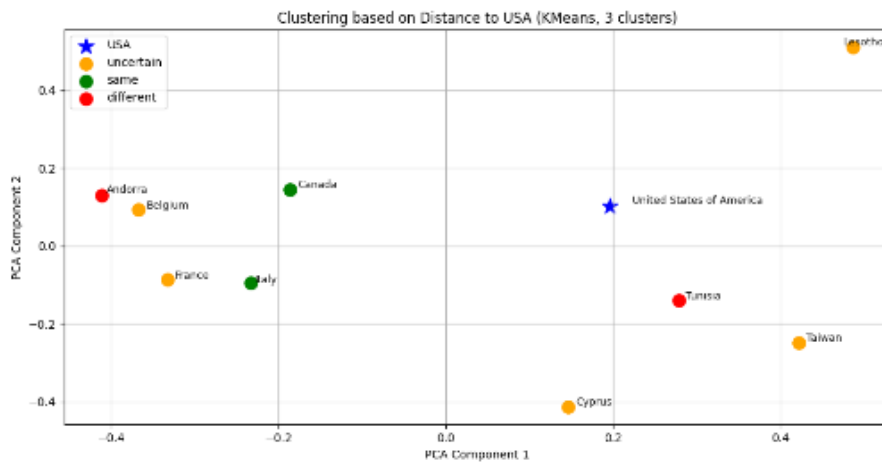
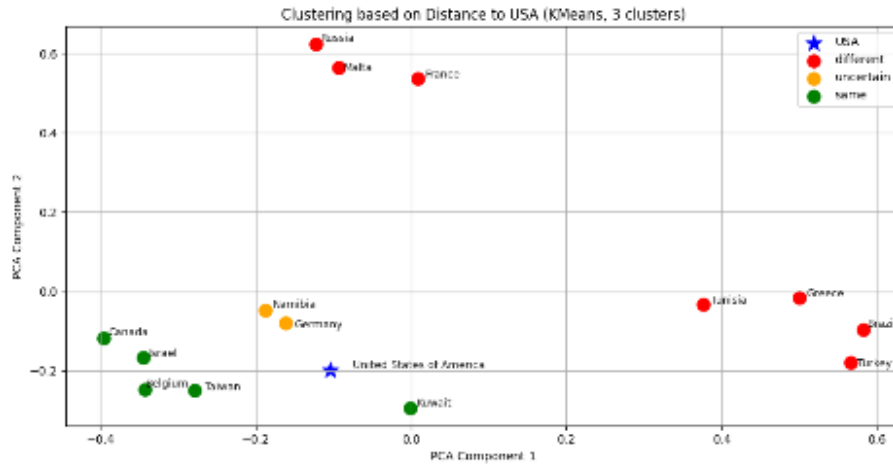


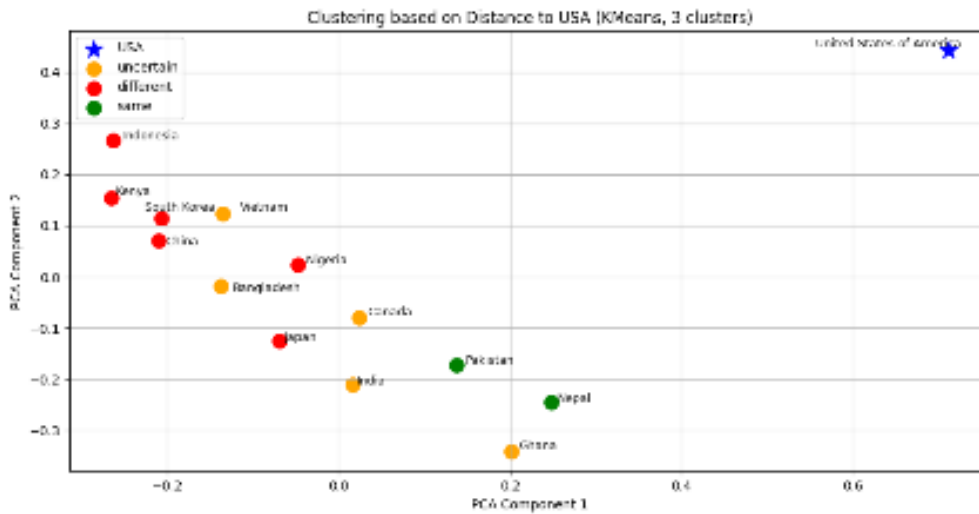
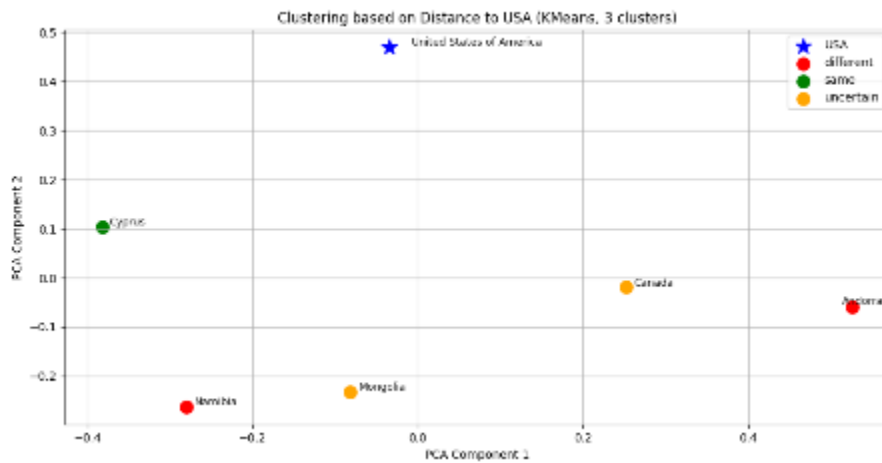
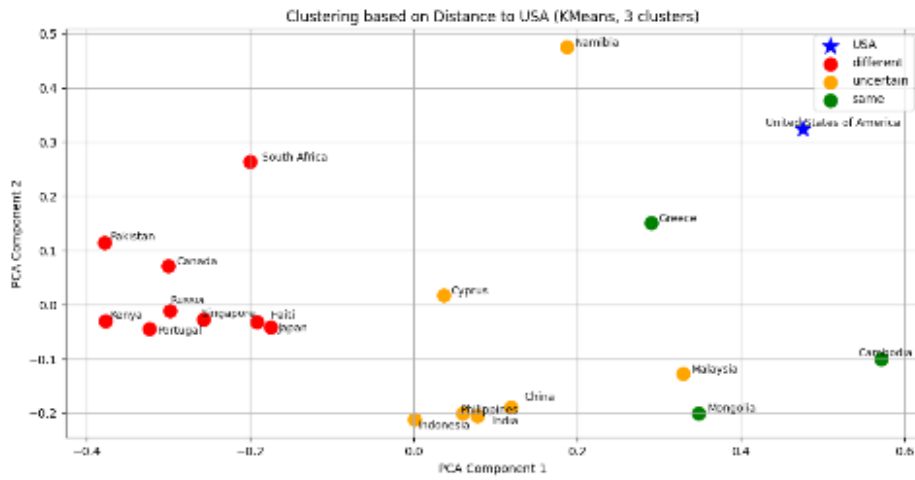


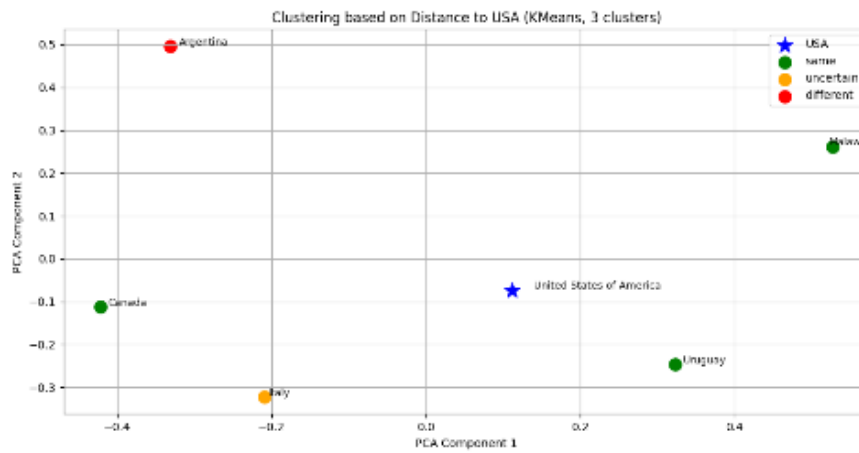
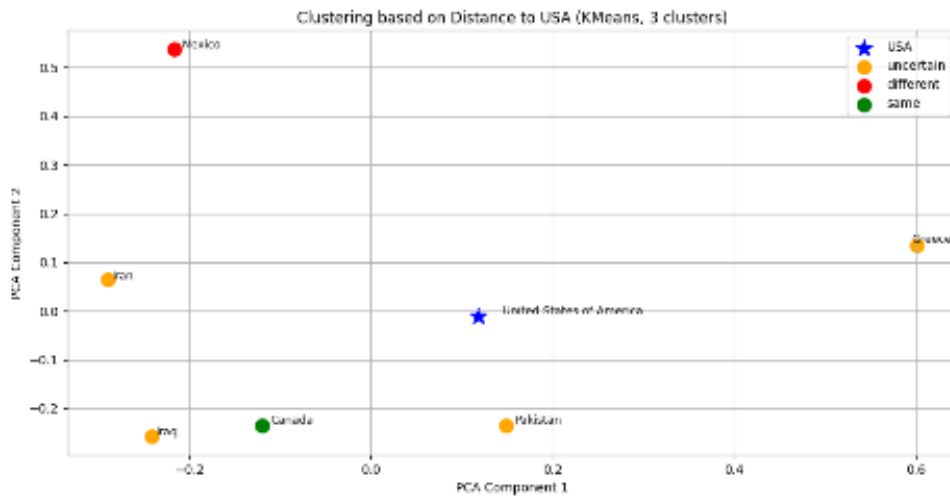
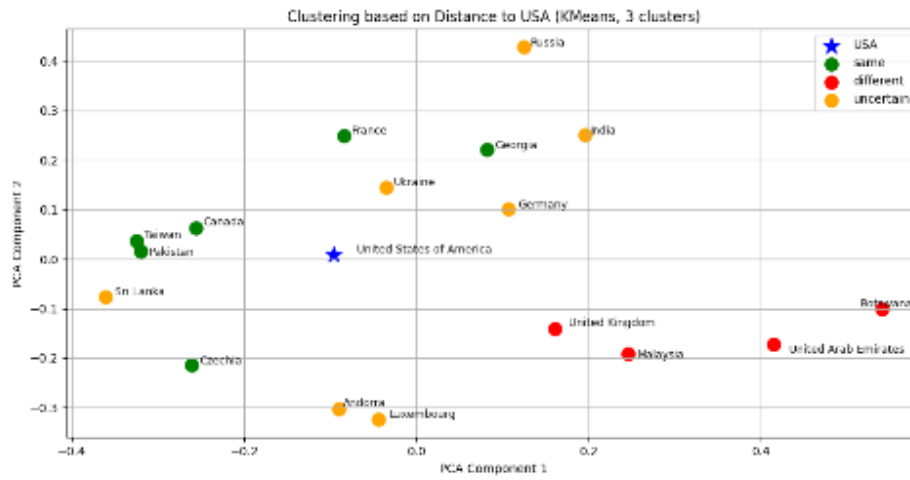


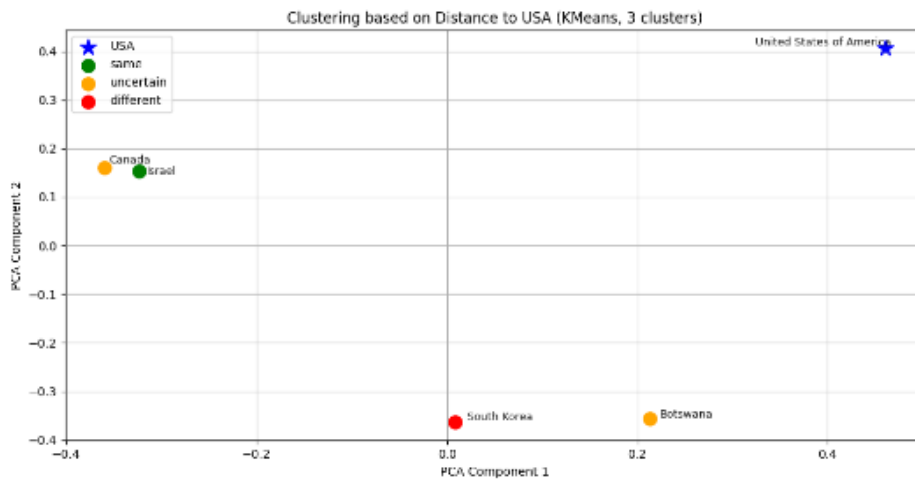
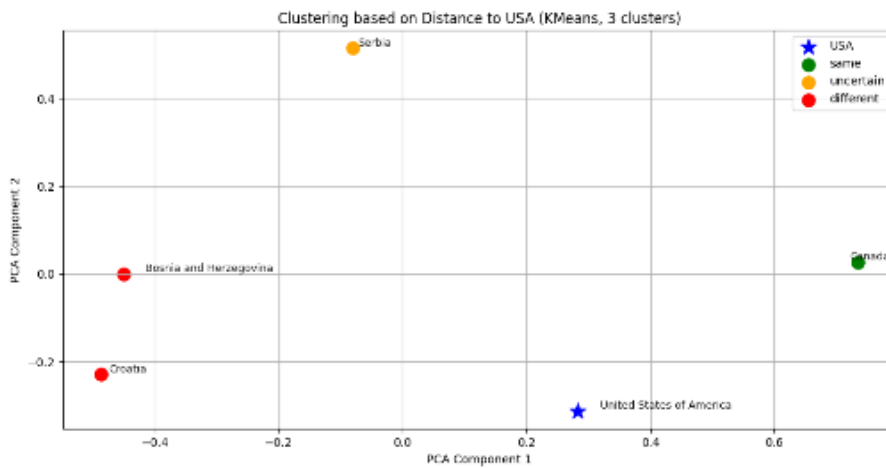
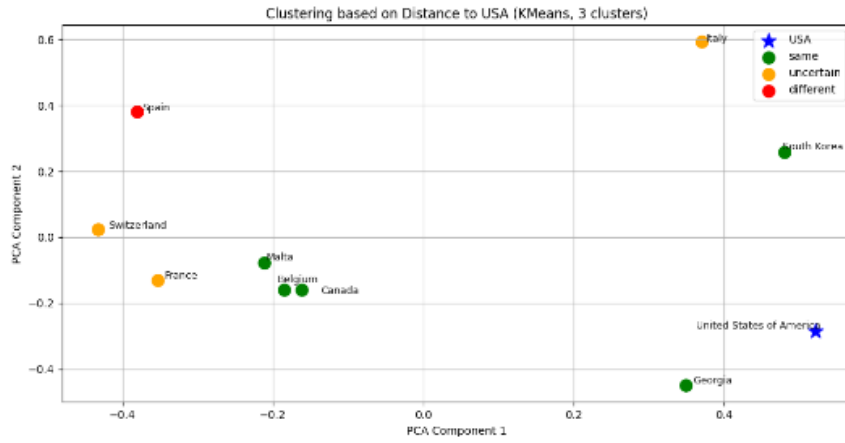


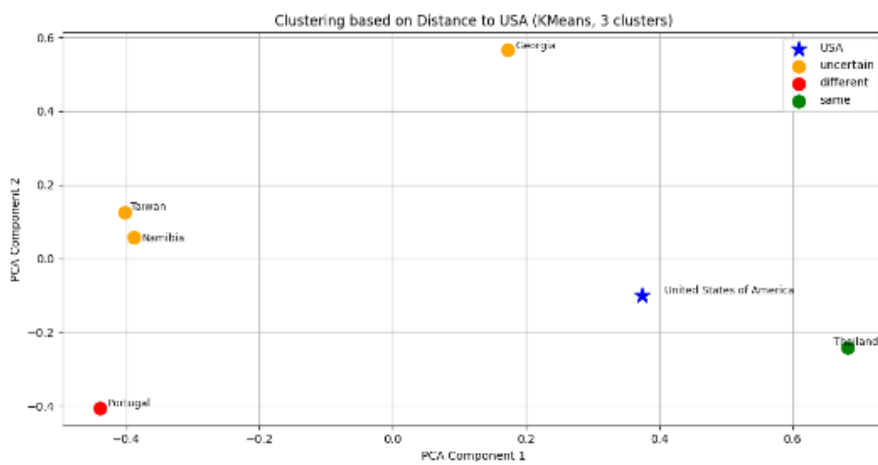
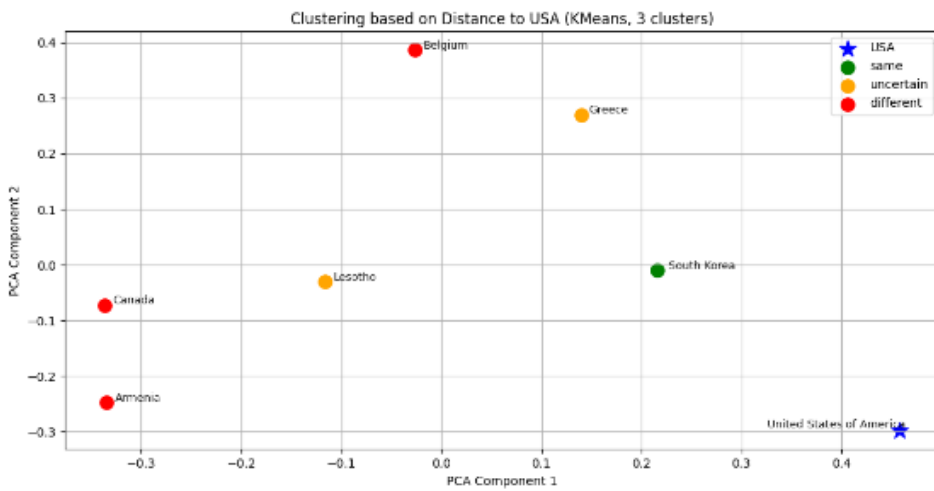
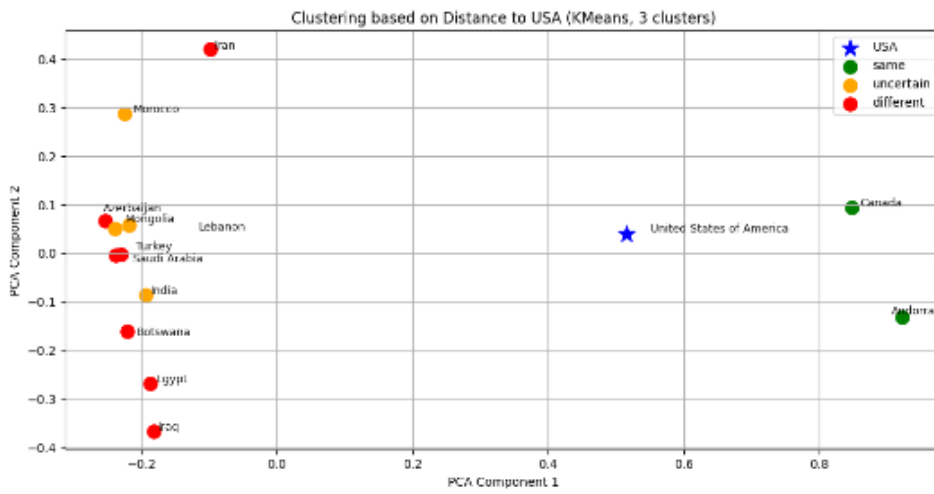
GPT-5. 1 英語プロンプト 代表的意味を出力

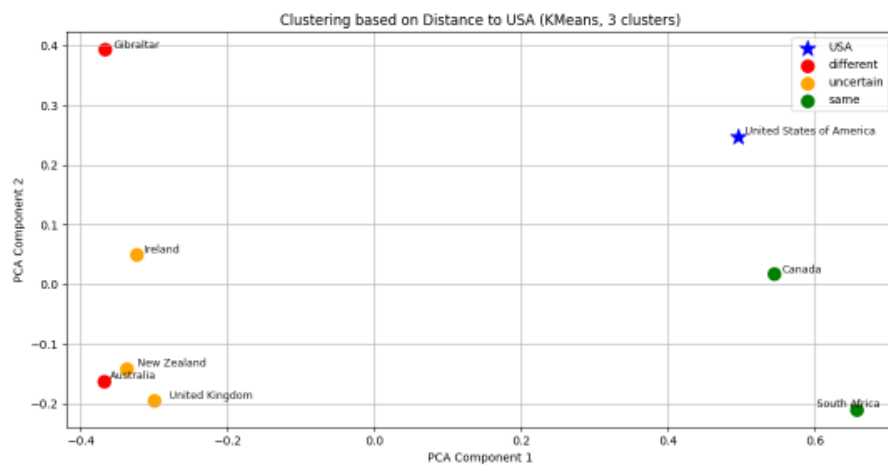
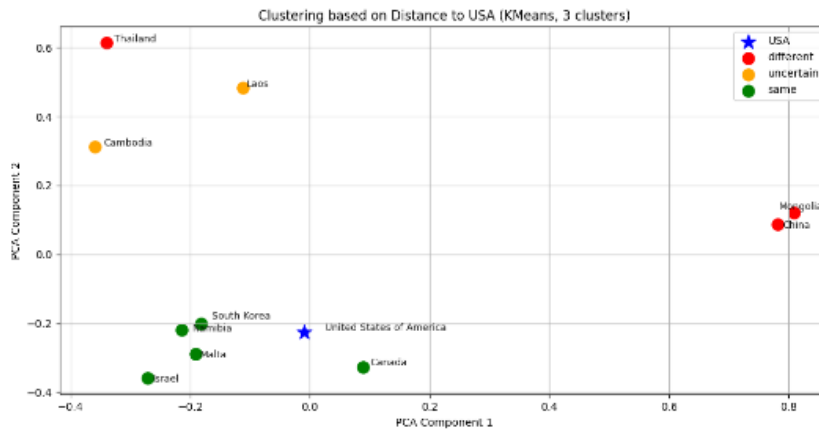
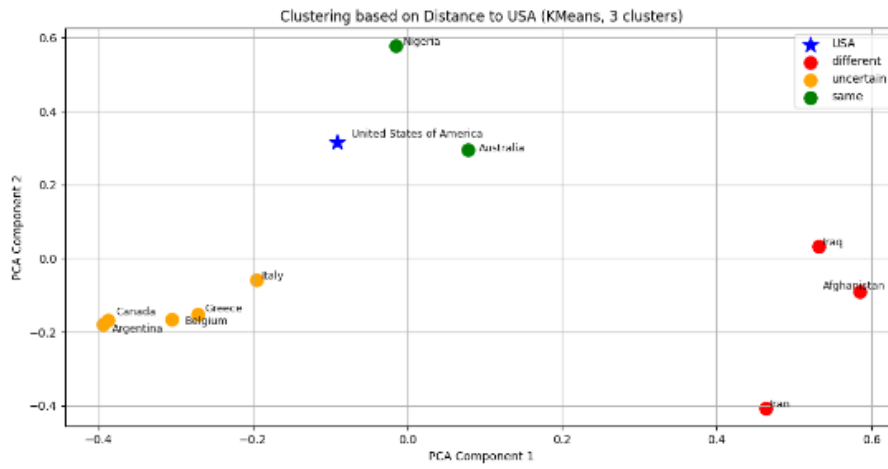


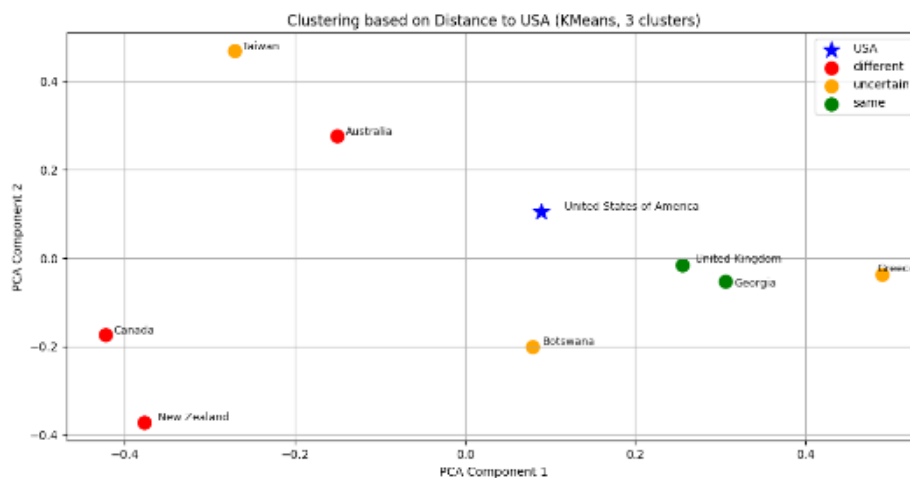




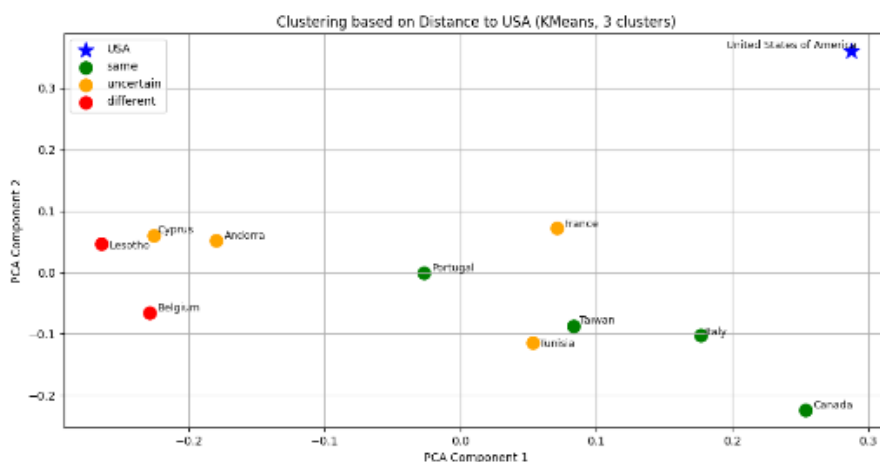
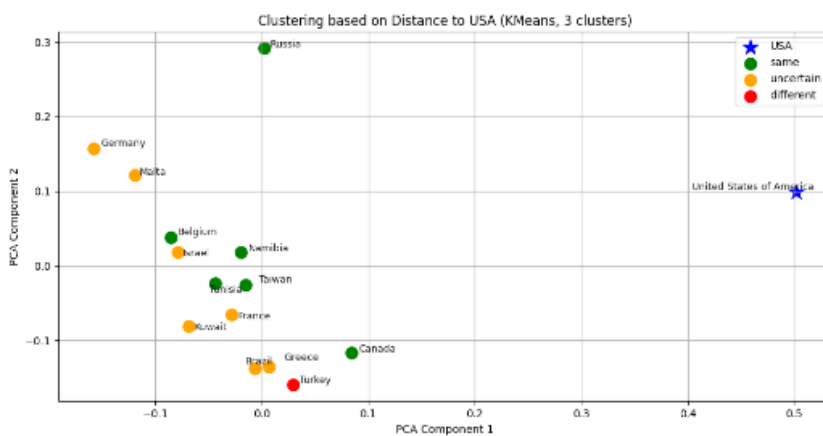


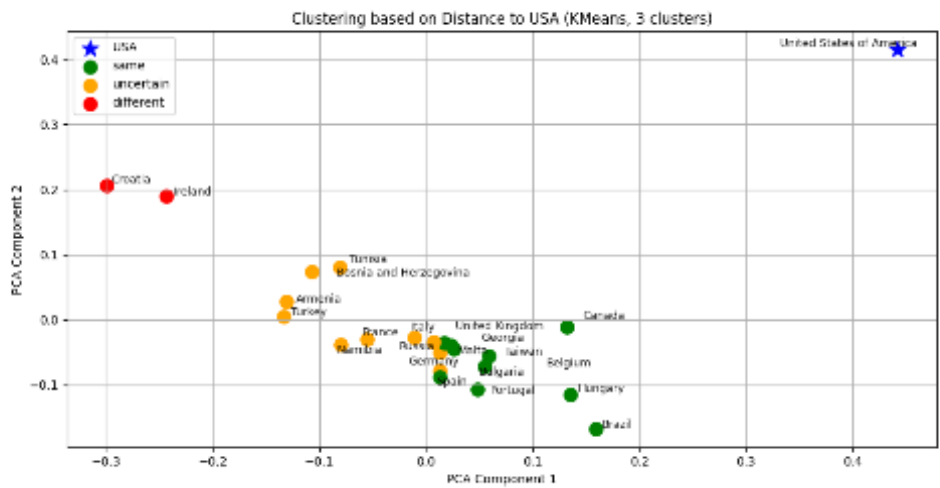
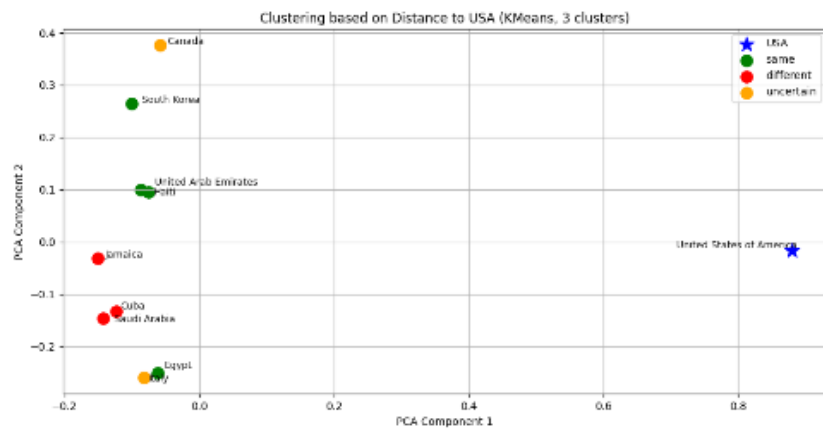
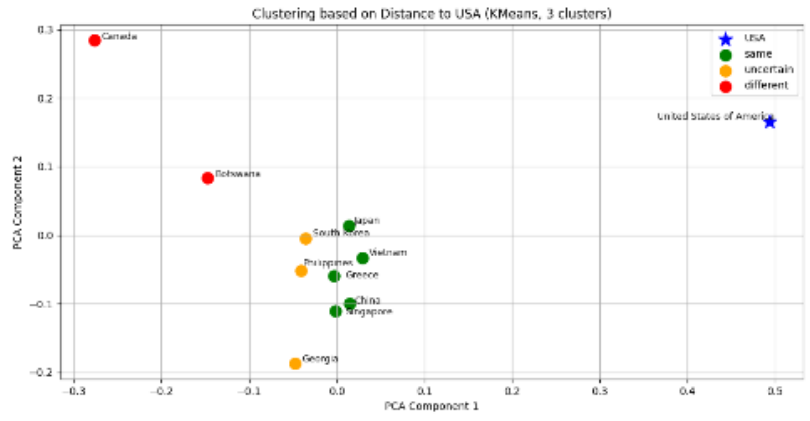


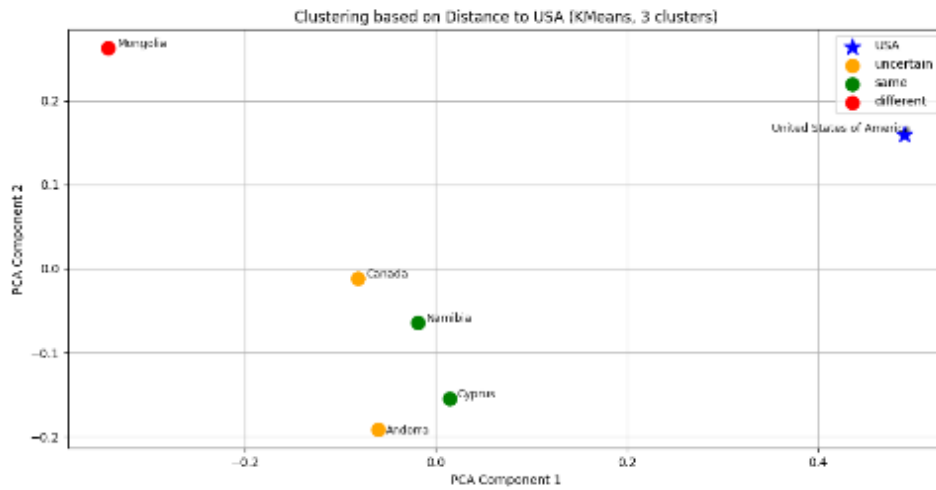
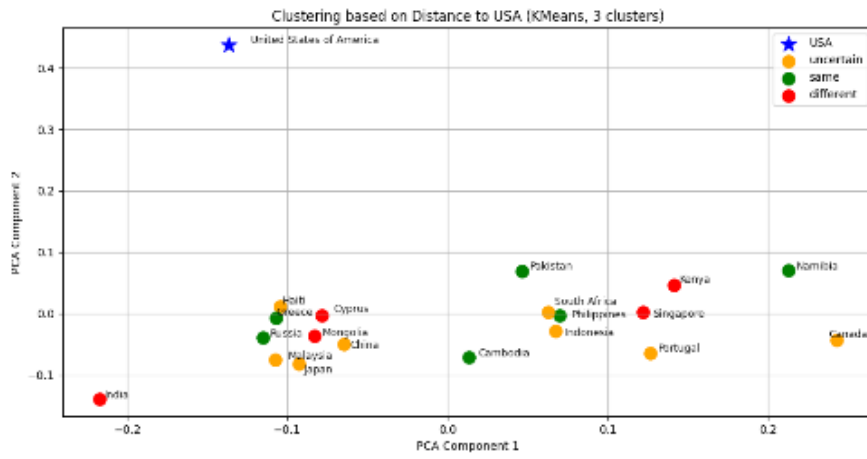
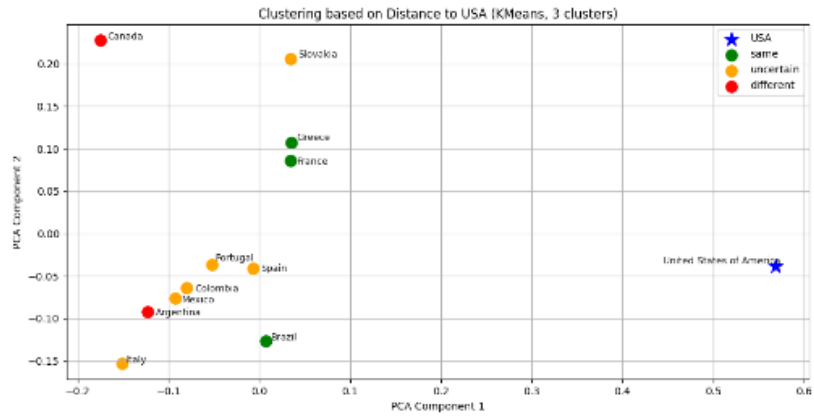


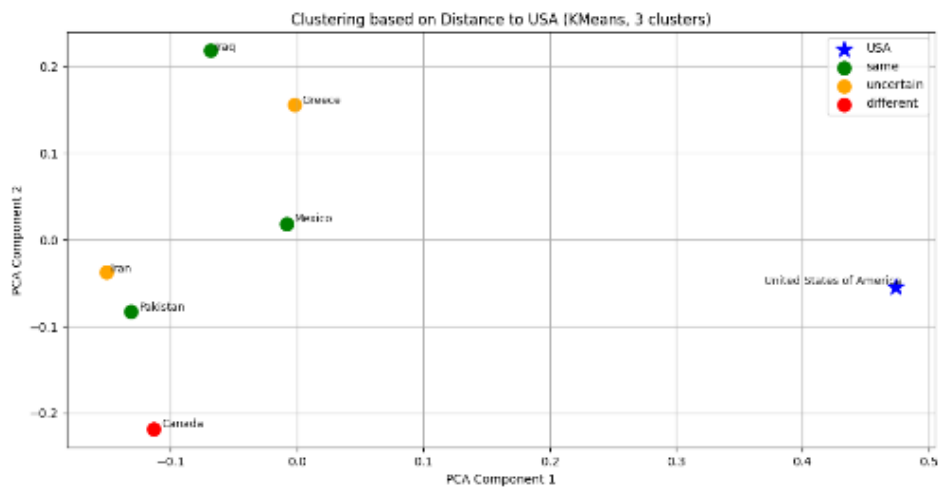
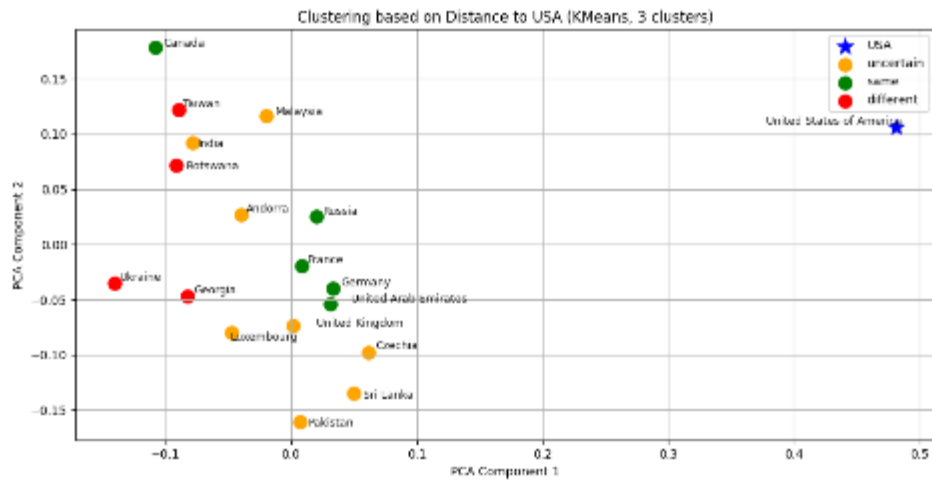
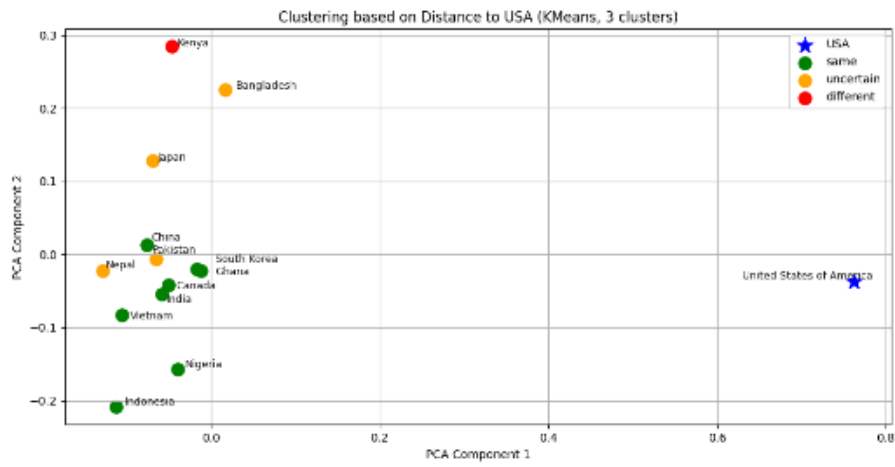


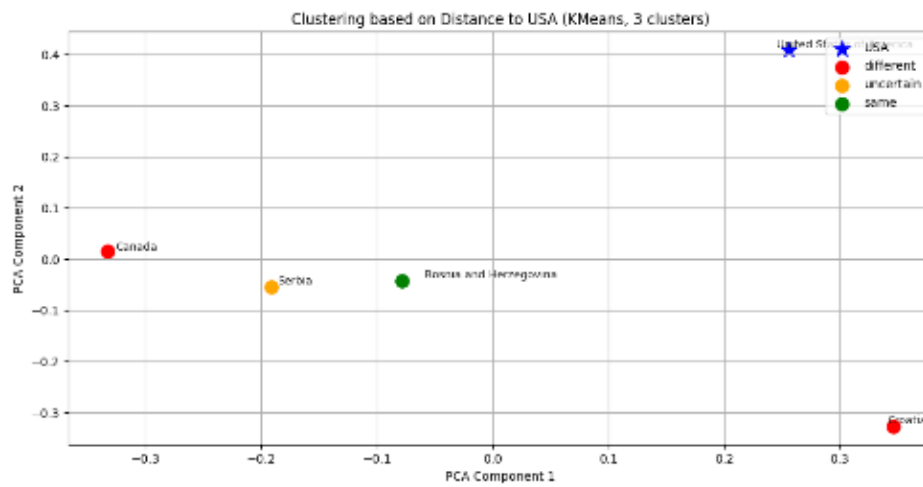
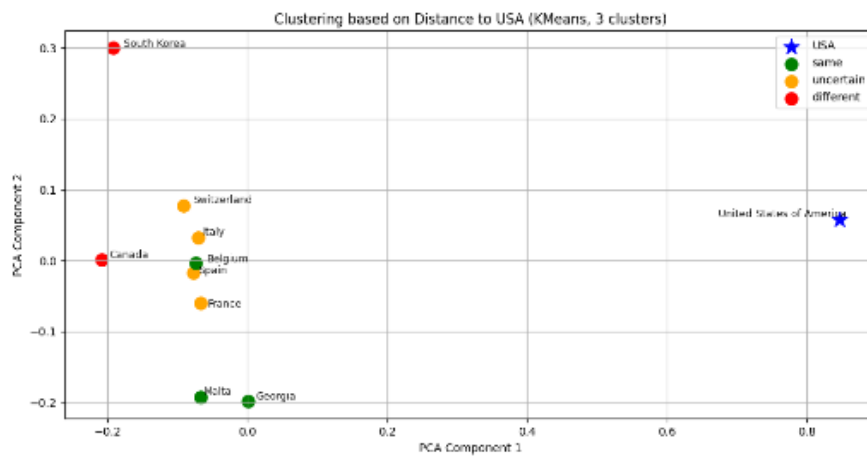
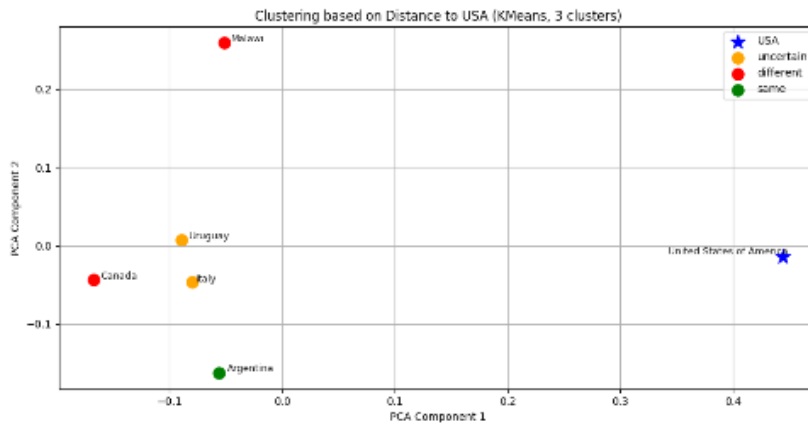
GPT-4. 1mini 公用語プロンプト 全ての意味を出力

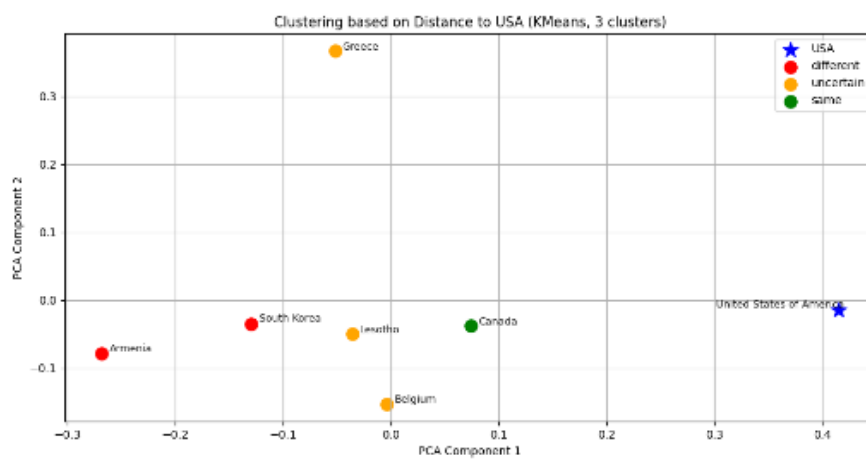
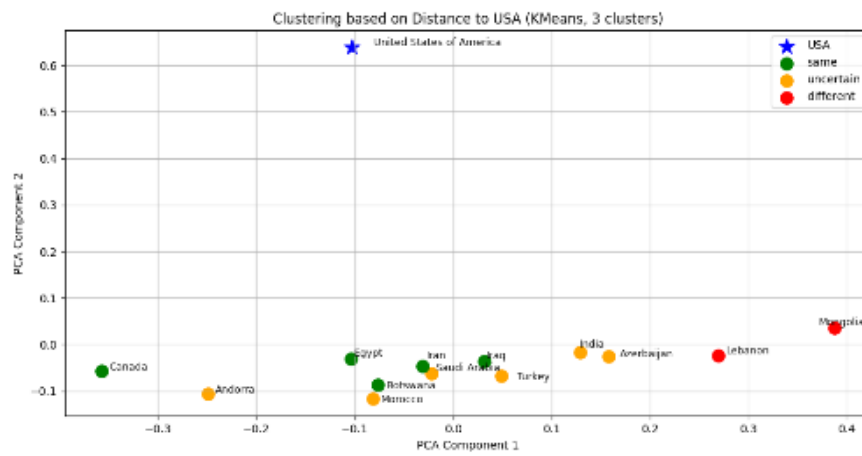
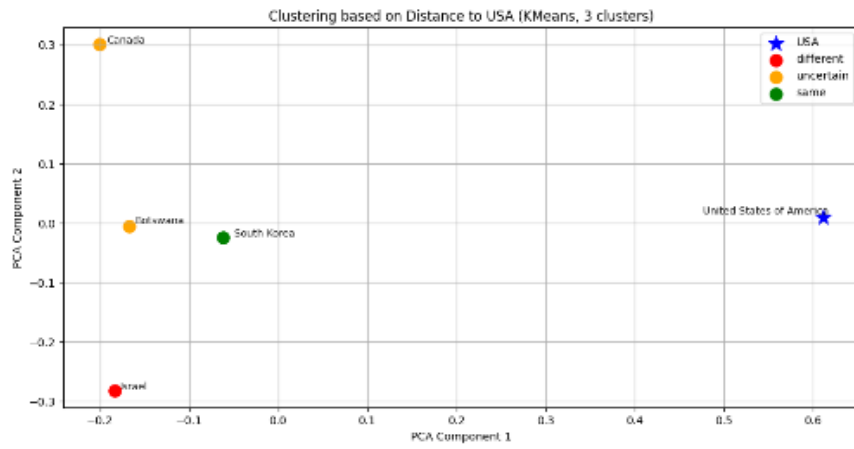


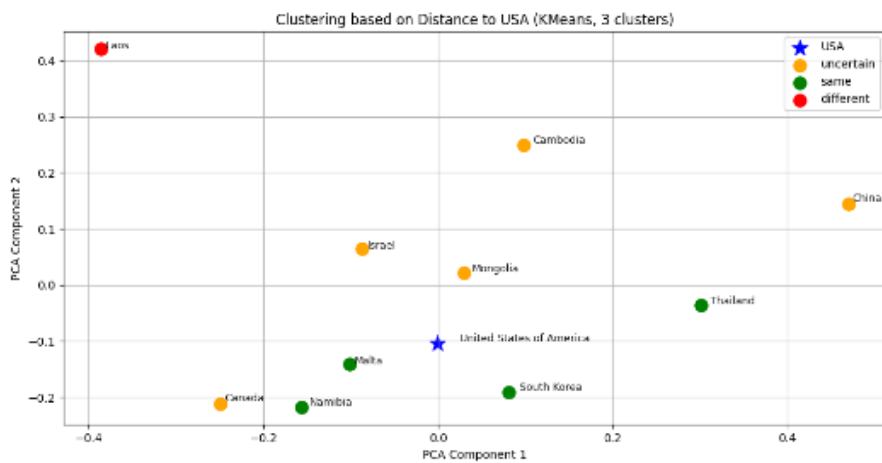
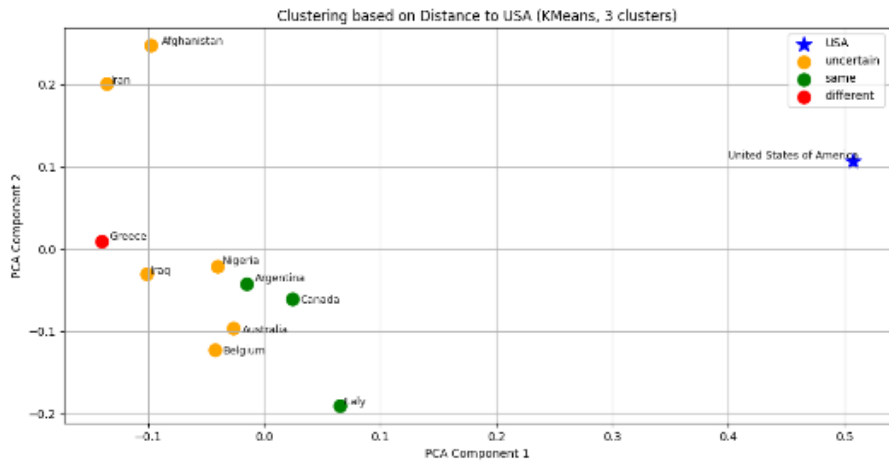
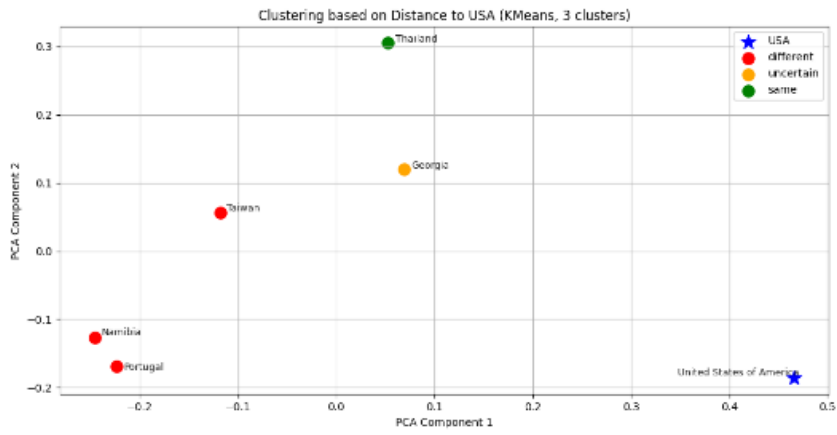


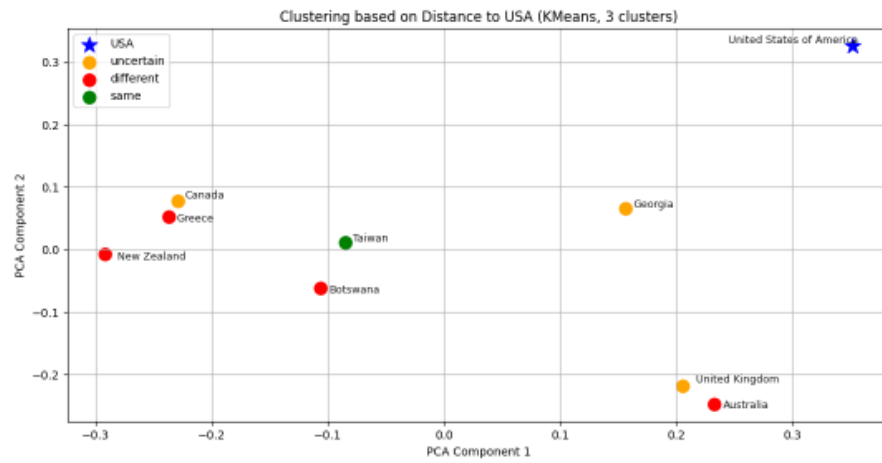
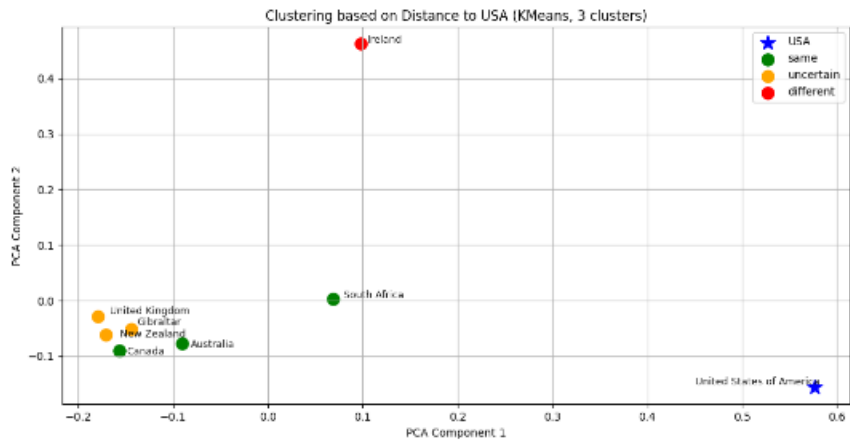




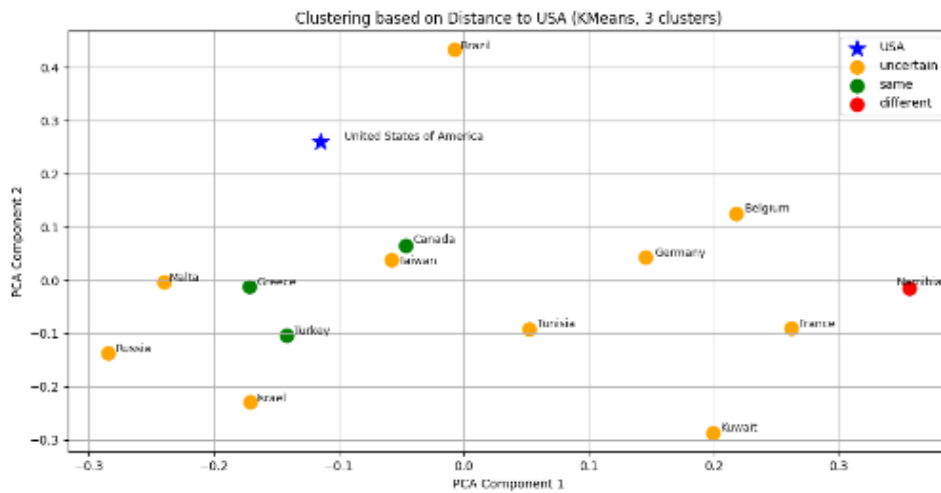


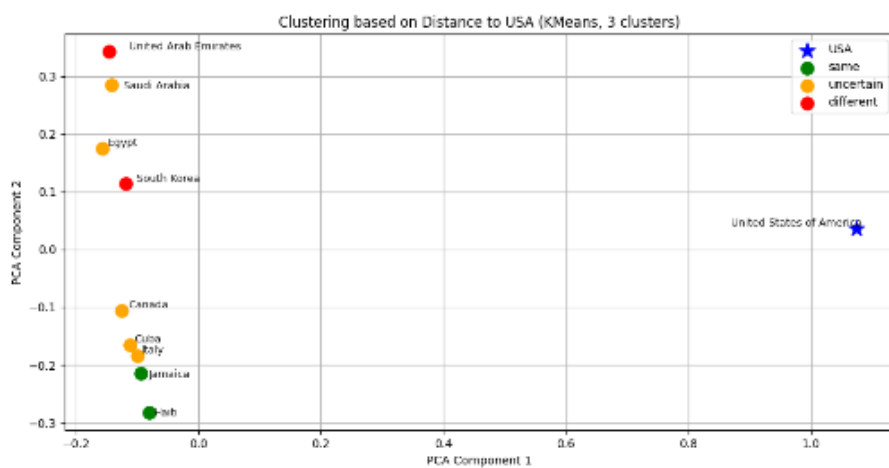
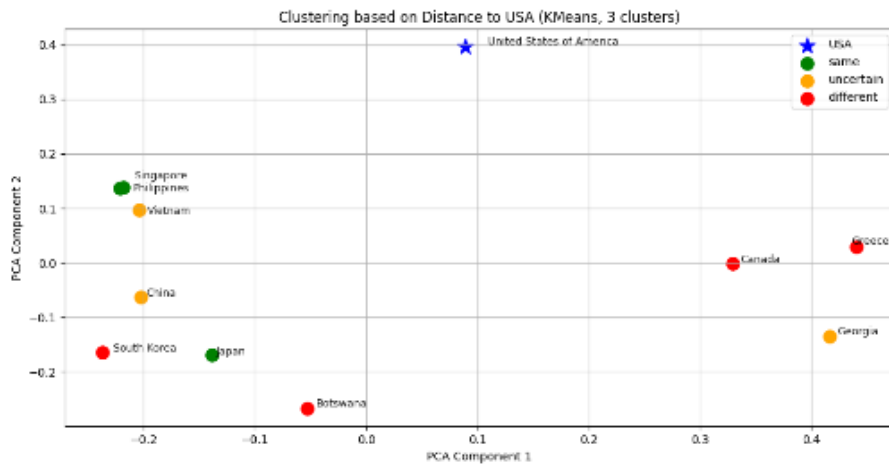
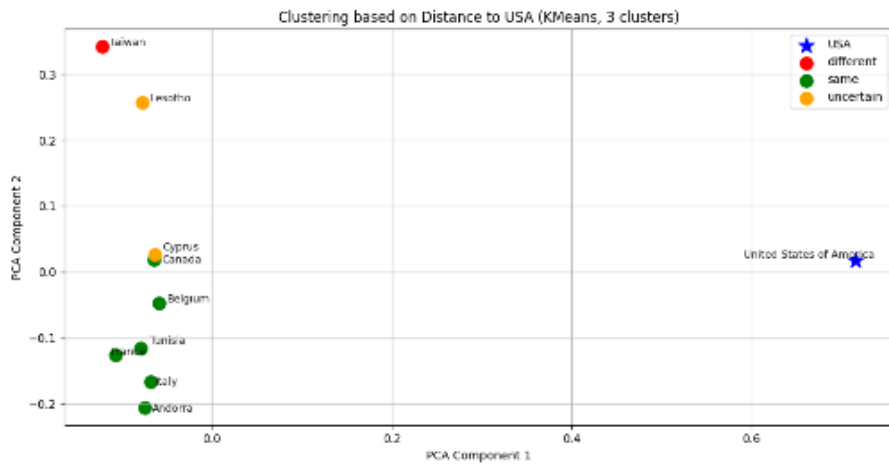


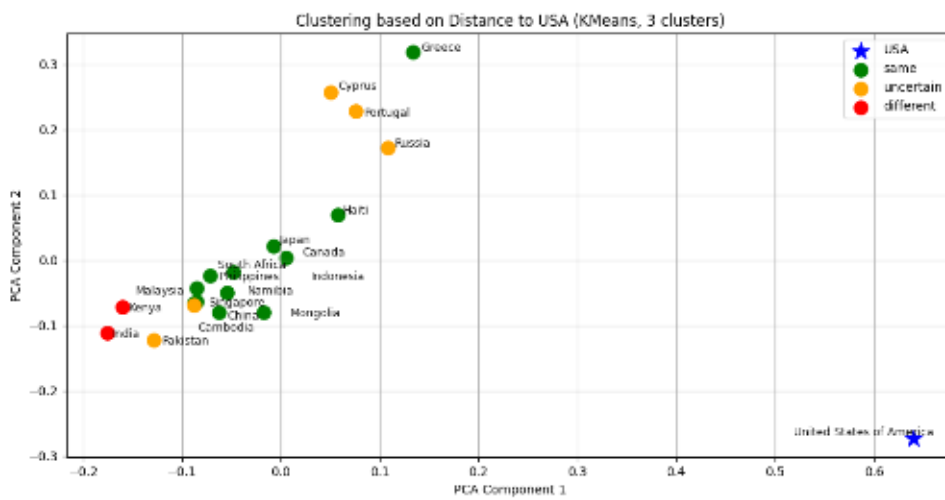
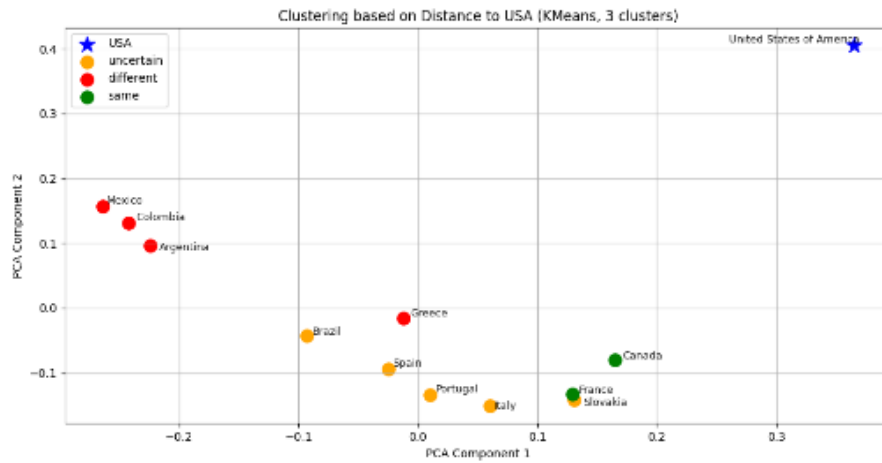
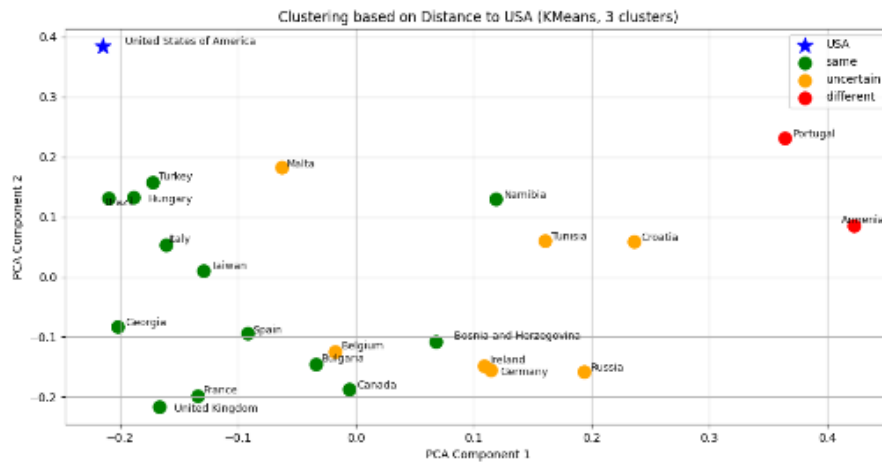


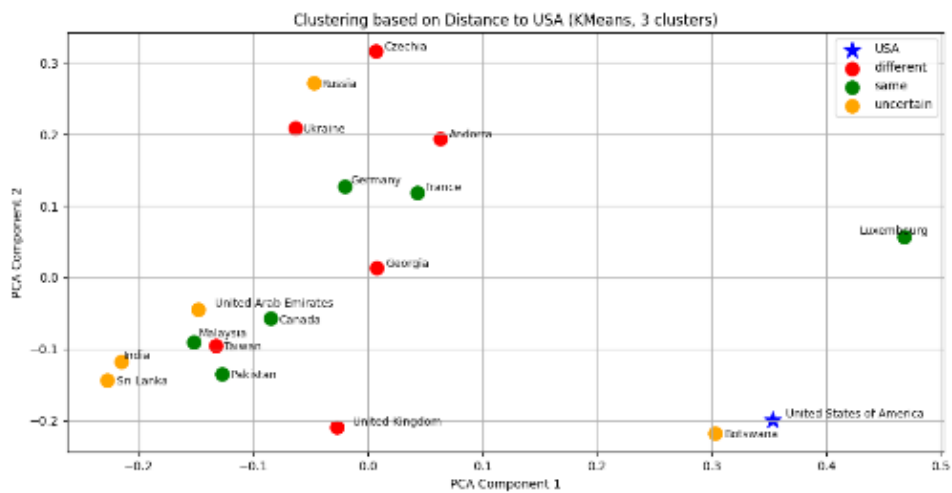
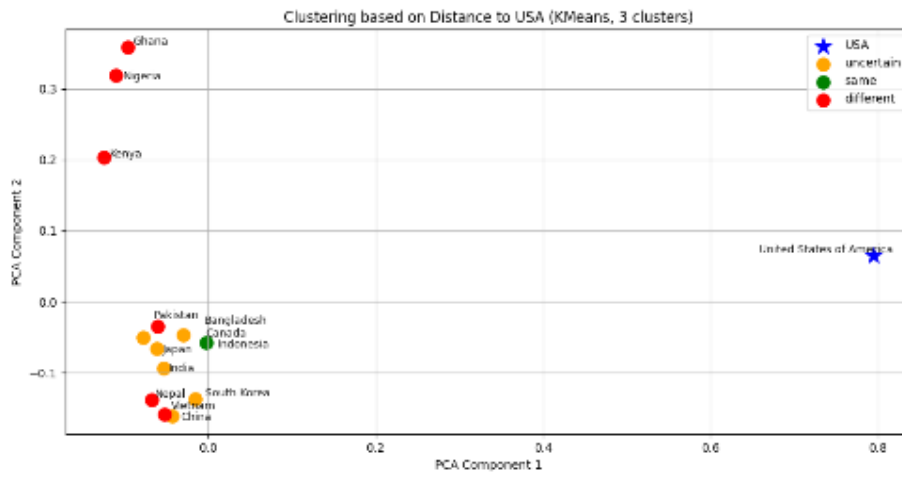
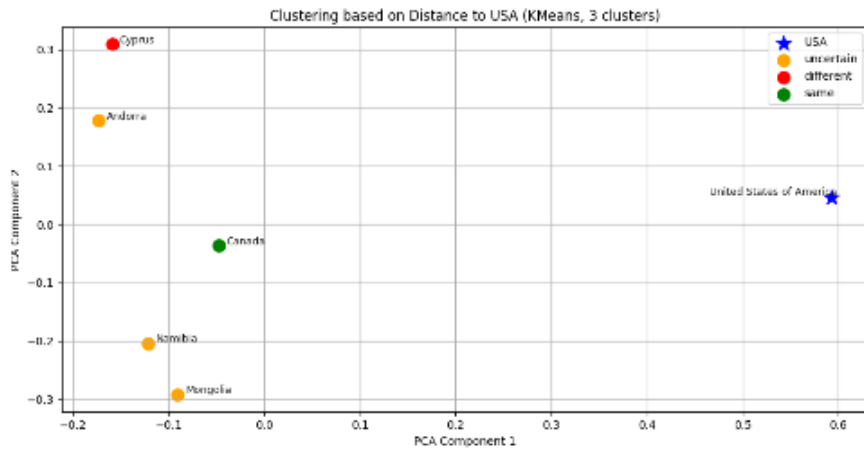


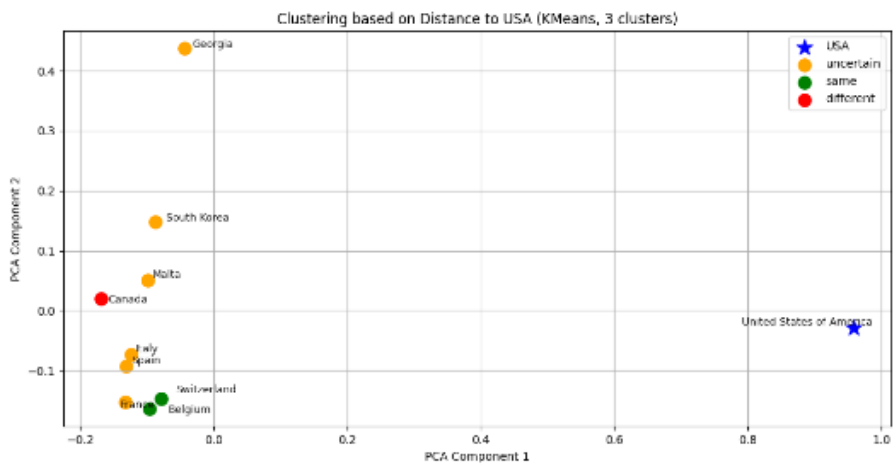
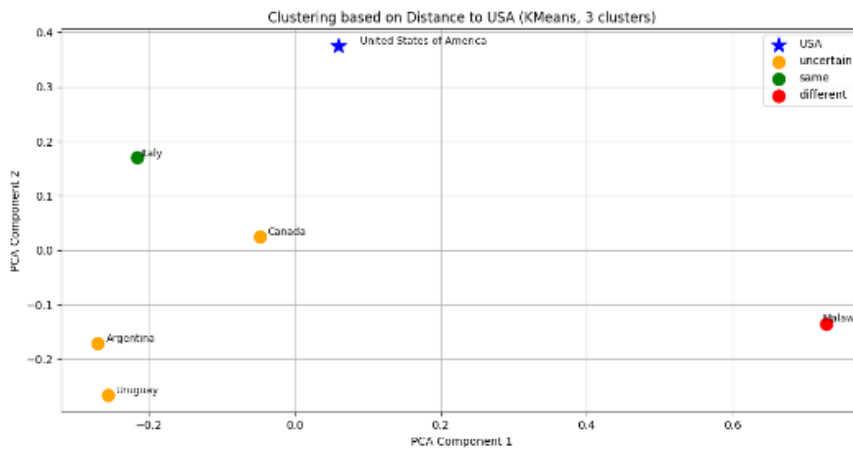
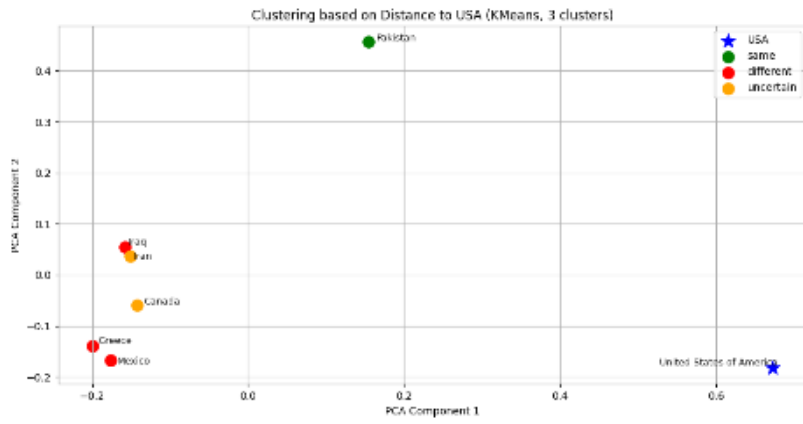
GPT-4. 1mini 公用語プロンプト 代表的の意味を出力

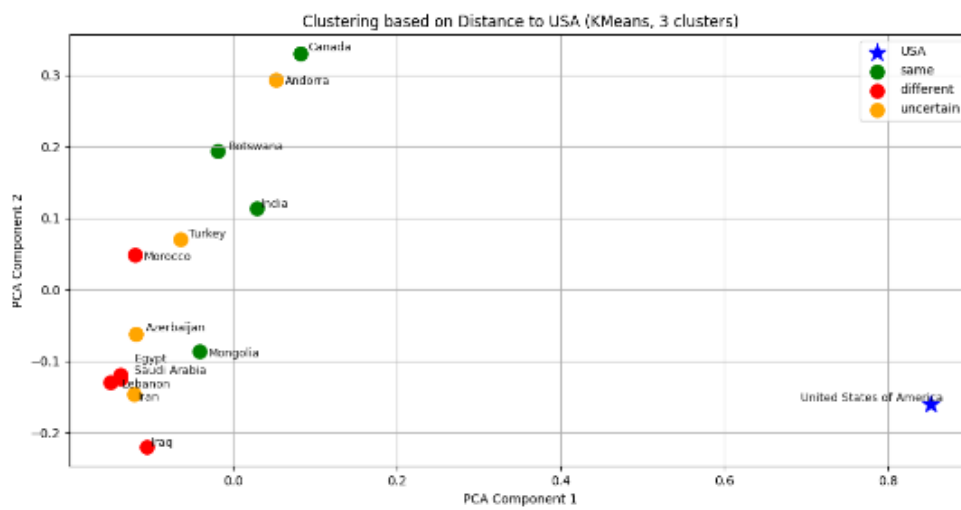
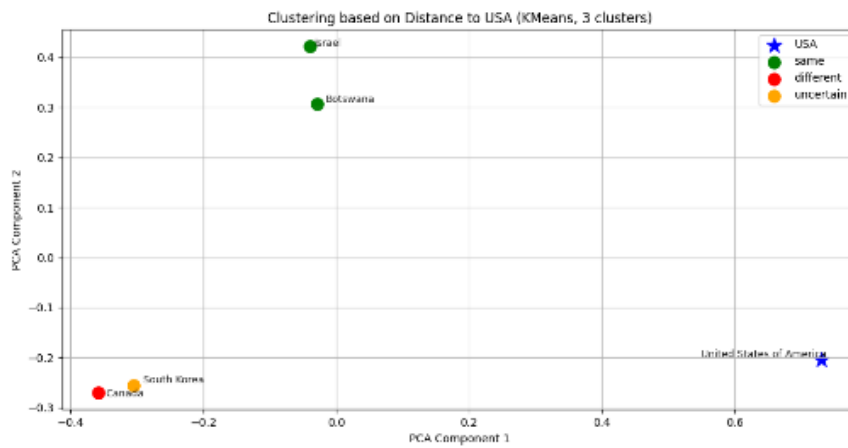
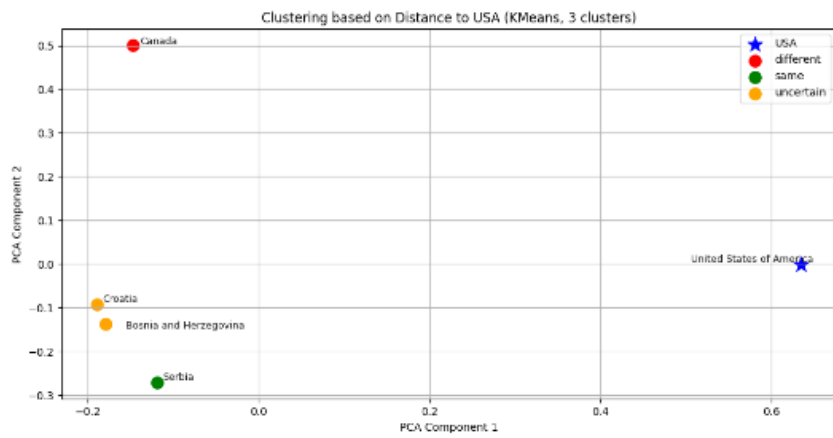


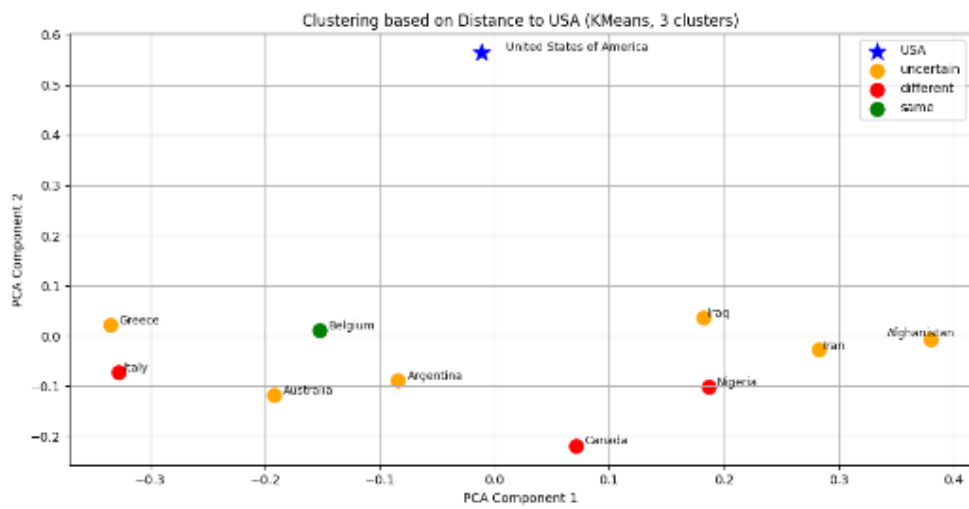
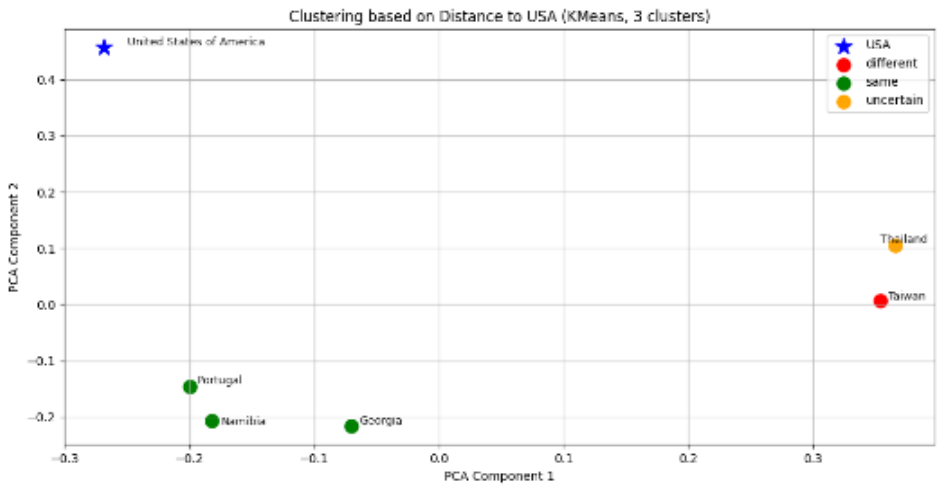
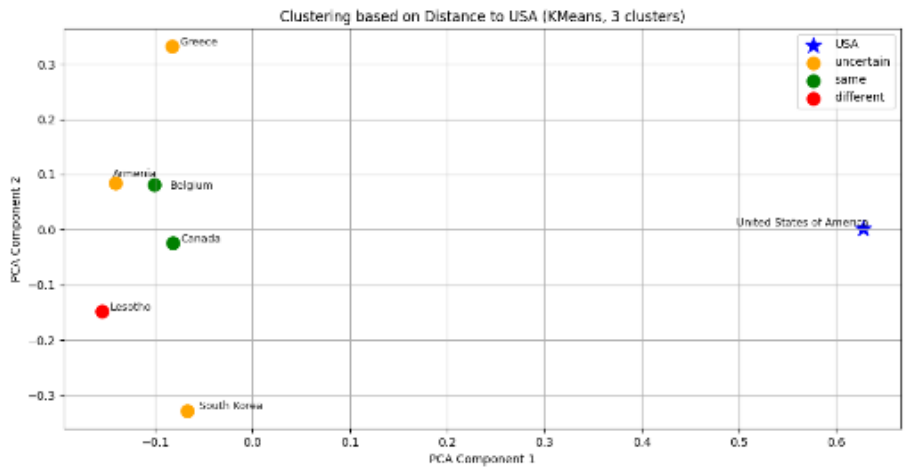


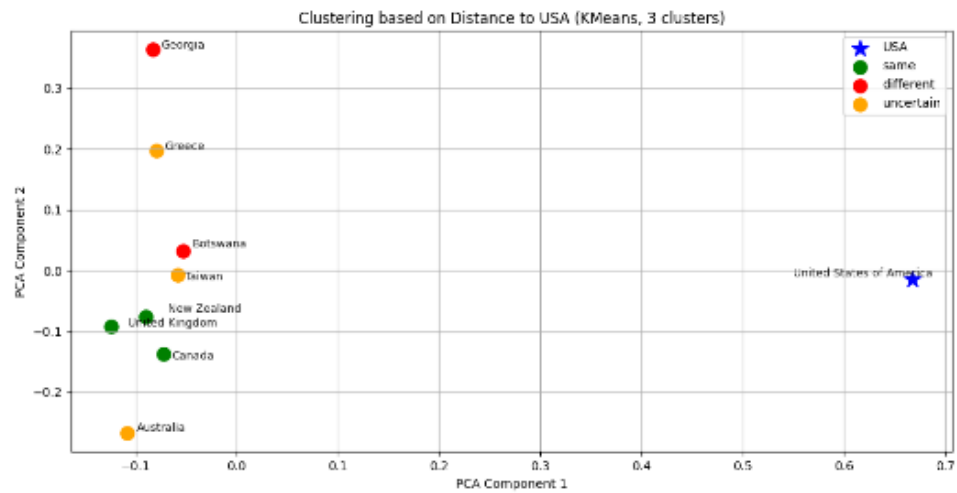
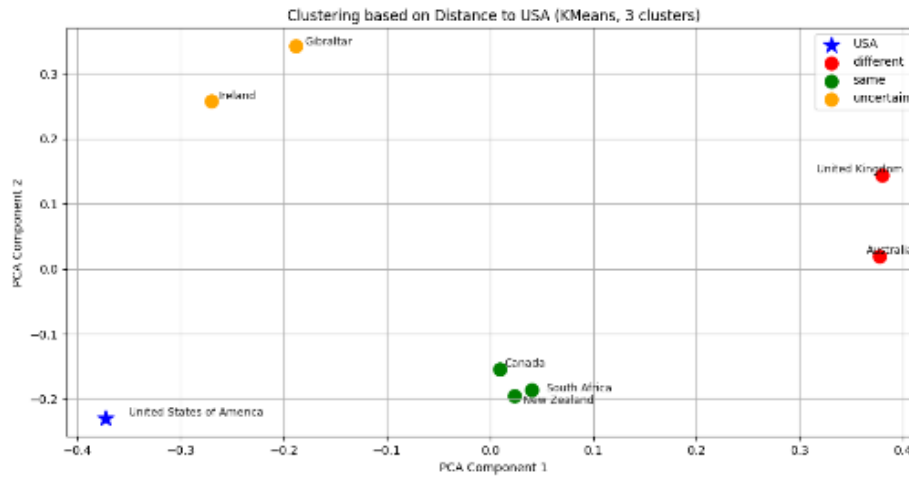
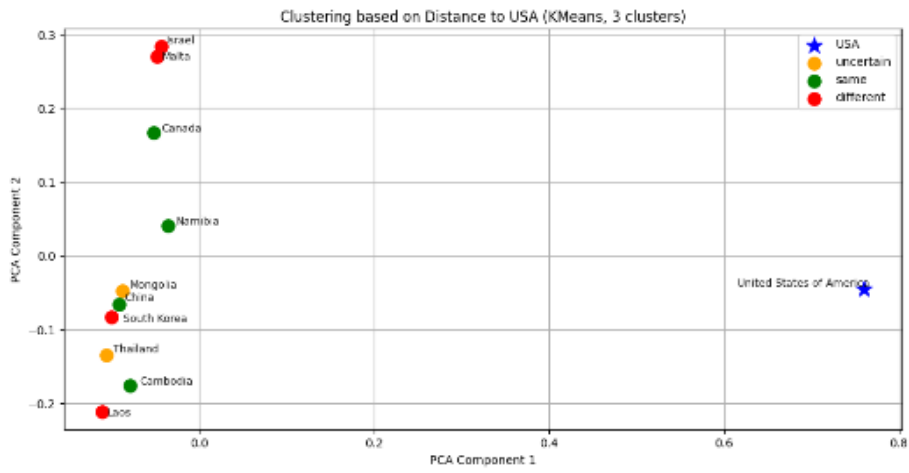




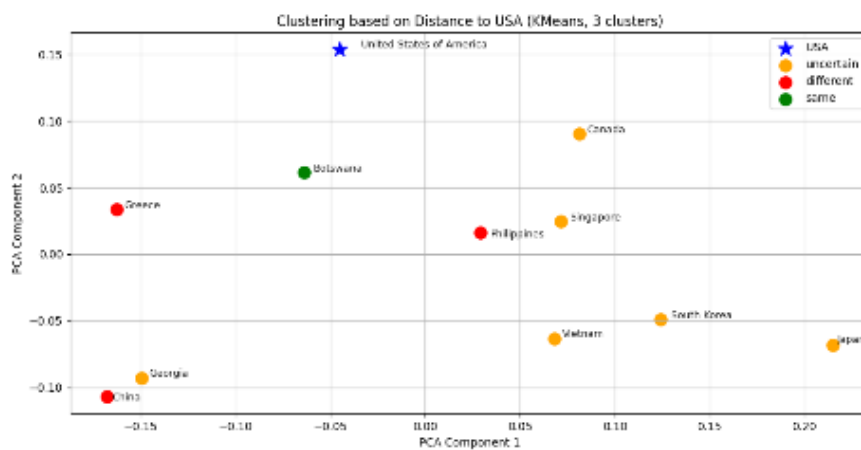
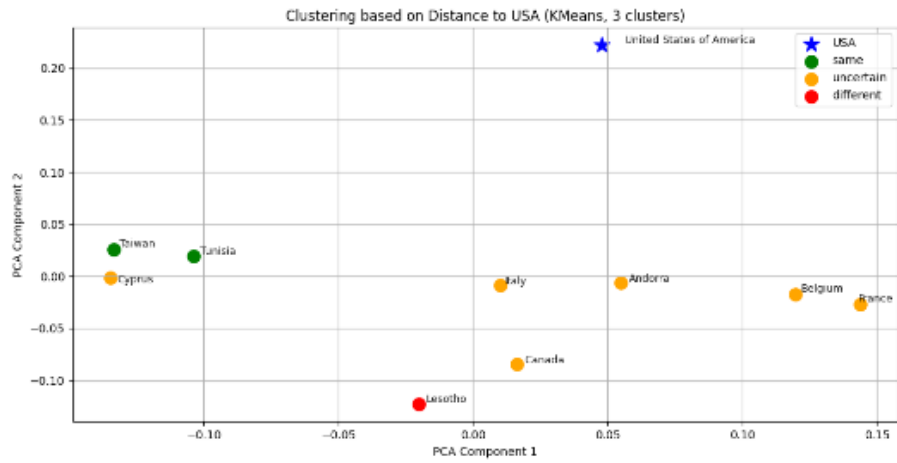
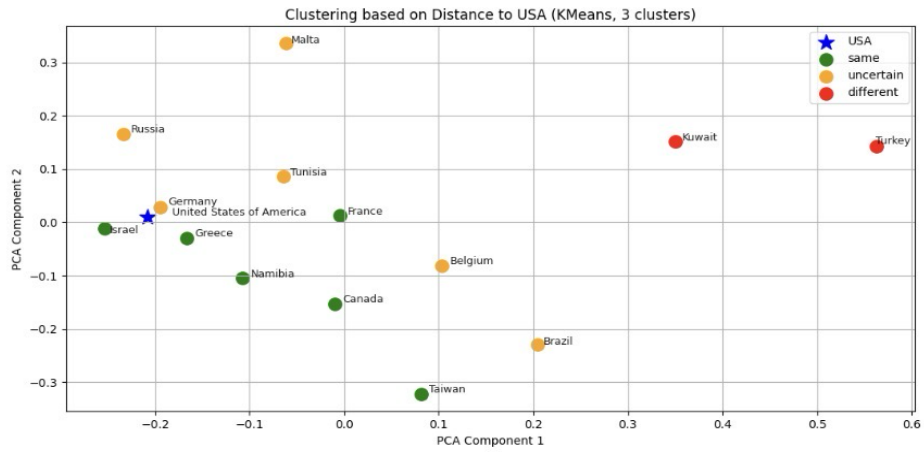


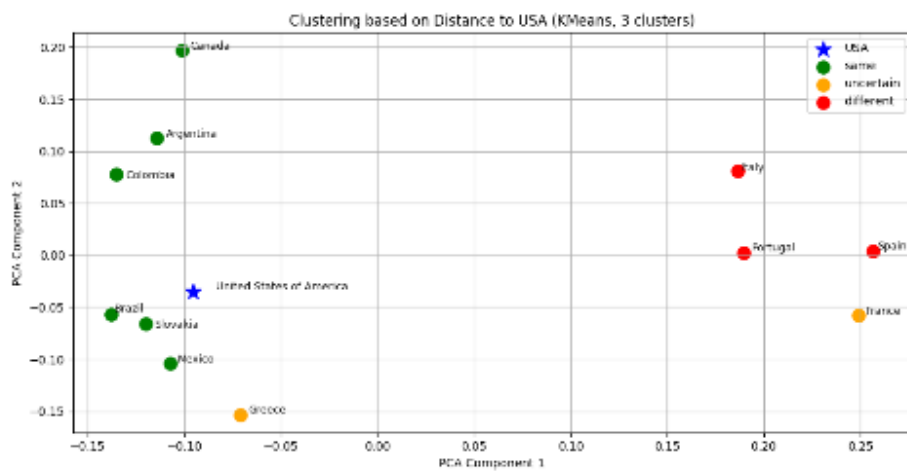
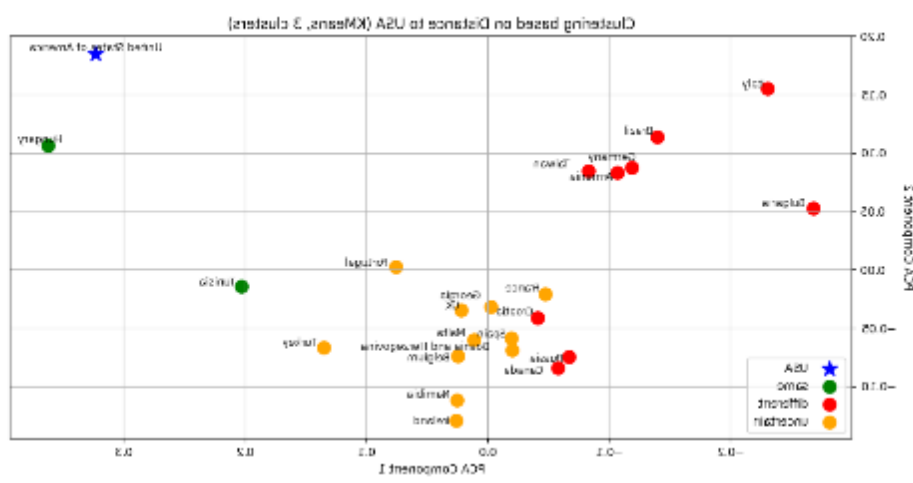
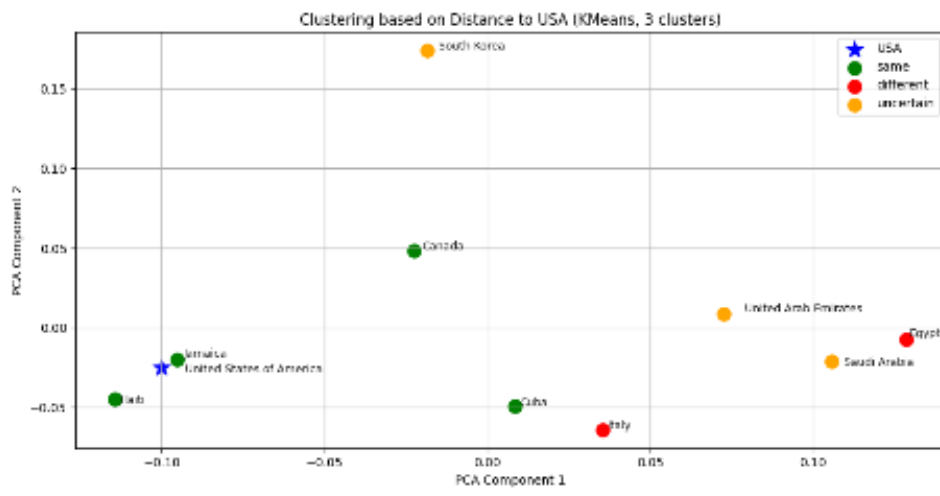


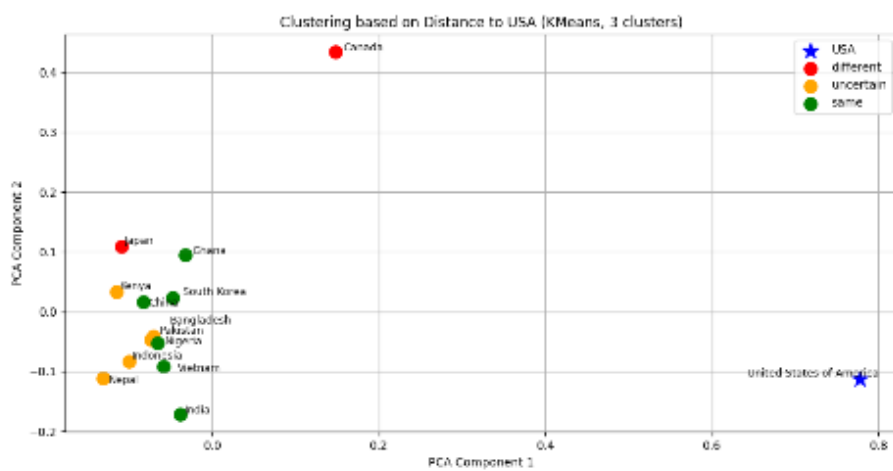
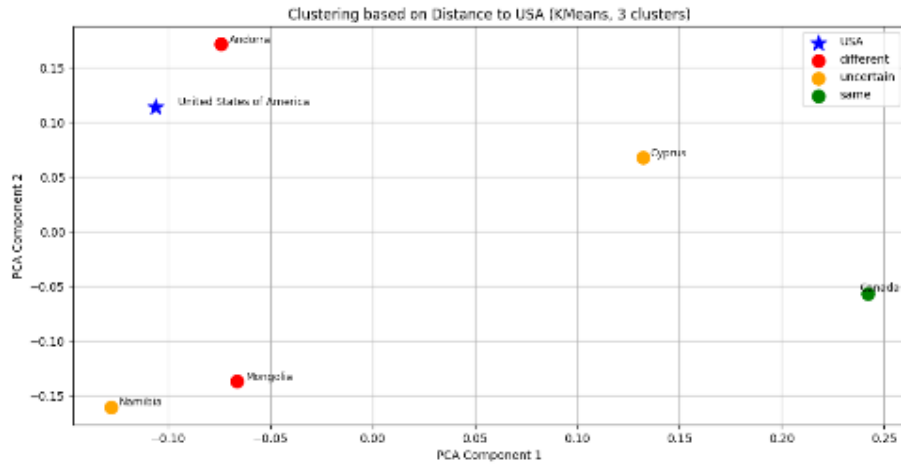
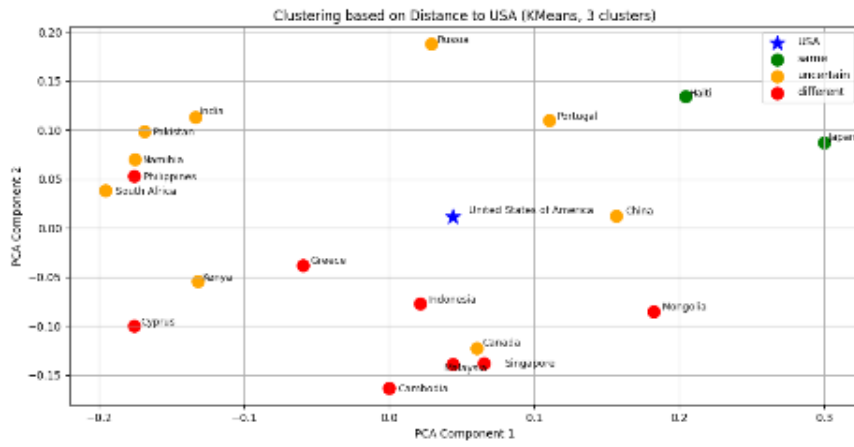


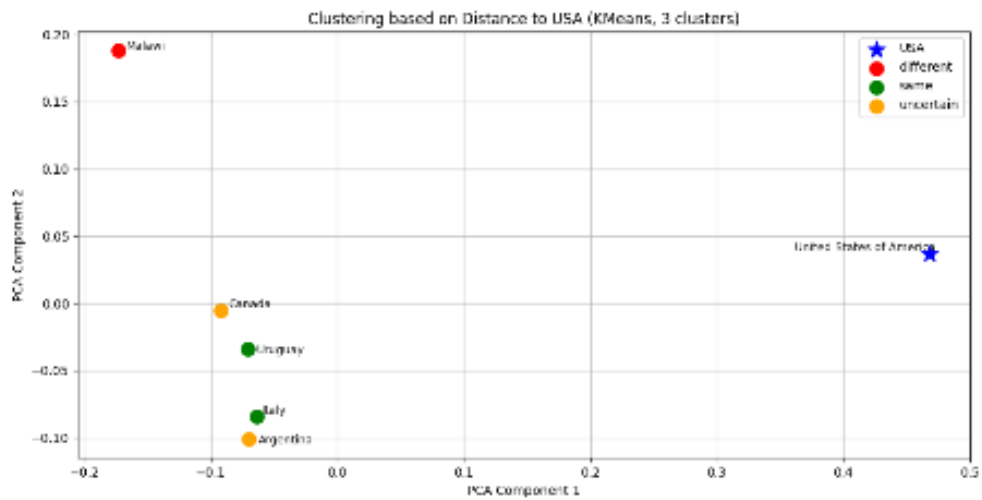
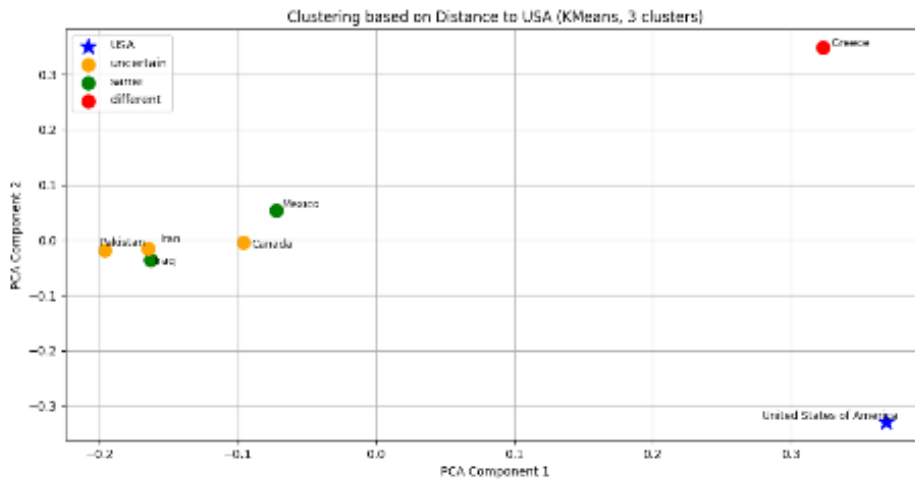
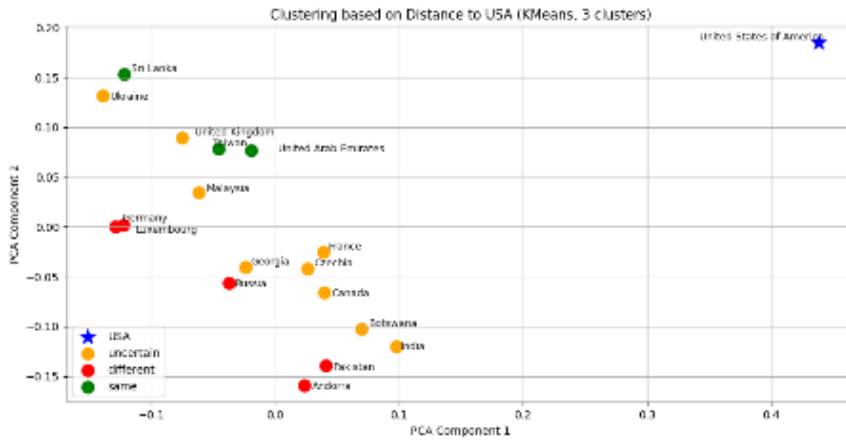


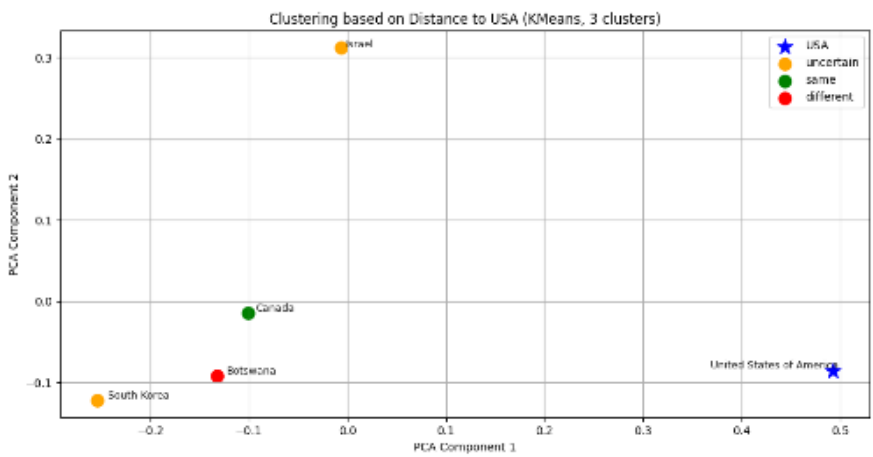
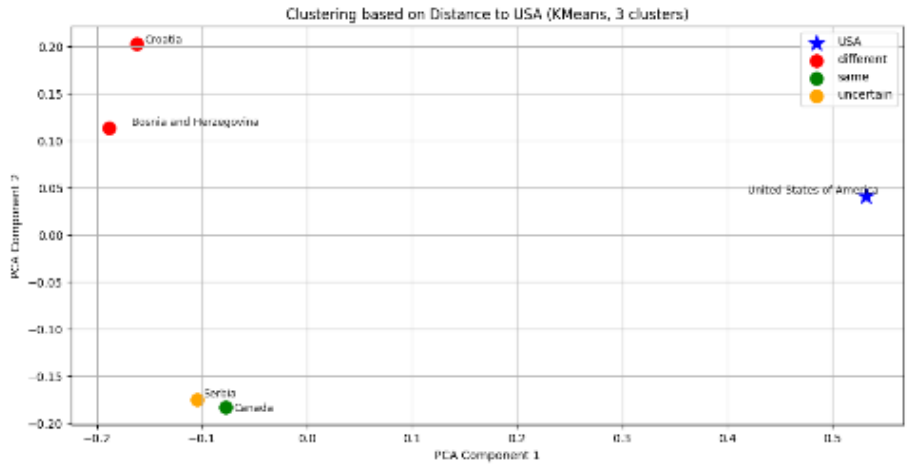
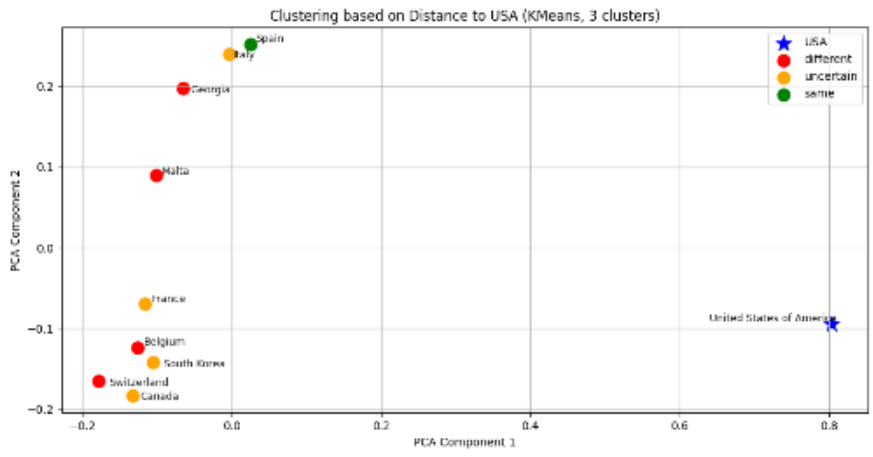
GPT-5. 1 公用語プロンプト 全ての意味を出力

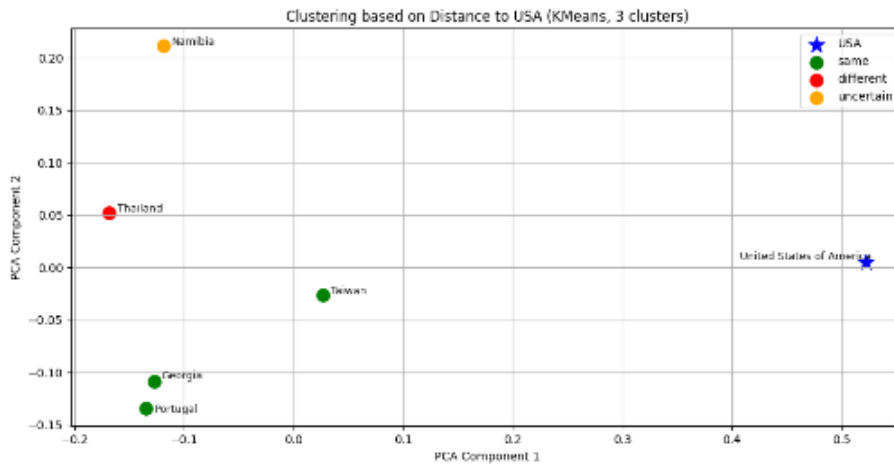
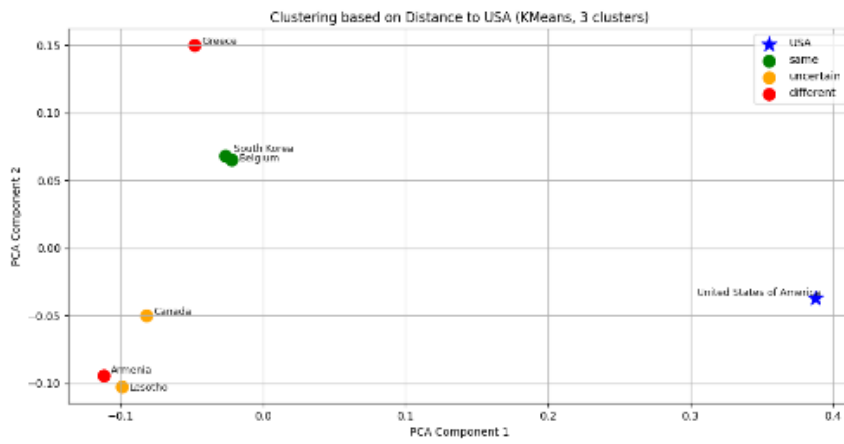
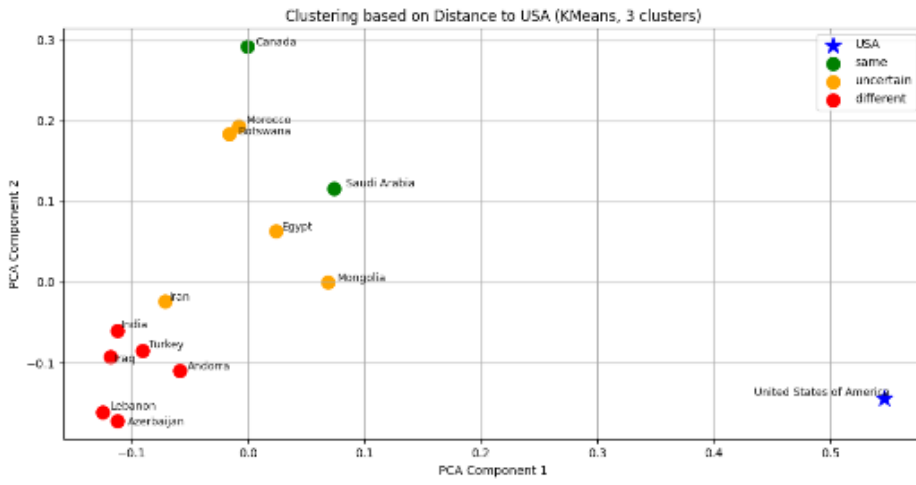


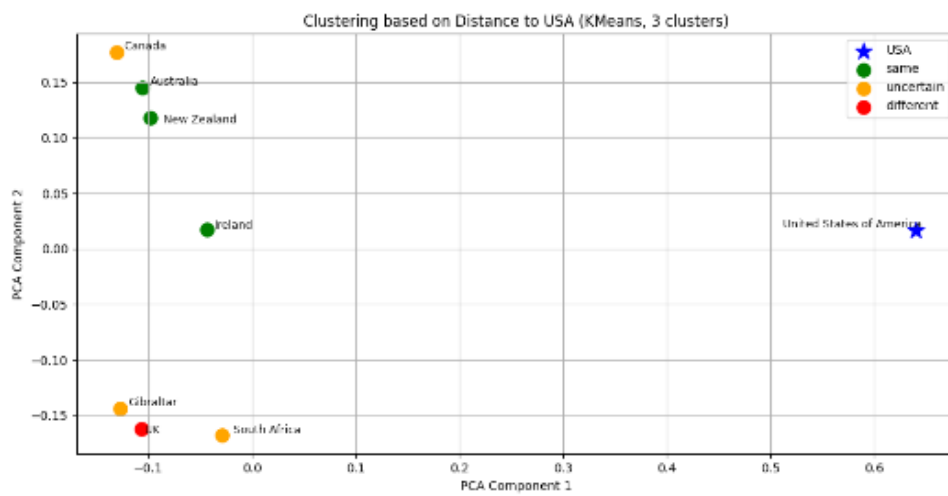
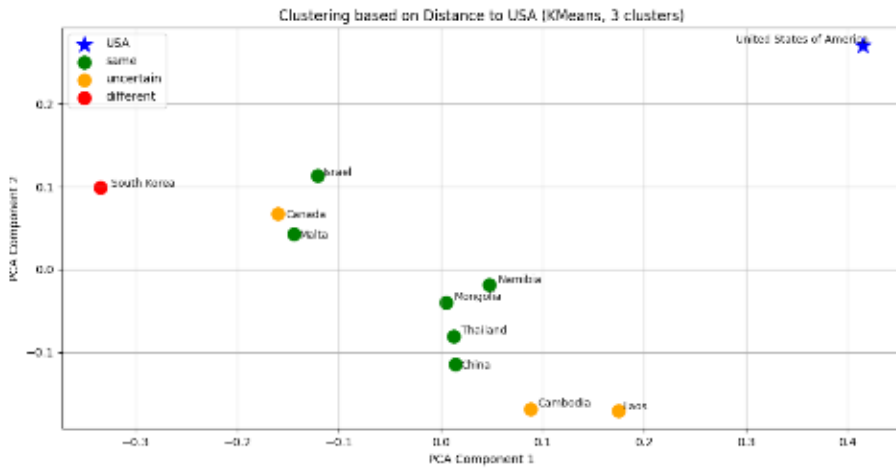
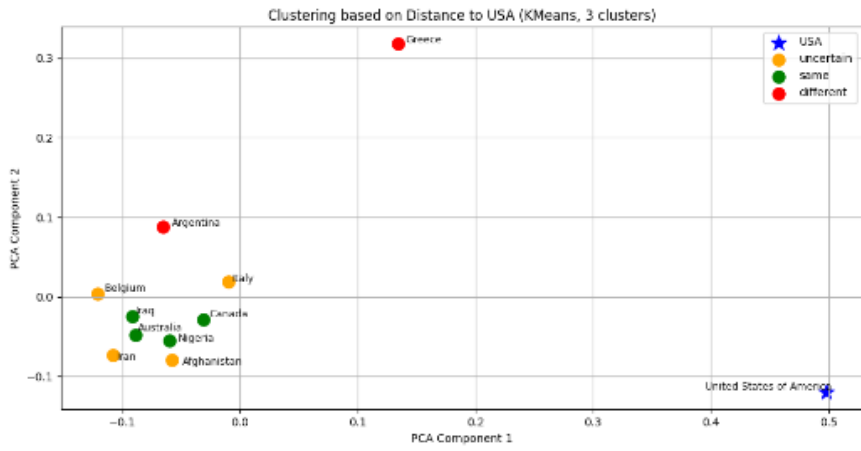


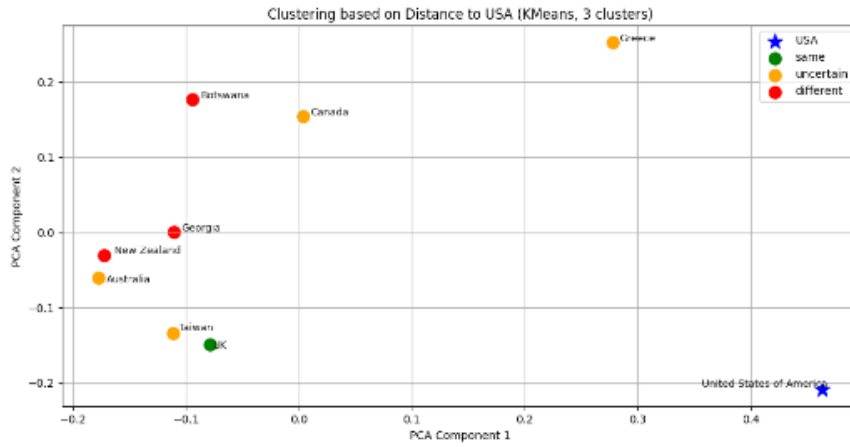




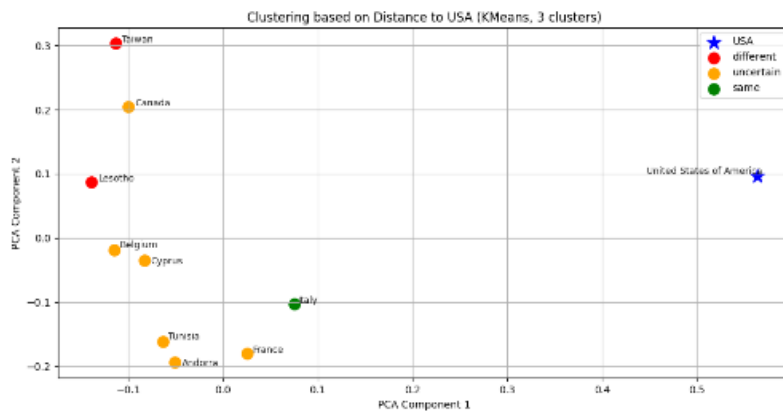
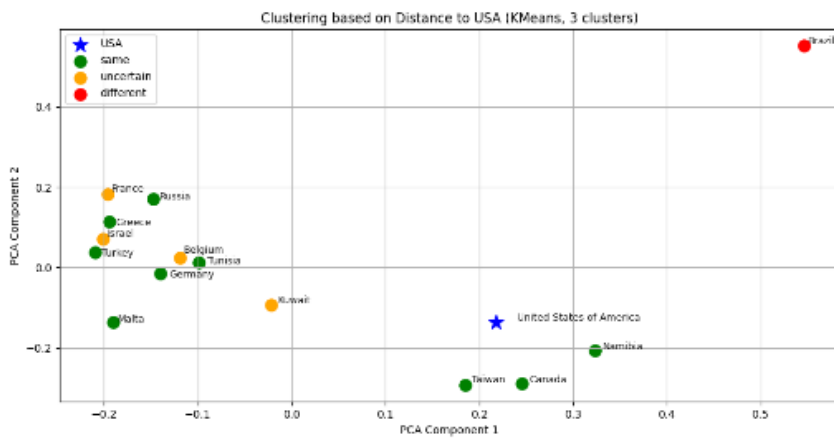


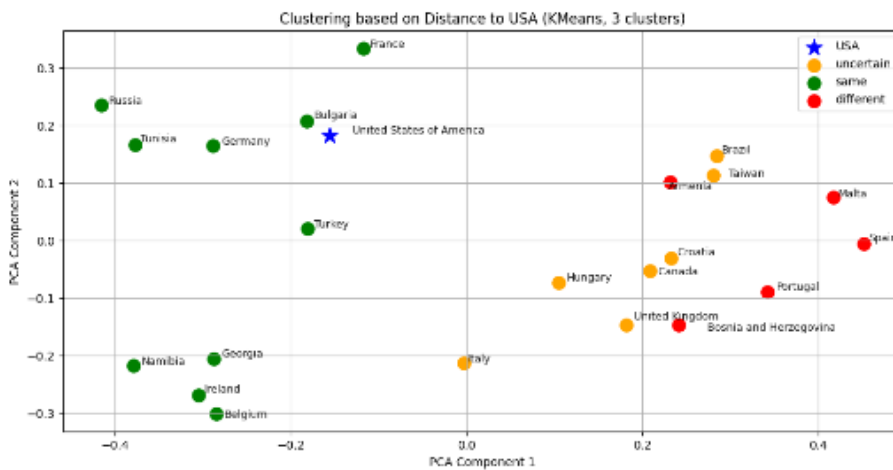
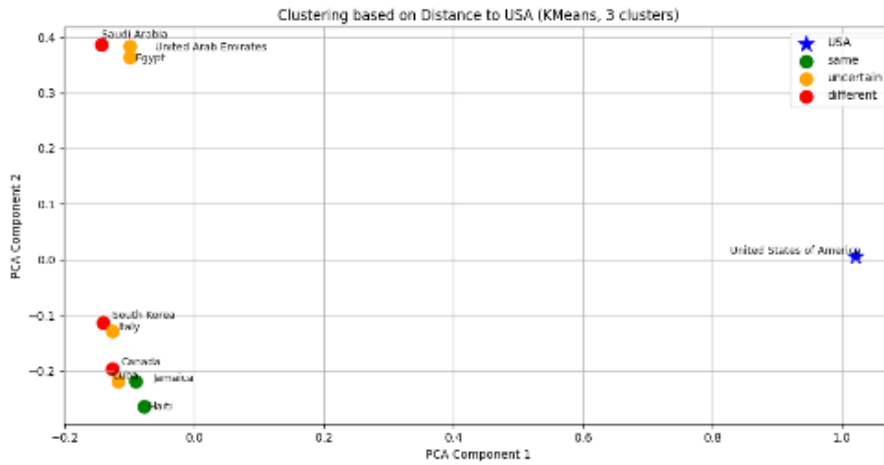
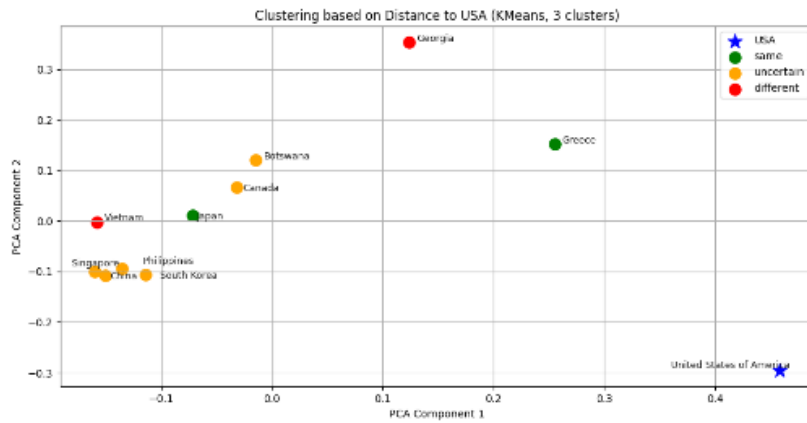


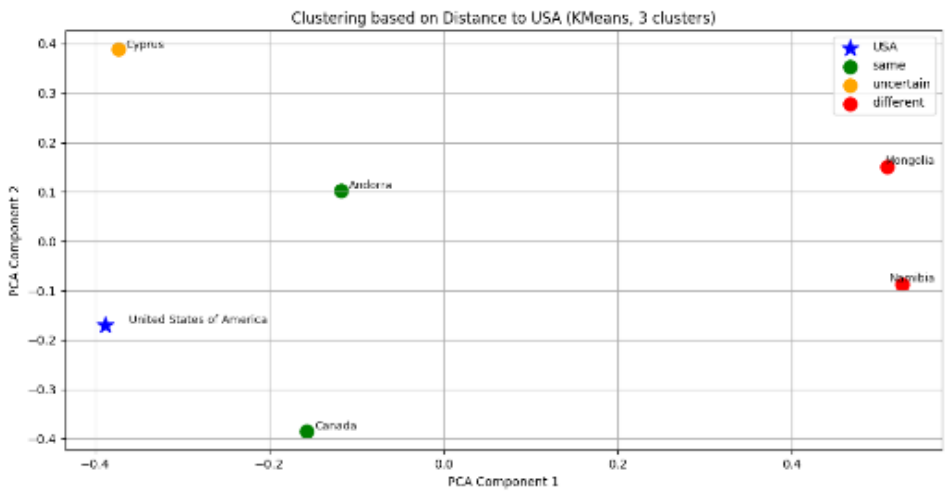
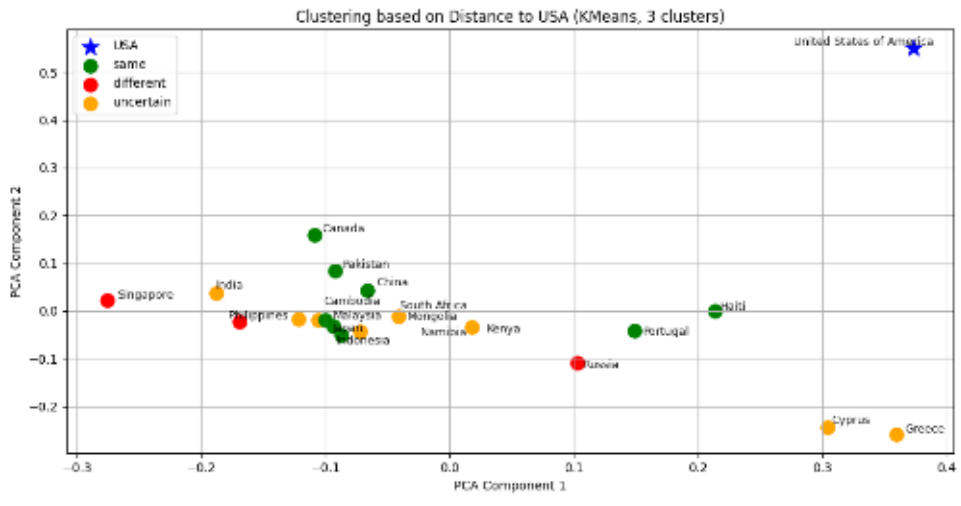
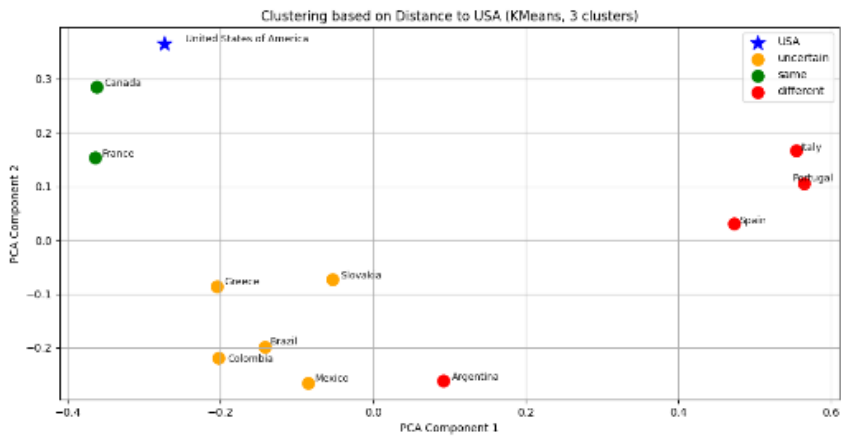


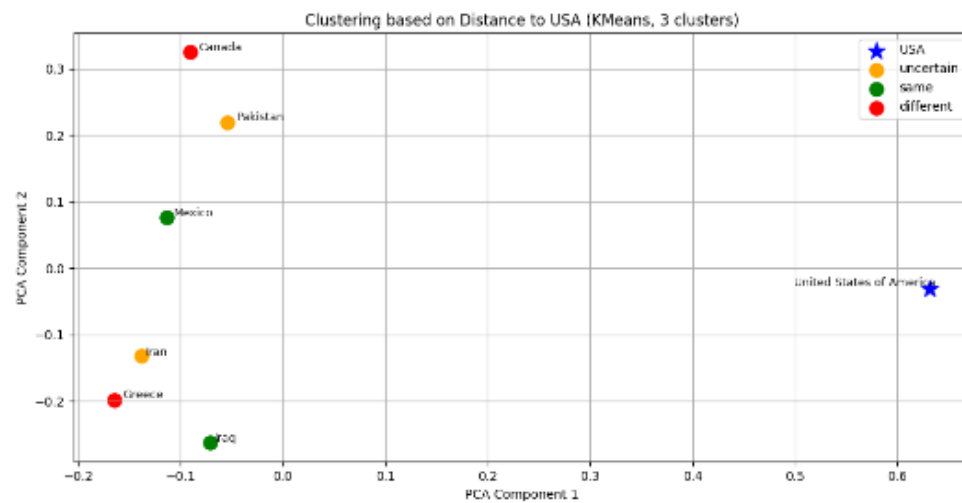
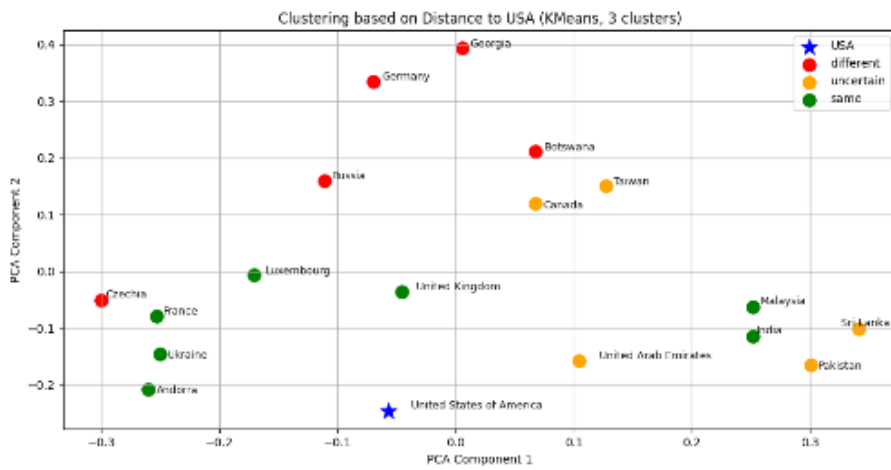
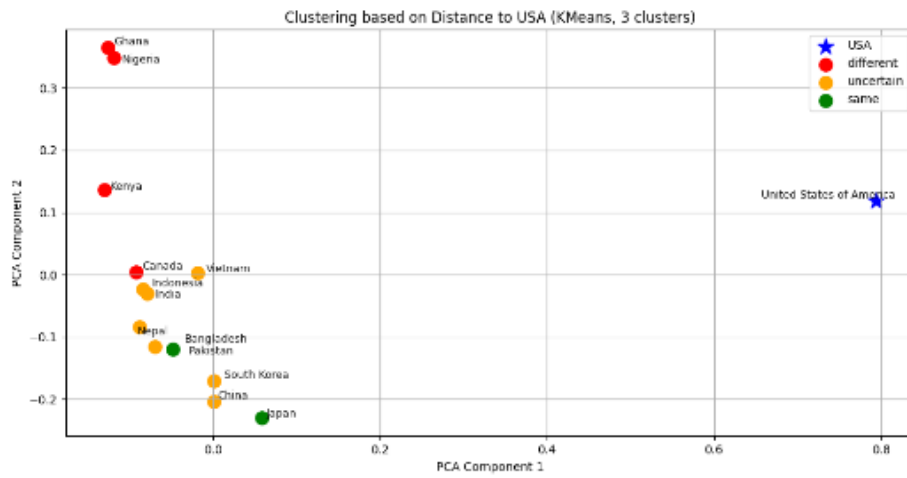


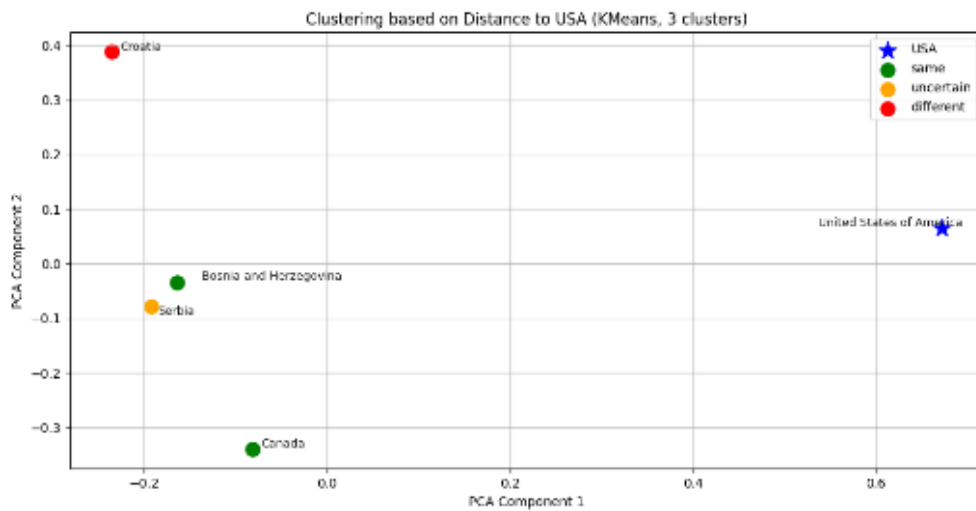
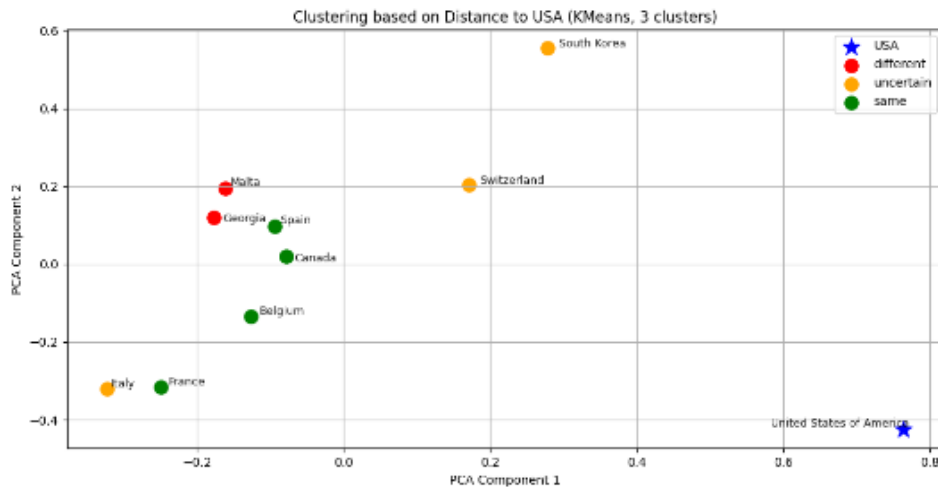
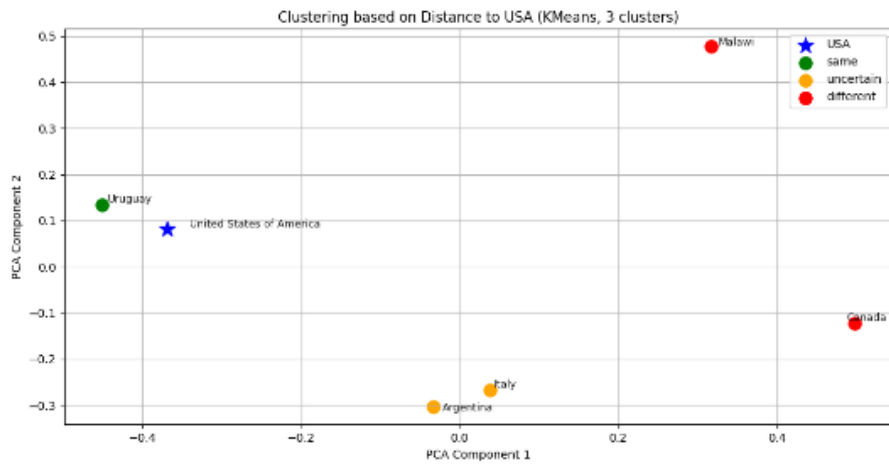
GPT-5. 1 公用語プロンプト 代表的意味を出力

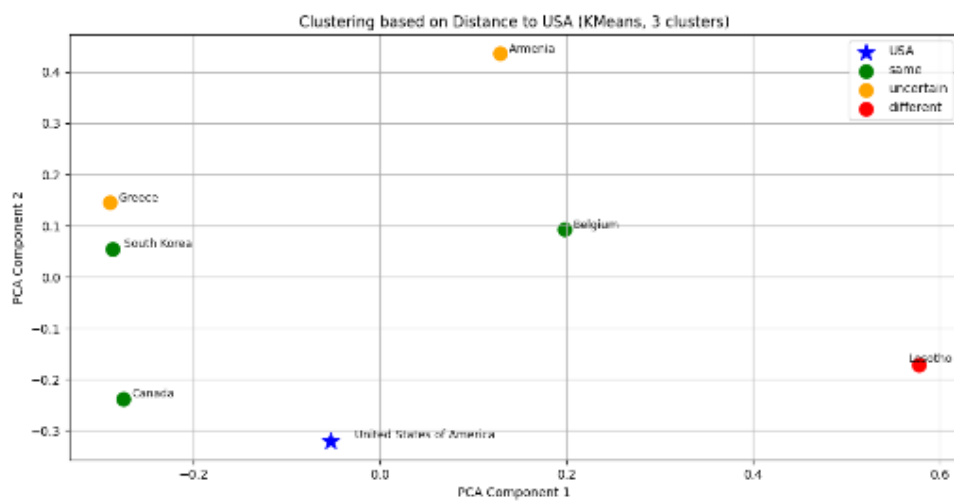
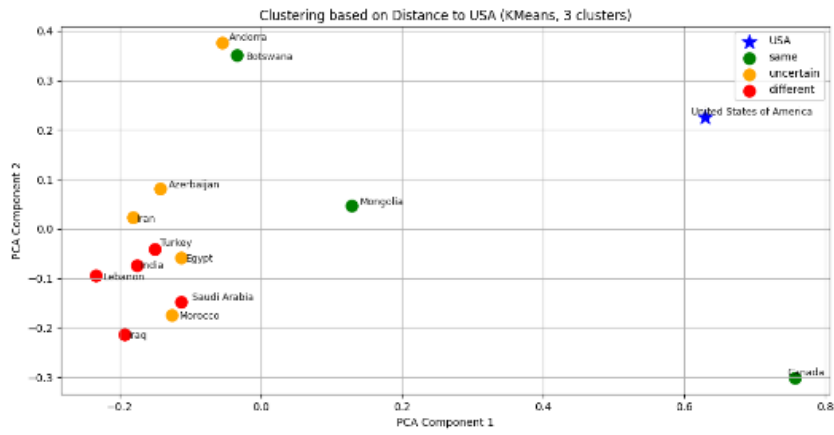
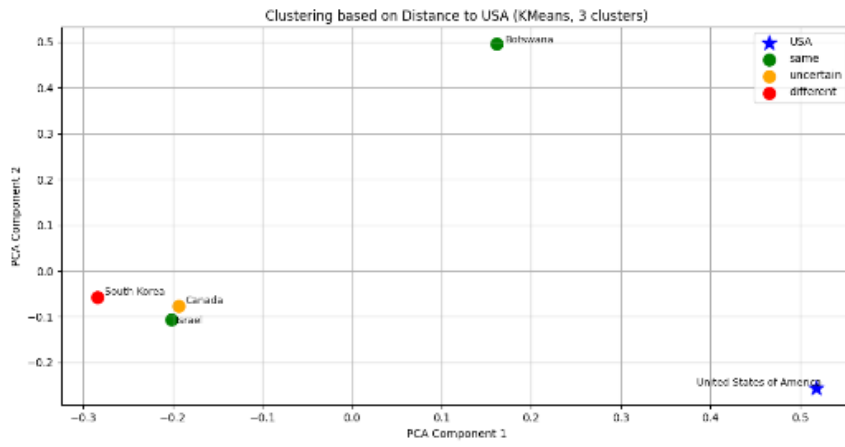


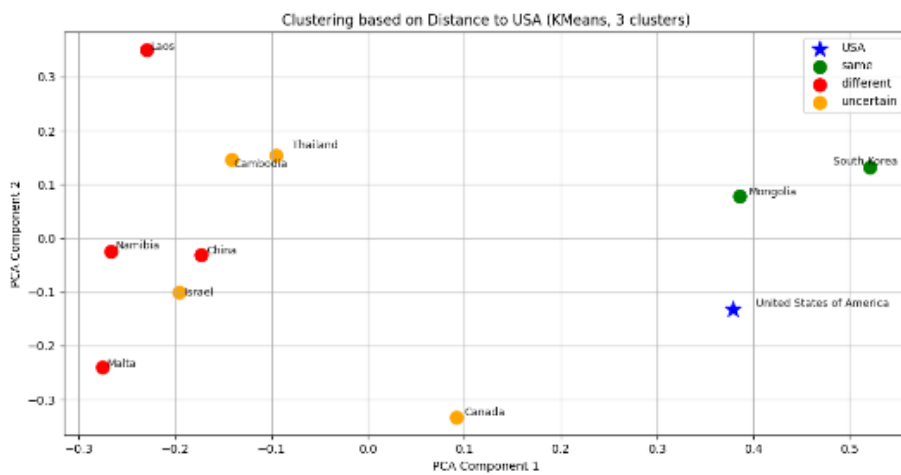
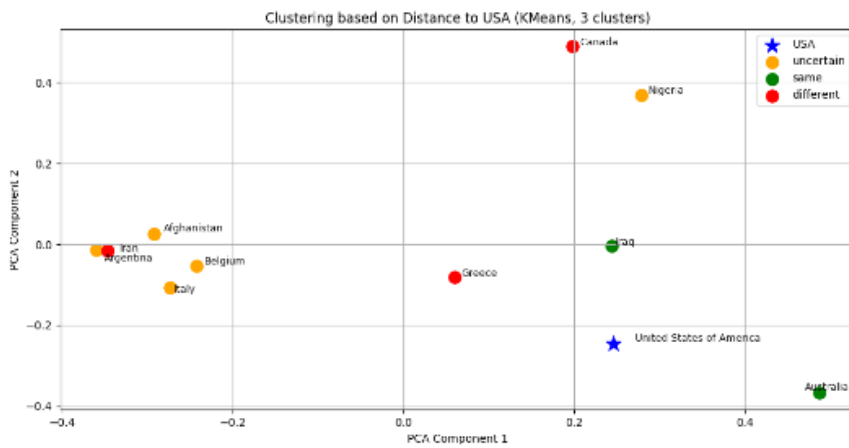
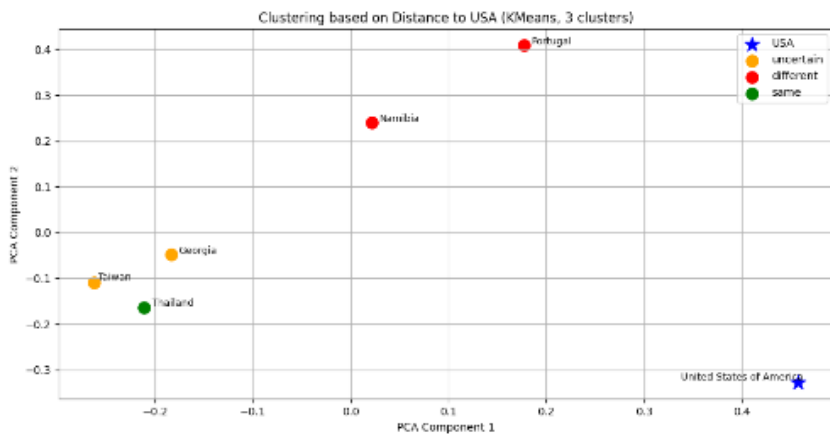


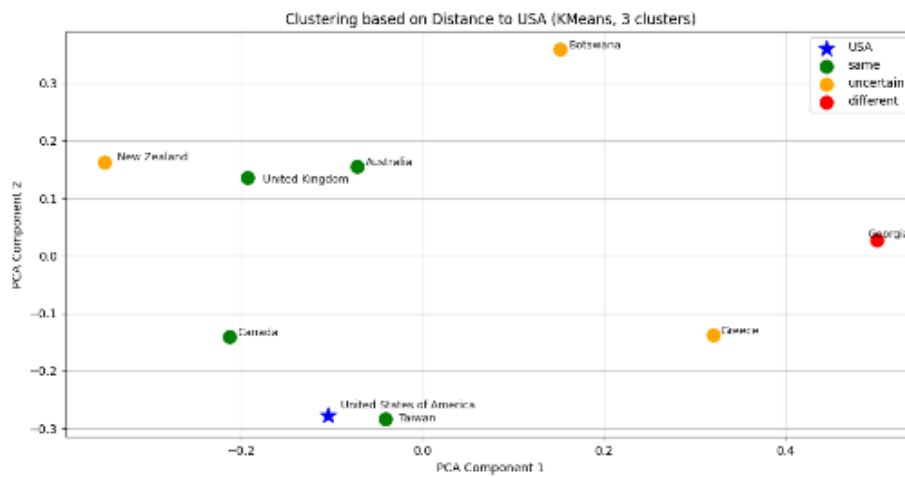
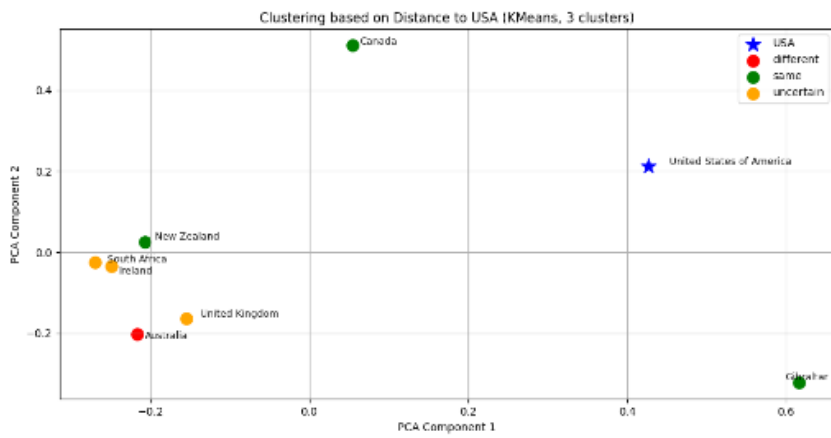












• ジェスチャー画像

• OK



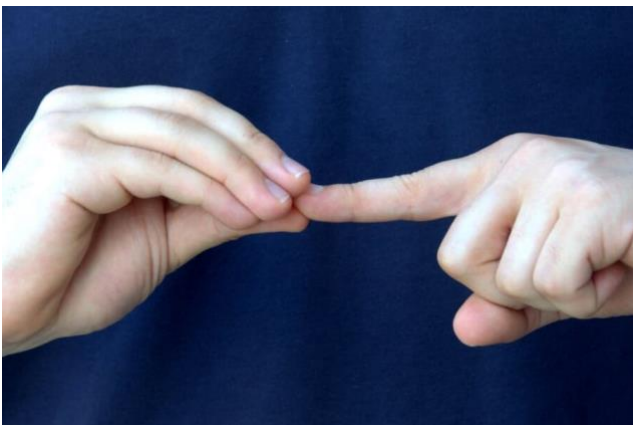
- Chin flick gesture



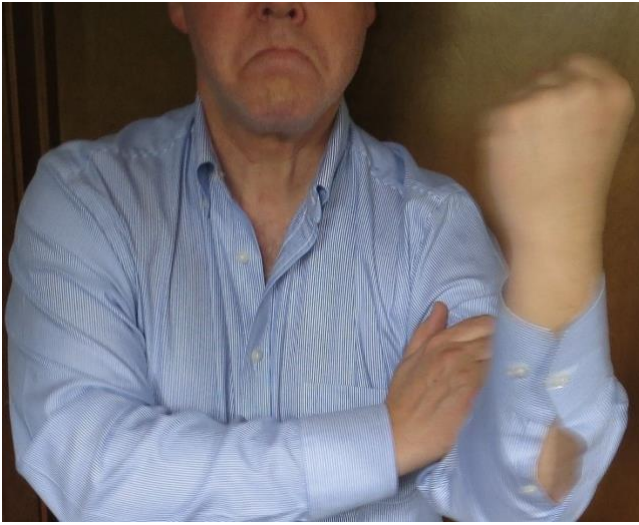
- curled finger gesture



- five fathers gesture



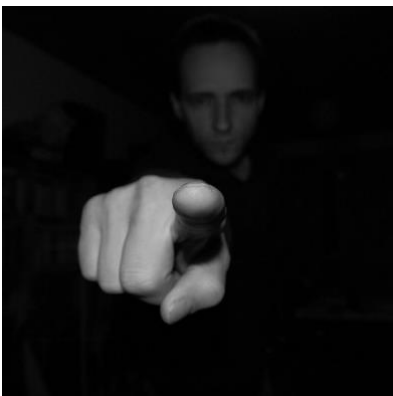
- Forearm_Jerk_gesture



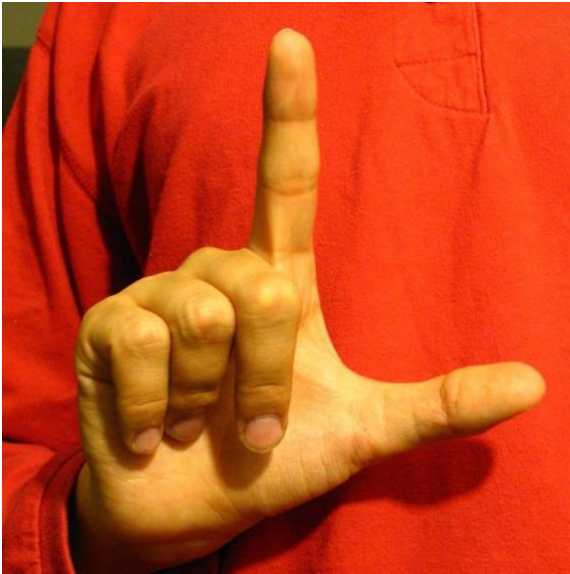
- Horns gesture



- Index_finger_pointing



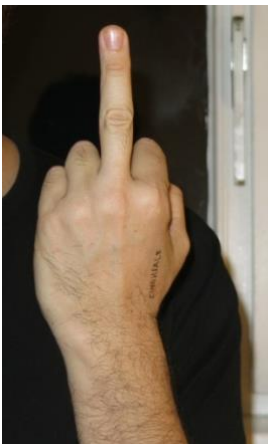
• L gesture



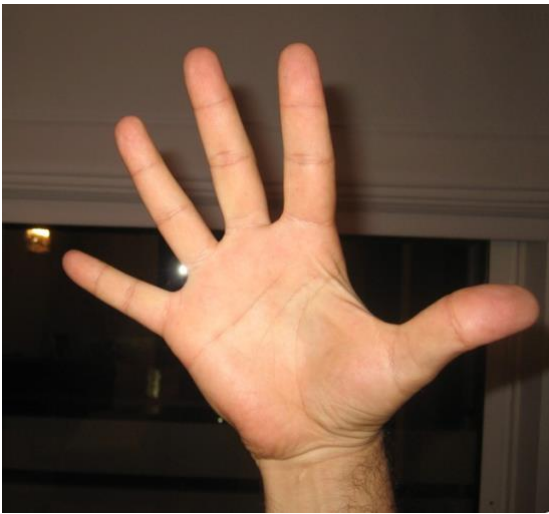
• Left_Hand gesture



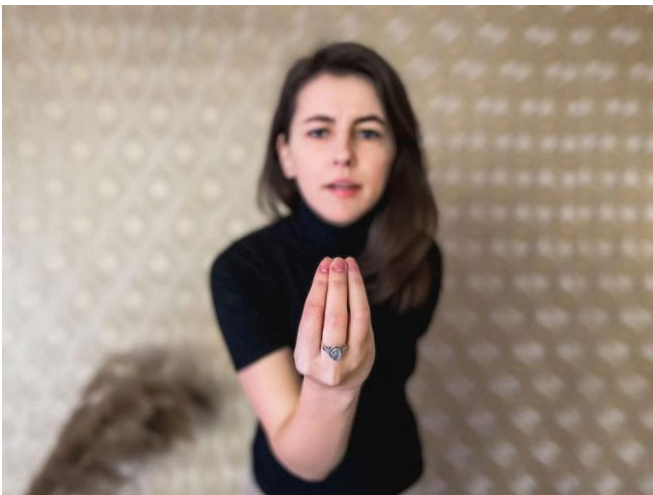
• Middle_Finger gesture



- Open_palm_with_fingers_spread gesture



- Pinched_Fingers gesture



- Quenelle gesture



- Serbian_Salute gesture



- Shocker gesture



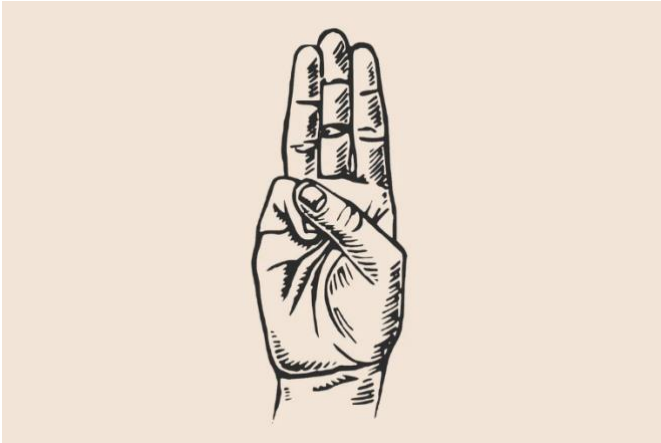
- Show_sole_of_shoe_or_feet gesture



- Snap_Fingers gesture



- Three_Finger_Salute gesture



- Thumbs_up gesture



- Touching_someones_head gesture



- V_sign gesture



- Wanker gesture

