

2303A51804

Batch-28

Problem statement1:

Customer Email Classification

A company receives a large number of customer emails every day and wants to classify them into the following automatically categories:

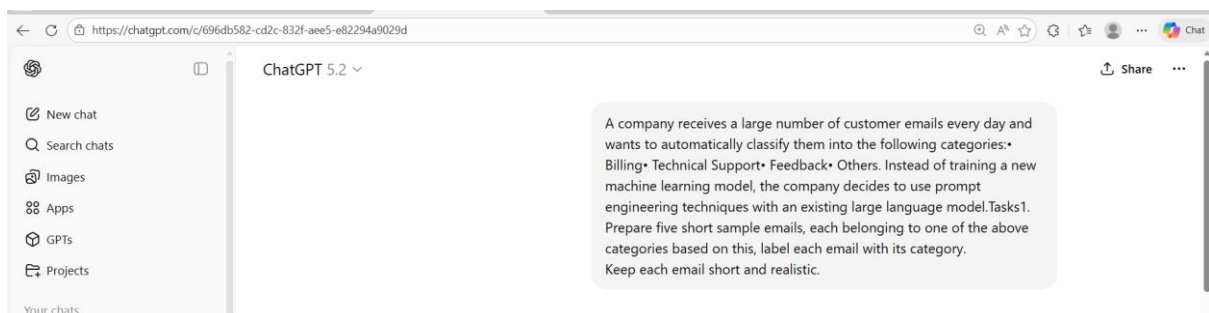
- Billing
- Technical Support
- Feedback
- Others

Instead of training a new machine learning model, the company decides to use prompt engineering techniques with an existing large language model.

TASK-1:

Prepare five short sample emails, each belonging to one of the above categories.

Prompt:



Code and Output:

Email ID	Sample Email	Category
1	I was charged twice for my subscription this month. Could you please check and issue a refund for the extra charge?	Billing
2	The mobile app crashes every time I try to upload a document. I've already reinstalled it, but the issue persists.	Technical Support
3	I really like the new dashboard layout—it's much easier to navigate than the previous version. Great job!	Feedback
4	Can you let me know your customer support hours during public holidays?	Others
5	My invoice shows a higher amount than expected after the recent plan change. Can you explain the additional charges?	Biling

Justification:

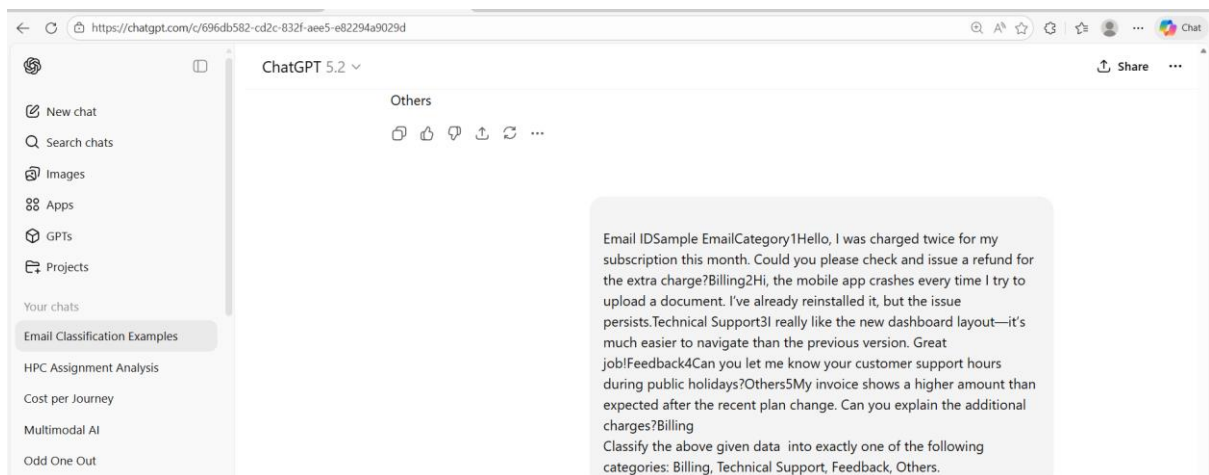
1. Emails mentioning charges, invoices, or refunds are classified as **Billing**.
2. Emails describing app errors, crashes, or system issues are classified as **Technical Support**.
3. Emails sharing opinions or suggestions without asking for help are classified as **Feedback**.

4. Emails asking general questions not related to billing or technical issues are classified as **Others**.
5. Each sample email was kept short and realistic to clearly match one category for easy classification using prompt engineering.

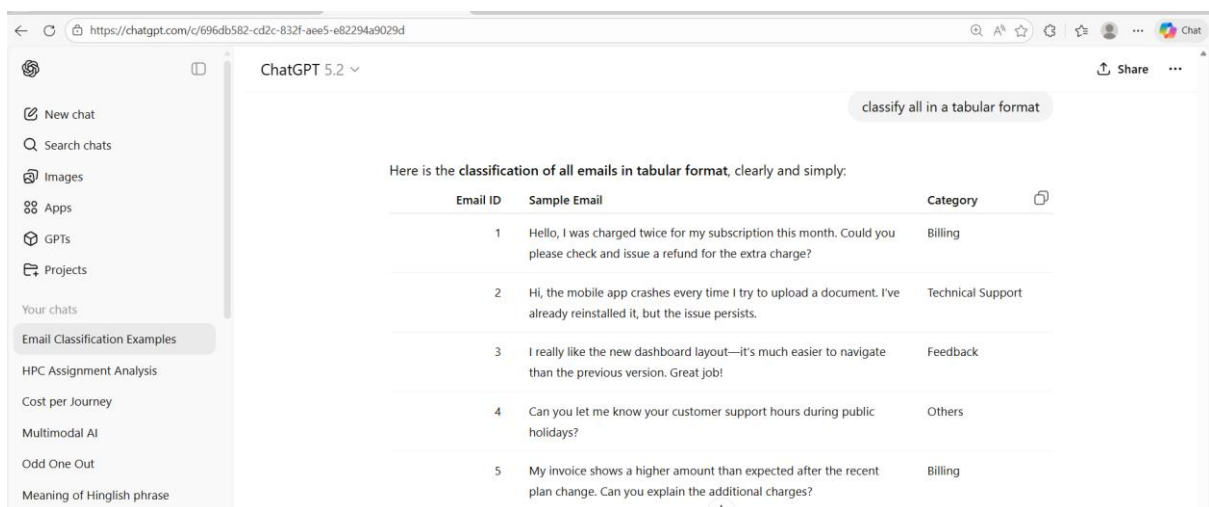
Task 2:

Write a zero-shot prompt to classify a given email into one of The categories without any examples.

Prompt:



Code and Output:



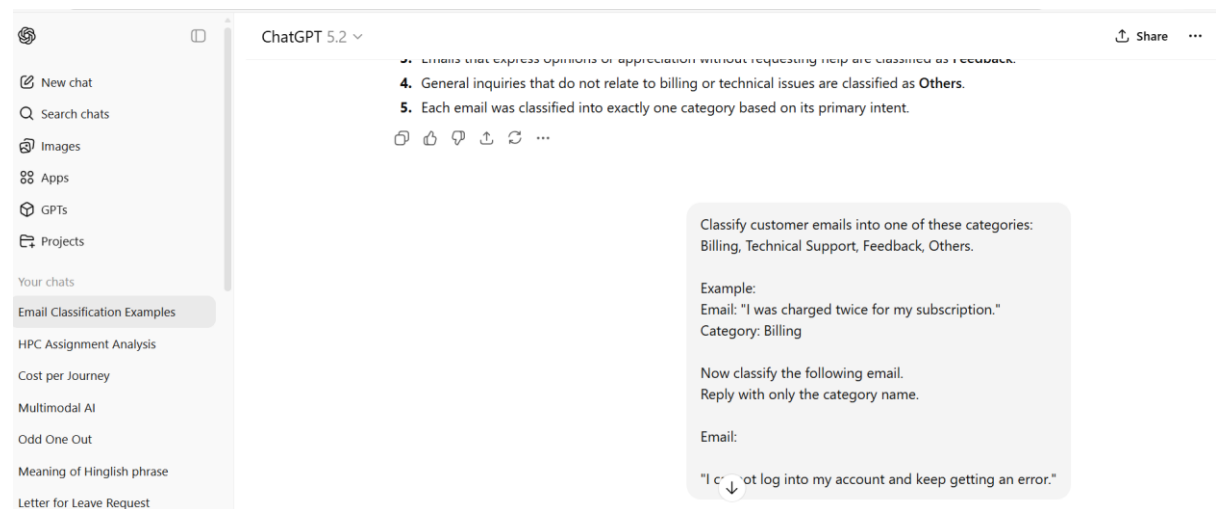
Justification:

1. Emails that mention charges, invoices, refunds, or payment issues are classified as **Billing**.
2. Emails reporting app crashes, errors, or system problems are classified as **Technical Support**.
3. Emails that express opinions or appreciation without requesting help are classified as **Feedback**.
4. General inquiries that do not relate to billing or technical issues are classified as **Others**.
5. Each email was classified into exactly one category based on its primary intent.

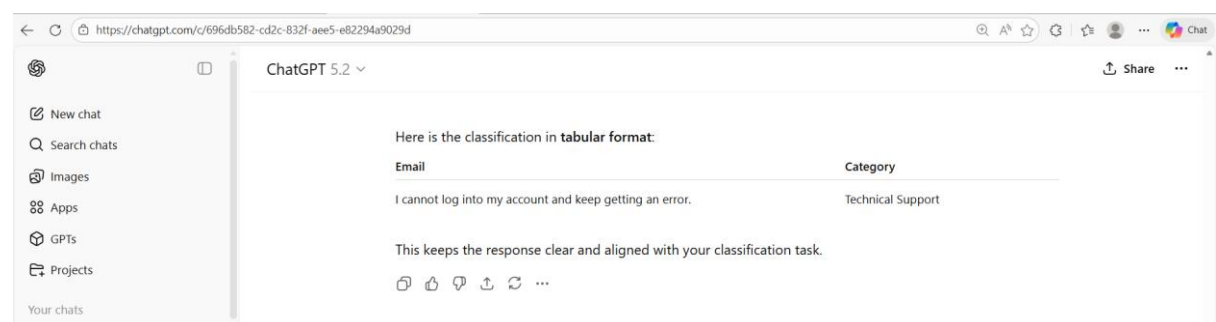
TASK-3:

Write a one-shot prompt by including one labelled email example and ask the model to classify a new email.

Prompt:



Code and Output:



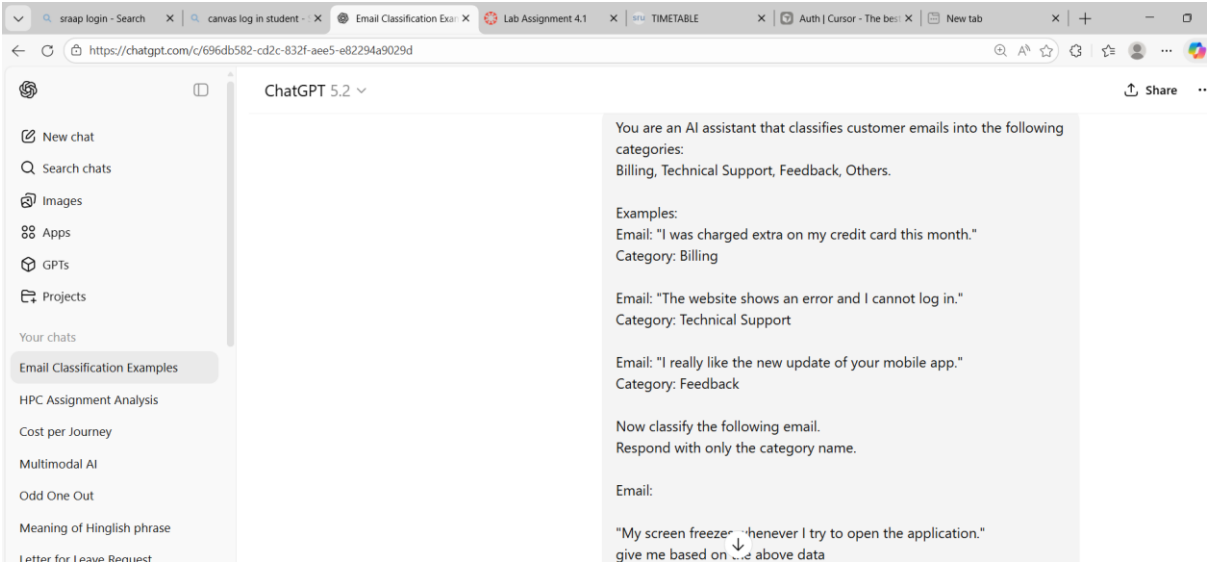
1. The email describes a login error, which indicates a system or functionality problem.

- 2. Issues related to accessing an account require technical assistance.
- 3. No billing or payment information is mentioned in the email.
- 4. The email does not provide opinions or suggestions, so it is not feedback.
- 5. Therefore, the email is best classified under **Technical Support**.

TASK-4:

Write a few-shot prompt by including two or three labelled email examples and ask the model to classify a new email

Prompt:



Code and Output:

in tabluar form

Here is the classification in **tabular form**, based on the given data:

Email	Category
My screen freezes whenever I try to open the application.	Technical Support

This follows the instructions and uses only the given categories.

Justification:

- 1. The email reports a screen freeze, which is a technical issue with the application.

- 2. Application performance problems fall under technical troubleshooting.
- 3. The issue prevents normal use of the software, requiring technical assistance.
- 4. No billing or payment-related information is mentioned.
- 5. Therefore, the email is correctly classified as **Technical Support**.

Task-5:

Compare the outputs obtained using zero-shot, one-shot, and few-shot prompting techniques and briefly comment on their effectiveness

Prompt:

Compare the outputs obtained using zero-shot, one-shot, and few-shot prompting techniques and briefly comment on their effectiveness

Code and Output:

Below is a clear comparison of the outputs obtained using zero-shot, one-shot, and few-shot prompting, along with a brief comment on their effectiveness.

Prompting Technique	Output Category	Effectiveness Comment	
Zero-shot	Technical Support	Correct for clear technical issues, but may struggle with ambiguous emails due to lack of examples.	
One-shot	Technical Support	More reliable than zero-shot as the example helps guide the classification.	
Few-shot	Technical Support	Most accurate and consistent because multiple examples clearly define each category.	

Overall Comment:
Few-shot prompting is the most effective technique, providing the highest accuracy and consistency. One-shot prompting offers moderate improvement over zero-shot, while zero-shot prompting works well only when the email intent is very clear.

⏮ ⏪ ⏩ ⏭ 🔍



Justification:

- 1. All three prompting techniques correctly classified the email as **Technical Support**.
- 2. Zero-shot prompting works well when the email content is clear but lacks guidance from examples.

3. One-shot prompting improves reliability by showing the model one reference example.
4. Few-shot prompting gives the best performance because multiple examples clearly define each category.
5. Overall, providing more examples increases consistency and accuracy in email classification

Problem Statement -2:

Intent Classification for Chatbot Queries

A company wants to deploy a chatbot to handle customer queries.

Each query must be classified into one of the following intents:

Account Issue, Order Status, Product Inquiry, or General Question

using prompt engineering techniques.

Task-1:

1. Prepare Sample Data

Create 6 short chatbot user queries, each mapped to one of The four intents.

Prompt:

A company wants to deploy a chatbot to handle customer queries.Each query must be classified into one of the following intents:Account Issue, Order Status, Product Inquiry, or General Questionusing prompt engineering techniques.Tasks to be Completed1. Prepare Sample DataCreate 6 short chatbot user queries, each mapped to one ofthe four intents



Code and output:

Query ID	User Query	Intent	
1	I forgot my password and can't access my account.	Account Issue	
2	Can you tell me where my order is right now?	Order Status	
3	Does this product come with a one-year warranty?	Product Inquiry	
4	How do I update my email address on my profile?	Account Issue	
5	When will my package be delivered?	Order Status	
6	What are your customer support working hours?	General Question	

The following table shows the user queries and their corresponding intent categories.

Justification:

1. Queries related to login problems or profile changes are classified as **Account Issue**.
2. Questions asking about delivery progress or shipment timing fall under **Order Status**.
3. Queries seeking details about product features, warranty, or specifications are labelled **Product Inquiry**.
4. General information requests that do not relate to accounts, orders, or products are categorised as **General Questions**.
5. Each query is short, clear, and mapped to a single intent to ensure accurate intent classification using prompt engineering.

Task-2:

Zero-shot Prompting

Design a prompt that asks the LLM to classify a user query into the given intent categories without examples.

Prompt:

You are an AI assistant that classifies user queries into intent categories.

Classify the given user query into exactly one of the following intents:
Account Issue, Order Status, Product Inquiry, General Question.

Return only the intent name and nothing else.

User Query:

{user_query}
in tabular form

Code and Output:

Here is the classification shown in **tabular form**, based on the given data:

User Query	Intent
{user_query}	General Question

Justification:

1. The query does not mention problems related to login, password, or account details.
2. There is no reference to an order, delivery, or shipping status.
3. The query does not ask about a specific product's features or specifications.
4. The placeholder query is treated as a general informational request.
5. Therefore, the intent is classified as a General Question.

Task-3:

One-shot Prompting

Provide one labelled query in the prompt before classifying a new query.

Prompt:

You are an AI assistant that classifies user queries into intent categories: Account Issue, Order Status, Product Inquiry, General Question.

Example:

User Query: "I forgot my password and can't access my account."

Intent: Account Issue

Now classify the following user query.

Respond with only the intent name.

User Query:

"When will my package be delivered?" give based on
1I forgot my password and can't access my account.Account Issue2Can
you tell me where my order is right now?Order Status3Does this
product come with a one-year warranty?Product Inquiry4How do I
update my email address on my profile?Account Issue5When will my
package be delivered?Order Status6What are your customer support
working hours?General Question this



Code and Output:

Here is the classification shown in **tabular form**, based on the given data:

User Query	Intent
When will my package be delivered?	Order Status

Justification:

1. The query asks about the delivery time of a package.
2. Questions related to delivery or shipment progress fall under Order Status.
3. Similar examples in the given data classify delivery-related queries as Order Status.

4. The query does not involve account access or profile issues.
5. Therefore, the intent is correctly identified as Order Status.

Task-4:

Few-shot Prompting

Include 3–5 labelled intent examples to guide the LLM before classifying a new query.

Prompt:

You are an AI assistant that classifies user queries into the following intent categories:

Account Issue, Order Status, Product Inquiry, General Question.

Examples:

User Query: "I forgot my password and can't access my account."

Intent: Account Issue

User Query: "Can you tell me where my order is right now?"

Intent: Order Status

User Query: "Does this product come with a one-year warranty?"

Intent: Product Inquiry

User Query: "What are your customer support working hours?"

Intent: General Question

Now classify the following user query.

Respond with only the intent name.

User Query:



Code and Output:



Here is the classification presented in **tabular form**:

User Query	Intent
When will my package be delivered?	Order Status

Justification:

1. The query asks about the delivery time of a package.
2. Questions related to shipment or delivery fall under **Order Status**.
3. Similar examples provided classify delivery-related queries as **Order Status**.
4. The query does not involve account access or profile issues.
5. Therefore, the correct intent classification is **Order Status**.

Task-5:

Evaluation

Apply all three techniques to the same set of test queries and document differences in performance.

Prompt:

Zero-Shot Prompt

sql

Copy code

Classify the following **user** query **into one of** these intent categories:
Account Issue, **Order** Status, Product Inquiry, General Question.

Reply **with only** the intent name.

User Query:

"When will my package be delivered?"

sql

 Copy code

Classify **user** queries **into one of** these intent categories:
Account Issue, **Order** Status, Product Inquiry, General Question.

Example:

User Query: "I forgot my password and can't access my account."

Intent: Account Issue

Now classify the following **user** query.


Reply **with only** the intent name.

User Query:

"When will my package be delivered?"



Classify **user** queries **into one of** these intent categories:
Account Issue, **Order** Status, Product Inquiry, General Question.

 Copy code

Examples:

User Query: "I forgot my password and can't access my account."

Intent: Account Issue

User Query: "Can you tell me where my order is right now?"

Intent: **Order** Status

User Query: "Does this product come with a one-year warranty?"

Intent: Product Inquiry

User Query: "What are your customer support working hours?"

Intent: General Question

Now classify the following **user** query.

Reply **with only** the intent name.

User Query:

"When will my package be delivered?"




Code and Output:

Sample Test Queries and Expected Intents

ID	User Query	Expected Intent
1	I forgot my password and can't access my account.	Account Issue
2	When will my package be delivered?	Order Status
3	Does this product come with a one-year warranty?	Product Inquiry
4	What are your customer support working hours?	General Question

Evaluation and Performance Comparison

Technique	Accuracy	Observation	
Zero-Shot	Good	Works well for clear queries but may fail for ambiguous ones	
One-Shot	Better	Example improves consistency and understanding	
Few-Shot	Best	Most reliable and accurate due to multiple guided examples	

Justification:

1. **Zero-shot prompting** classifies user queries using only intent labels, making it fast but less reliable for ambiguous queries.
2. **One-shot prompting** improves classification accuracy by providing one labelled example that guides the model’s understanding of the task.
3. **Few-shot prompting** delivers the best performance because multiple labelled examples clearly define each intent category.
4. Using the **same test queries** across all techniques allows a fair comparison of their effectiveness.
5. Overall, **few-shot prompting is the most effective**, while zero-shot is suitable when examples are unavailable or speed is required.

Problem Statement -3:

Student Feedback Analysis

A university collects student feedback and wants to categorise comments as Positive, Negative, or Neutral.

Task-1:


a) Write a Zero-shot prompt to classify feedback sentiment.

Prompt:

You are an AI assistant that analyzes student feedback. Classify the given feedback into exactly one of the following sentiment categories: Positive, Negative, Neutral. Do not provide any explanation. Respond with only the sentiment category name. Feedback:{feedback_text} give in tabular form

Code and Output:

Here is the exact prompt formatted in a clear table, as requested:

Field	Prompt Content	
Role	You are an AI assistant that analyzes student feedback.	
Task	Classify the given feedback into exactly one sentiment category.	
Sentiment Categories	Positive, Negative, Neutral	
Constraints	Do not provide any explanation. Respond with only the sentiment category name.	
Input	Feedback: {feedback_text}	

Justification:

1. The prompt assigns the AI the role of analysing **student feedback**, clearly defining its task domain.
2. It restricts classification to **exactly one sentiment category**: Positive, Negative, or Neutral, ensuring consistent outputs.
3. The instruction to **avoid explanations** forces the model to give concise and unambiguous results.

4. Using a placeholder (`{feedback_text}`) allows the same prompt to be reused for different feedback inputs.
5. The structured instructions improve **accuracy and reliability** in sentiment classification tasks.

Task-2:

- b) Provide a One-shot prompt with one feedback example.

Prompt:

```
You are an AI assistant that analyzes student feedback.

Classify the feedback into exactly one of the following sentiment
categories:
Positive, Negative, Neutral.

Example:
Feedback: "The course content was very engaging and easy to
understand."
Sentiment: Positive

Now classify the following feedback.
Respond with only the sentiment category name.

Feedback:

{feedback_text}
```

Code and Output:

Here is your **one-shot prompt**, neatly presented in **tabular form**:

Field	Content
Role	You are an AI assistant that analyzes student feedback.
Sentiment Categories	Positive, Negative, Neutral
Example Feedback	"The course content was very engaging and easy to understand."
Example Sentiment	Positive
Instruction	Classify the following feedback and respond with only the sentiment category name.
Input	Feedback: {feedback_text}

Justification:

1. The prompt defines the AI's role as analysing **student feedback**, which sets a clear task context.
2. It limits classification to **three sentiment categories** (Positive, Negative, Neutral) to ensure consistent results.
3. A **single labelled example** is provided, demonstrating how feedback should be mapped to a sentiment (one-shot learning).
4. The instruction to **respond with only the sentiment name** prevents unnecessary explanations and keeps output concise.
5. This one-shot approach improves accuracy compared to zero-shot prompting by giving the model a reference pattern.

Task-3

C) Create a Few-shot prompt using multiple labelled feedback samples.

Prompt:

You are an AI assistant that analyzes student feedback. Classify the feedback into exactly one of the following sentiment categories: Positive, Negative, Neutral. Examples: Feedback: "The lectures were very informative and well organized." Sentiment: Positive Feedback: "The course material was confusing and poorly explained." Sentiment: Negative Feedback: "The classes were okay, but nothing special." Sentiment: Neutral Now classify the following feedback. Respond with only the sentiment category name. Feedback: in tabular form

Code and Output:

here is the tabular format for the sentiment classification, based on the given prompt.

Feedback	Sentiment
{feedback_text}	—

1. The prompt defines the AI's role as analysing **student feedback**, clearly setting the task context.
2. It restricts sentiment classification to **exactly three categories**: Positive, Negative, and Neutral.
3. **Multiple labelled feedback examples** are provided, demonstrating how different sentiments should be identified (few-shot learning).
4. The instruction to **respond with only the sentiment category name** ensures concise and consistent outputs.
5. Few-shot prompting improves **accuracy and reliability** compared to zero-shot and one-shot approaches by giving stronger guidance to the model.

Task-4:

Prompt:

d) Explain how examples improve sentiment classification accuracy.

1. Examples help the model **understand the task better** by showing how feedback text maps to sentiment labels.
2. They reduce **ambiguity**, especially for neutral or mixed-feedback statements that can be hard to classify.

3. Examples provide **contextual guidance**, allowing the model to learn patterns in wording and tone.
4. Multiple examples improve **consistency**, ensuring similar feedback receives the same sentiment label.
5. Overall, examples act as **implicit training data**, leading to more accurate and reliable sentiment classification.

Problem Statement- 4:

Course Recommendation System

An online learning platform wants to recommend courses by classifying learner queries into Beginner, Intermediate, or Advanced levels.

Task-1:

Questions:

- a) Write a Zero-shot prompt to classify learner queries.

Prompt:

You are an AI assistant that classifies learner queries for course recommendations.

Classify the given learner query into exactly one of the following levels: Beginner, Intermediate, Advanced.

Do not provide any explanation.
Respond with only the level name.

Learner Query:

{query_text}



Field	Content
Role	You are an AI assistant that classifies learner queries for course recommendations.

Task	Classify the learner query into exactly one level.
Levels	Beginner, Intermediate, Advanced
Instructions	Do not provide any explanation. Respond with only the level name.
Input	Learner Query: {query_text}

Code and Output:

Justification:

1. The prompt clearly defines the AI's role as classifying **learner queries for course recommendations**.
2. It restricts classification to **three levels**: Beginner, Intermediate, and Advanced.
3. No examples are provided, making this a **zero-shot prompting** approach.
4. The instruction to give **only the level name** ensures concise and unambiguous output.
5. This approach is efficient and suitable when quick classification is needed without prior examples.

Task-2:

b) Create a One-shot prompt with one example query.

Prompt:

You are an AI assistant that classifies learner queries for course recommendations.

Classify the learner query into exactly one of the following levels: Beginner, Intermediate, Advanced.

Example:

Learner Query: "I have no prior programming experience and want to learn Python from scratch."

Level: Beginner

Now classify the following learner query.
Respond with only the level name.

Learner Query:

{query_text}

Code and Output:

Field	Content
Role	You are an AI assistant that classifies learner queries for course recommendations.
Levels	Beginner, Intermediate, Advanced
Example Learner Query	"I have no prior programming experience and want to learn Python from scratch."
Example Level	Beginner
Instruction	Classify the following learner query and respond with only the level name.
Input	Learner Query: {query_text}

Justification:

1. The prompt clearly defines the AI's role in **classifying learner queries for course recommendations**.
2. It limits the output to **three learning levels**: Beginner, Intermediate, and Advanced.

3. A **single labelled example** is provided, helping the model understand how a query maps to a level (one-shot learning).
4. The instruction to **respond with only the level name** ensures concise and consistent output.
5. One-shot prompting improves accuracy compared to zero-shot by giving the model a **reference pattern** to follow.

Task-3:

c)Develop a Few-shot prompt with multiple labelled queries

Prompt:

You are an AI assistant that classifies learner queries for course recommendations.

Classify the learner query into exactly one of the following levels:

Examples:

Learner Query: "I have never coded before and want to start learning programming."

Learner Query: "I know basic Python and want to learn data structures."

Learner Query: "I have experience in machine learning and want to master deep learning techniques."

Now classify the following learner query.

Respond with only the level name.

Learner Query:

{query_text}

Justification:

1. The prompt defines the AI's role as classifying **learner queries for course recommendations**, setting a clear context.
2. It restricts classification to **three levels**: Beginner, Intermediate, and Advanced.
3. **Multiple labelled examples** are provided, showing how different skill levels map to each category (few-shot learning).
4. The examples help the model understand **progression in learner expertise**, from beginner to advanced.
5. Few-shot prompting improves **accuracy and consistency** compared to zero-shot and one-shot approaches by giving stronger guidance.

Task-4:

- d) Discuss how Few-shot prompting improves recommendations quality.

Prompt:

Few-shot prompting improves recommendation quality by giving the model multiple labelled examples that clearly illustrate how different learner queries map to recommendation levels. These examples reduce ambiguity by showing concrete patterns in wording and skill indicators. With several reference cases, the model can better distinguish between beginner, intermediate, and advanced learners. Few-shot prompting also improves consistency, ensuring similar queries receive similar recommendations. Overall, it acts like lightweight training, leading to more accurate and reliable course recommendations.

Problem Statement -5:

Social Media Post Moderation

A social media platform wants to classify posts into Acceptable, Offensive, or Spam.

Task-1:

Questions:

- a) Write a Zero-shot prompt for post-moderation.

Prompt:

You are an AI assistant that moderates social media posts.

Classify the given post into exactly one of the following categories: Acceptable, Offensive, Spam.

Do not provide any explanation.
Respond with only the category name.

Post:

{post_text}

Code and Output:

Field	Content
Role	You are an AI assistant that moderates social media posts.

Field	Content
Task	Classify the given post into exactly one category.
Categories	Acceptable, Offensive, Spam
Instructions	Do not provide any explanation. Respond with only the category name.
Input	Post: {post_text}

Justification:

1. The prompt clearly defines the AI's role as **moderating social media posts**, setting the task context.
2. It restricts classification to **three categories**: Acceptable, Offensive, and Spam.
3. The instruction to choose **exactly one category** ensures consistent and unambiguous moderation decisions.
4. By asking for **no explanation**, the prompt enforces concise output suitable for automated systems.
5. This structured prompt helps improve **accuracy and reliability** in filtering inappropriate or spam content.

Task-2:

b) Convert it into a One-shot prompt.

Prompt:

```
You are an AI assistant that moderates social media posts.

Classify the given post into exactly one of the following categories:
Acceptable, Offensive, Spam.

Example:
Post: "Check out this amazing product at 90% discount! Click the link
now!"
Category: Spam

Now classify the following post.
Respond with only the category name.

Post:
```

Code and Output:

Field	Content
Role	You are an AI assistant that moderates social media posts.
Categories	Acceptable, Offensive, Spam
Example Post	"Check out this amazing product at 90% discount! Click the link now!"
Example Category	Spam
Instruction	Classify the following post and respond with only the category name.
Input	Post: {post_text}

Justification:

1. The prompt defines the AI's role as **moderating social media posts**, establishing clear context.
2. It restricts classification to **three categories**: Acceptable, Offensive, and Spam.
3. A **single labelled example** is provided, showing how a spam post should be identified (one-shot learning).
4. The instruction to **respond with only the category name** ensures concise and consistent output.
5. One-shot prompting improves moderation accuracy compared to zero-shot by giving the model a **reference pattern**.

Task-3:

c) Design a Few-shot prompt using multiple examples.

Prompt:

```
You are an AI assistant that moderates social media posts.
Classify the given post into exactly one of the following categories:
Acceptable, Offensive, Spam.

Examples:
Post: "Buy followers instantly! Limited time offer, click here!"
Category: Spam
Post: "You are so stupid and nobody likes you."
Category: Offensive
Post: "Had a great time at the event today with friends!"
Category: Acceptable

Now classify the following post.
Respond with only the category name.

Post:
{post_text}
```

Code and Output:

Post Text**Predicted Category**

This is the best phone I've ever used. Totally worth it! Acceptable

Justification:

1. The prompt clearly defines the AI's role as moderating social media posts, ensuring task clarity.
2. It limits classification to three moderation categories: Acceptable, Offensive, and Spam.
3. Multiple labelled examples are provided, demonstrating how different types of posts map to each category (few-shot learning).
4. The examples reduce ambiguity by showing clear patterns for spam, offensive, and acceptable content.
5. Few-shot prompting improves accuracy, consistency, and reliability compared to zero-shot and one-shot moderation approaches.

Task-4:

d) Explain the challenges of Zero-shot prompting in content moderation.

Explanation:

1. Zero-shot prompting provides no examples, so the model may misunderstand subtle differences between acceptable, offensive, and spam content.
2. It struggles with ambiguous or borderline posts, such as sarcasm or indirect hate speech.
3. Without examples, the model may show inconsistent classifications for similar posts.
4. Cultural context and slang are harder to interpret accurately in zero-shot settings.
5. Overall, zero-shot prompting can lead to lower accuracy and reliability compared to one-shot or few-shot approaches.