

# Kelp Global Junior Data Scientist - Practice Papers

Prepared by Perplexity AI

## Contents

<b>1 KELP GLOBAL - JUNIOR DATA SCIENTIST PRACTICE PAPERS</b>	<b>3</b>
1.1 Understanding the Format (Based on Senior's Feedback) . . . . .	3
<b>2 PRACTICE PAPER SET 1</b>	<b>3</b>
2.1 Section A: JSON Data Debugging (20 marks) . . . . .	3
2.2 Section B: SQL Debugging and Normalization (25 marks) . . . . .	3
2.3 Section C: Python Script Debugging (20 marks) . . . . .	4
2.4 Section D: RAG Workflow Design (15 marks) . . . . .	5
2.5 Section E: Model Selection Scenarios (20 marks) . . . . .	5
<b>3 PRACTICE PAPER SET 1 - ANSWERS</b>	<b>5</b>
3.1 Section A: JSON Data Debugging - ANSWERS . . . . .	5
3.2 Section B: SQL Debugging and Normalization - ANSWERS . . . . .	7
3.3 Section C: Python Script Debugging - ANSWERS . . . . .	8
3.4 Section D: RAG Workflow Design - ANSWERS . . . . .	9
3.5 Section E: Model Selection Scenarios - ANSWERS . . . . .	10
<b>4 PRACTICE PAPER SET 2</b>	<b>11</b>
4.1 Section A: JSON Data Debugging (20 marks) . . . . .	11
4.2 Section B: SQL Debugging and Normalization (25 marks) . . . . .	12
4.3 Section C: Python Script Debugging (20 marks) . . . . .	12
4.4 Section D: RAG Workflow Design (15 marks) . . . . .	13
4.5 Section E: Model Selection Scenarios (20 marks) . . . . .	13
<b>5 PRACTICE PAPER SET 2 - ANSWERS</b>	<b>13</b>
5.1 Section A: JSON Data Debugging - ANSWERS . . . . .	13
5.2 Section B: SQL Normalization - ANSWERS . . . . .	15
5.3 Section C: Python Script Debugging - ANSWERS . . . . .	16
5.4 Section D: RAG Workflow - ANSWERS . . . . .	17
5.5 Section E: Model Selection - ANSWERS . . . . .	18
<b>6 PRACTICE PAPER SET 3</b>	<b>19</b>
6.1 Section A: JSON Data Debugging (20 marks) . . . . .	19
6.2 Section B: SQL Debugging and Normalization (25 marks) . . . . .	20
6.3 Section C: Python Script Debugging (20 marks) . . . . .	20
6.4 Section D: Embedding Demonstration (15 marks) . . . . .	21
6.5 Section E: Model Selection Scenarios (20 marks) . . . . .	21
<b>7 PRACTICE PAPER SET 3 - ANSWERS</b>	<b>22</b>
7.1 Section A: JSON Data Debugging - ANSWERS . . . . .	22
7.2 Section B: SQL Normalization - ANSWERS . . . . .	23

7.3	Section C: Python Script Debugging - ANSWERS . . . . .	24
7.4	Section D: Embedding Demonstration - ANSWERS . . . . .	26
7.5	Section E: Model Selection - ANSWERS . . . . .	27
<b>8</b>	<b>PRACTICE PAPER SET 4</b>	<b>28</b>
8.1	Section A: JSON Data Debugging (20 marks) . . . . .	28
8.2	Section B: SQL Debugging and Normalization (25 marks) . . . . .	29
8.3	Section C: Python Script Debugging (20 marks) . . . . .	29
8.4	Section D: RAG Workflow Design (15 marks) . . . . .	30
8.5	Section E: Model Selection Scenarios (20 marks) . . . . .	30
<b>9</b>	<b>PRACTICE PAPER SET 4 - ANSWERS</b>	<b>31</b>
9.1	Section A: JSON Data Debugging - ANSWERS . . . . .	31
9.2	Section B: SQL Normalization - ANSWERS . . . . .	32
9.3	Section C: Python Script Debugging - ANSWERS . . . . .	34
9.4	Section D: RAG Workflow - ANSWERS . . . . .	35
9.5	Section E: Model Selection - ANSWERS . . . . .	36
<b>10</b>	<b>PRACTICE PAPER SET 5</b>	<b>37</b>
10.1	Section A: JSON Data Debugging (20 marks) . . . . .	37
10.2	Section B: SQL Debugging and Normalization (25 marks) . . . . .	38
10.3	Section C: Python Script Debugging (20 marks) . . . . .	39
10.4	Section D: Complete RAG + Embedding Task (15 marks) . . . . .	40
10.5	Section E: Model Selection Scenarios (20 marks) . . . . .	40
<b>11</b>	<b>PRACTICE PAPER SET 5 - ANSWERS</b>	<b>41</b>
11.1	Section A: JSON Data Debugging - ANSWERS . . . . .	41
11.2	Section B: SQL Normalization - ANSWERS . . . . .	42
11.3	Section C: Python Script Debugging - ANSWERS . . . . .	43
11.4	Section D: RAG + Embedding System - ANSWERS . . . . .	45
11.5	Section E: Model Selection - ANSWERS . . . . .	46

# 1 KELP GLOBAL - JUNIOR DATA SCIENTIST PRACTICE PAPERS

## 1.1 Understanding the Format (Based on Senior's Feedback)

Your seniors mentioned that the interview was **tough** and included:

- Dataset with tasks like bug finding, providing solutions, and explaining approach
- Embedding demonstration
- RAG workflow design
- SQL query debugging and normalization
- JSON structural flaw detection and fixing
- Python script error identification and correction
- Scenario-based model selection with rules definition
- **No DSA questions**

## 2 PRACTICE PAPER SET 1

### 2.1 Section A: JSON Data Debugging (20 marks)

**Question 1:** You are given the following JSON data representing customer orders. Identify ALL structural flaws, explain what errors they can cause, and provide the corrected JSON.

```
1 {  
2     'customer_id': 1001,  
3     "name": "Rahul Sharma"  
4     "orders": [  
5         {"order_id": "ORD001", "amount": 2500.00, "status": "delivered"},  
6         {"order_id": "ORD002", "amount": undefined, "status": "pending"},  
7         {"order_id": "ORD002", "amount": 1800, "status": "shipped"}  
8     ],  
9     "email": 'rahul@example.com',  
10    "phone": +919876543210,  
11    "address": {  
12        "city": "Mumbai"  
13        "pincode": 400001  
14    }  
15 }
```

#### Tasks:

- List all structural flaws and what errors they will cause (8 marks)
- Provide the corrected JSON (6 marks)
- How would you validate this JSON programmatically in Python? (6 marks)

### 2.2 Section B: SQL Debugging and Normalization (25 marks)

**Question 2:** Given the following table with data:

EmpID	EmpName	Dept	DeptHead	Projects	Salary	JoinDate
101	Amit Kumar	IT	Priya Singh	Proj1,Proj2	50000	2023-15-01
102	Sneha Patel	HR	Rahul Joshi	Proj3	NULL	01-Mar-2022
101	Amit Kumar	IT	Priya Singh	Proj1,Proj2	55000	2023-15-01

#### Tasks:

- Identify incorrect data and structural issues in this table (5 marks)
- Explain the anomalies present (insertion, update, deletion) (5 marks)
- Normalize this table to 3NF with proper schema design (8 marks)
- Write SQL query to create normalized tables with appropriate constraints (7 marks)

### 2.3 Section C: Python Script Debugging (20 marks)

**Question 3:** The following Python script has multiple errors. Identify the line numbers where errors will occur, explain why, and provide the corrected script.

```

1 1. import pandas as pd
2 2. import numpy as np
3 3.
4 4. def analyze_sales(filepath)
5 5.     df = pd.read_csv(filepath)
6 6.
7 7.     # Clean the data
8 8.     df['revenue'] = df['price'] * df[quantity]
9 9.
10 10.    # Handle missing values
11 11.    df.fillna(method='ffill', inplace=True)
12 12.
13 13.    # Calculate statistics
14 14.    avg_revenue = df['revenue'].mean
15 15.    max_sale = df.loc[df['revenue'] == df['revenue'].max()]
16 16.
17 17.    # Group by category
18 18.    category_sales = df.groupby('category').agg({
19 19.        'revenue': ['sum', 'mean']
20 20.        'quantity': 'count'
21 21.    })
22 22.
23 23.    return avg_revenue max_sale, category_sales
24 24.
25 25. if __name__ == "__main__":
26 26.     result = analyze_sales("sales_data.csv")
27 27.     Print(result)

```

#### Tasks:

- Identify all errors with line numbers and reasons (12 marks)
- Rewrite the corrected script (8 marks)

## 2.4 Section D: RAG Workflow Design (15 marks)

**Question 4:** You are building a customer support chatbot for Kelp Global that answers questions about private equity investments using company documents.

**Tasks:**

- (a) Draw a complete RAG workflow diagram showing all components (8 marks)
- (b) Explain the role of each component (embedding model, vector database, retriever, LLM) (7 marks)

## 2.5 Section E: Model Selection Scenarios (20 marks)

**Question 5:** For each scenario below, specify which model/technique you would use, provide 3 rules/reasons for your choice, and mention any preprocessing required.

Scenario	Your Answer
S1: Predicting customer churn (yes/no) based on transaction history	
S2: Grouping customers based on purchasing behavior (no labels available)	
S3: Predicting house prices based on features like area, location, bedrooms	
S4: Detecting anomalies in credit card transactions	
S5: Classifying customer reviews into 5 sentiment categories	
S6: Building a recommendation system for products	

**Format for each scenario:**

- Model/Technique chosen
- 3 Rules/Reasons for selection
- Key preprocessing steps

## 3 PRACTICE PAPER SET 1 - ANSWERS

### 3.1 Section A: JSON Data Debugging - ANSWERS

**(a) Structural Flaws Identified:**

1. **Line 2: Single quotes used for key** - 'customer\_id' uses single quotes instead of double quotes
  - Error: JSON parsing will fail with "Unexpected token"
2. **Line 3: Missing comma** - No comma after "Rahul Sharma"
  - Error: Syntax error, unexpected string
3. **Line 5: Trailing comma** - Comma after "delivered",
  - Error: JSON doesn't allow trailing commas

- 4. **Line 6: Undefined value** - "amount": undefined
  - Error: undefined is not a valid JSON data type
- 5. **Line 7: Duplicate key** - "order\_id": "ORD002" appears twice
  - Error: Data ambiguity, unpredictable parsing behavior
- 6. **Line 9: Single quotes for value** - 'rahul@example.com'
  - Error: JSON requires double quotes for strings
- 7. **Line 10: Unquoted number with special character** - +919876543210
  - Error: Phone number should be string, '+' makes it invalid number
- 8. **Line 12: Missing comma** - No comma between "Mumbai" and "pincode"
  - Error: Syntax error

**(b) Corrected JSON:**

```
1 {
2     "customer_id": 1001,
3     "name": "Rahul Sharma",
4     "orders": [
5         {"order_id": "ORD001", "amount": 2500.00, "status": "delivered"},
6         {"order_id": "ORD002", "amount": null, "status": "pending"},
7         {"order_id": "ORD003", "amount": 1800, "status": "shipped"}
8     ],
9     "email": "rahul@example.com",
10    "phone": "+919876543210",
11    "address": {
12        "city": "Mumbai",
13        "pincode": 400001
14    }
15 }
```

**(c) Python Validation Code:**

```
1 import json
2
3 def validate_json(json_string):
4     try:
5         data = json.loads(json_string)
6         print("Valid JSON")
7         return True, data
8     except json.JSONDecodeError as e:
9         print(f"Invalid JSON: {e.msg} at line {e.lineno}, column {e.colno}")
10        return False, None
11
12 # Using jsonschema for structure validation
13 from jsonschema import validate, ValidationError
14
15 schema = {
16     "type": "object",
17     "required": ["customer_id", "name", "orders"],
18     "properties": {
19         "customer_id": {"type": "integer"},
20         "name": {"type": "string"},
21         "orders": {
```

```

22         "type": "array",
23         "items": {
24             "type": "object",
25             "properties": {
26                 "order_id": {"type": "string"},
27                 "amount": {"type": ["number", "null"]},
28                 "status": {"type": "string"}
29             }
30         }
31     }
32 }
33 }

```

---

## 3.2 Section B: SQL Debugging and Normalization - ANSWERS

### (a) Incorrect Data and Structural Issues:

1. **Duplicate EmpID:** EmpID 101 appears twice with different salaries
2. **Multi-valued attribute:** Projects column contains "Proj1,Proj2" - violates 1NF
3. **Invalid date format:** "2023-15-01" has invalid month (15)
4. **Inconsistent date formats:** Different formats used (YYYY-MM-DD vs DD-Mon-YYYY)
5. **NULL value:** Salary is NULL for EmpID 102
6. **Transitive dependency:** DeptHead depends on Dept, not directly on EmpID

### (b) Anomalies Present:

- **Insertion Anomaly:** Cannot add new department without adding employee
- **Update Anomaly:** If DeptHead changes, must update multiple rows
- **Deletion Anomaly:** Deleting last employee in dept loses dept info

### (c) Normalized Schema (3NF):

#### Employees Table:

```
1 Employees(EmpID PK, EmpName, DeptID FK, Salary, JoinDate)
```

#### Departments Table:

```
1 Departments(DeptID PK, DeptName, DeptHead)
```

#### Projects Table:

```
1 Projects(ProjectID PK, ProjectName)
```

#### Employee\_Projects Table:

```
1 Employee_Projects(EmpID FK, ProjectID FK) - Composite PK
```

### (d) SQL to Create Normalized Tables:

```

1  -- Departments Table
2  CREATE TABLE Departments (
3      DeptID INT PRIMARY KEY AUTO_INCREMENT,
4      DeptName VARCHAR(50) NOT NULL UNIQUE,
5      DeptHead VARCHAR(100)
6  );
7
8  -- Employees Table
9  CREATE TABLE Employees (
10     EmpID INT PRIMARY KEY,
11     EmpName VARCHAR(100) NOT NULL,
12     DeptID INT,
13     Salary DECIMAL(10,2) NOT NULL DEFAULT 0,
14     JoinDate DATE NOT NULL,
15     FOREIGN KEY (DeptID) REFERENCES Departments(DeptID),
16     CONSTRAINT chk_salary CHECK (Salary >= 0)
17 );
18
19 -- Projects Table
20 CREATE TABLE Projects (
21     ProjectID INT PRIMARY KEY AUTO_INCREMENT,
22     ProjectName VARCHAR(100) NOT NULL
23 );
24
25 -- Employee_Projects Junction Table
26 CREATE TABLE Employee_Projects (
27     EmpID INT,
28     ProjectID INT,
29     PRIMARY KEY (EmpID, ProjectID),
30     FOREIGN KEY (EmpID) REFERENCES Employees(EmpID),
31     FOREIGN KEY (ProjectID) REFERENCES Projects(ProjectID)
32 );

```

### 3.3 Section C: Python Script Debugging - ANSWERS

#### (a) Errors Identified:

Line	Error	Reason
4	Missing colon	Function definition needs ':' at end
8	Missing quotes	'df[quantity]' should be 'df['quantity']'
11	Missing comma	Between 'ffill' and 'inplace=True'
14	Missing parentheses	'mean' is a method, needs '()'
19-20	Missing comma	Between dictionaries in 'agg()'
23	Missing comma	Between return values
25	Single '=' used	Should be '==' for comparison
27	Capital P	'Print' should be 'print'

#### (b) Corrected Script:

```

1  import pandas as pd
2  import numpy as np

```



```

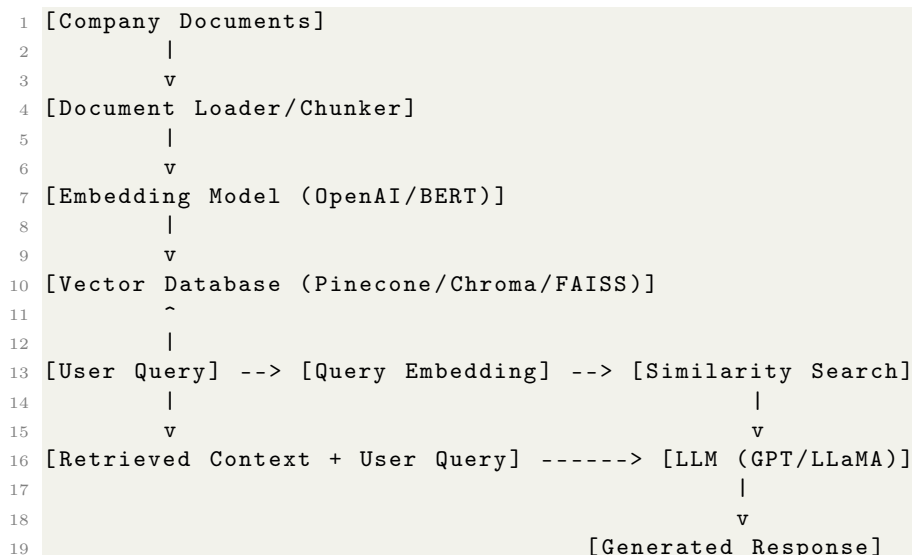
3
4 def analyze_sales(filepath):
5     df = pd.read_csv(filepath)
6
7     # Clean the data
8     df['revenue'] = df['price'] * df['quantity']
9
10    # Handle missing values
11    df.fillna(method='ffill', inplace=True)
12
13    # Calculate statistics
14    avg_revenue = df['revenue'].mean()
15    max_sale = df.loc[df['revenue'] == df['revenue'].max()]
16
17    # Group by category
18    category_sales = df.groupby('category').agg({
19        'revenue': ['sum', 'mean'],
20        'quantity': 'count'
21    })
22
23    return avg_revenue, max_sale, category_sales
24
25 if __name__ == "__main__":
26     result = analyze_sales("sales_data.csv")
27     print(result)

```

---

### 3.4 Section D: RAG Workflow Design - ANSWERS

#### (a) RAG Workflow Diagram:



#### (b) Component Roles:

1. **Embedding Model:** Converts text chunks and queries into dense vector representations (embeddings) that capture semantic meaning. Similar concepts have vectors close in space.
2. **Vector Database:** Stores document embeddings and enables efficient similarity search.

Uses approximate nearest neighbor algorithms for fast retrieval.

3. **Retriever:** Fetches the most relevant document chunks based on query embedding similarity. Determines the context window for the LLM.

4. **LLM:** Takes the retrieved context and user query to generate coherent, contextually appropriate responses. Grounds responses in retrieved information.

### 3.5 Section E: Model Selection Scenarios - ANSWERS

#### S1: Customer Churn Prediction (Binary Classification)

- **Model:** Logistic Regression, Random Forest, or XGBoost

- **Rules:**

1. Binary outcome (churn/no churn) suits classification algorithms
2. Random Forest handles non-linear relationships and feature importance
3. XGBoost provides high accuracy with imbalanced datasets

- **Preprocessing:** Handle class imbalance (SMOTE), feature scaling, encode categoricals

#### S2: Customer Grouping (Clustering)

- **Model:** K-Means Clustering or DBSCAN

- **Rules:**

1. No labels available - unsupervised learning required
2. K-Means works well with spherical clusters of similar sizes
3. Elbow method helps determine optimal number of clusters

- **Preprocessing:** Feature scaling (StandardScaler), dimensionality reduction (PCA)

#### S3: House Price Prediction (Regression)

- **Model:** Linear Regression, Ridge/Lasso, or Gradient Boosting Regressor

- **Rules:**

1. Continuous target variable (price) requires regression
2. Ridge/Lasso handles multicollinearity and feature selection
3. Gradient Boosting captures non-linear patterns

- **Preprocessing:** Handle missing values, encode location, feature engineering

#### S4: Anomaly Detection in Transactions

- **Model:** Isolation Forest or One-Class SVM

- **Rules:**

1. Anomalies are rare events - one-class classification needed
2. Isolation Forest efficiently isolates outliers
3. Doesn't require labeled fraud examples for training

- **Preprocessing:** Feature engineering on transaction patterns, scaling

#### S5: Sentiment Classification (Multi-class)

- **Model:** BERT or Random Forest with TF-IDF

- **Rules:**

1. Multi-class text classification - 5 sentiment categories
2. BERT captures contextual word embeddings for better understanding
3. Pre-trained models reduce training data requirements

- **Preprocessing:** Text cleaning, tokenization, handling class imbalance

### S6: Product Recommendation System

- **Model:** Collaborative Filtering or Content-Based + Matrix Factorization

- **Rules:**

1. Collaborative filtering uses user-item interaction patterns
2. Matrix factorization handles sparse user-item matrices
3. Hybrid approach combines strengths of multiple techniques

- **Preprocessing:** Create user-item matrix, handle cold start problem

## 4 PRACTICE PAPER SET 2

### 4.1 Section A: JSON Data Debugging (20 marks)

**Question 1:** Analyze the following JSON representing API response for employee data. Find all flaws and fix them.

```
1 {
2   "api_version": 2.1,
3   "timestamp": 2024-12-02T10:30:00,
4   "status_code": "200",
5   "employees": [
6     {
7       "id": 001,
8       "name": "Priya Mehta",
9       "department": "Data Science",
10      "skills": ["Python", "SQL", "ML",]
11      "salary": 75000,
12      "manager_id": null
13    },
14    {
15      "id": 002,
16      "name": 'Vikram Singh',
17      "department": "Data Science",
18      "skills": ["Python", "Tableau"],
19      "salary": NaN,
20      "manager_id": 001
21    }
22  ]
23  "metadata": {
24    "total_count": 2,
25    "page": 1,
26  }
27 }
```

#### Tasks:

- (a) List all structural flaws with potential error messages (8 marks)
- (b) Provide corrected JSON (6 marks)
- (c) Write a Python function to detect and report common JSON issues (6 marks)

## 4.2 Section B: SQL Debugging and Normalization (25 marks)

**Question 2:** Analyze the following product inventory table:

ProductID	ProductName	Category	SupplierName	SupplierContact	OrderHistory
P001	Laptop Dell	Electronics	TechMart	9876543210, tech@mart.com	ORD1:10,000
P002	iPhone 15	Electronics,Mobile	AppleStore	NULL	ORD3:2
P001	Laptop Dell	Electronics	TechMart	9876543210, tech@mart.com	ORD1:10,000

### Tasks:

- Identify data integrity and structural issues (6 marks)
- Explain why this violates normal forms (5 marks)
- Design normalized schema up to 3NF with relationships (7 marks)
- Write SQL queries for creating tables and a query to find products with negative stock (7 marks)

## 4.3 Section C: Python Script Debugging (20 marks)

**Question 3:** Debug the following data preprocessing script:

```
1 1. import pandas as pd
2 2. from sklearn.preprocessing import StandardScaler
3 3. from sklearn.model_selection import train_test_split
4 4.
5 5. class DataProcessor
6 6.     def __init__(self, data_path):
7 7.         self.df = pd.read_csv(data_path)
8 8.         self.scaler = StandardScaler[]
9 9.
10 10.     def clean_data(self):
11 11.         # Remove duplicates
12 12.         self.df.drop_duplicates(inplace=True)
13 13.
14 14.         # Handle missing values
15 15.         for col in self.df.columns:
16 16.             if self.df[col].dtype == 'object':
17 17.                 self.df[col].fillna('Unknown', inplace=True)
18 18.             else:
19 19.                 self.df[col].fillna(self.df[col].median, inplace=True)
20 20.
21 21.     def encode_features(self, categorical_cols):
22 22.         for col in categorical_cols:
23 23.             self.df[col] = pd.get_dummies(self.df[col])
24 24.
25 25.     def split_data(self, target_col, test_size=0.2):
26 26.         X = self.df.drop(target_col)
27 27.         y = self.df[target_col]
28 28.         return train_test_split(X, y, test_size=test_size, random_state=42)
29 29.
```

```

30 30. if __name__ == '__main__':
31     processor = DataProcessor('data.csv')
32     processor.clean_data()
33 33. X_train, X_test, y_train, y_test = processor.split_data['target']

```

#### Tasks:

- Identify all errors with line numbers and explanations (12 marks)
- Provide the corrected script (8 marks)

### 4.4 Section D: RAG Workflow Design (15 marks)

**Question 4:** Design a RAG system for a financial document Q&A application that handles PDF reports, Excel data, and regulatory documents.

#### Tasks:

- Draw a detailed RAG architecture showing document processing pipeline (8 marks)
- Explain chunking strategies and their trade-offs for financial documents (4 marks)
- How would you handle numerical data and tables in this RAG system? (3 marks)

### 4.5 Section E: Model Selection Scenarios (20 marks)

**Question 5:** For each scenario, choose appropriate technique and justify with 3 rules.

Scenario	Description
S1	Predicting loan default amount (continuous value) for defaulting customers
S2	Identifying topics from a large collection of news articles (no predefined topics)
S3	Real-time fraud detection where model needs to update with new patterns
S4	Predicting next month's sales from historical time-series data
S5	Building a spam classifier with very limited labeled data (100 examples)
S6	Image classification to identify product defects in manufacturing

## 5 PRACTICE PAPER SET 2 - ANSWERS

### 5.1 Section A: JSON Data Debugging - ANSWERS

#### (a) Structural Flaws:

- Line 3:** Unquoted timestamp - 2024-12-02T10:30:00 should be string "2024-12-02T10:30:00"
- Line 5:** Leading zeros in number - 001 is invalid; use 1 or "001"
- Line 8:** Trailing comma in array - ["Python", "SQL", "ML",,]
- Line 8:** Missing comma - After skills array, before salary

5. **Line 13:** Single quotes - 'Vikram Singh' should use double quotes
6. **Line 13:** Missing comma - After name value
7. **Line 16:** NaN is not valid - Use null instead
8. **Line 19:** Missing comma - After employees array, before metadata
9. **Line 22:** Trailing comma - After "page": 1,

(b) Corrected JSON:

```

1 {
2     "api_version": 2.1,
3     "timestamp": "2024-12-02T10:30:00",
4     "status_code": "200",
5     "employees": [
6         {
7             "id": 1,
8             "name": "Priya Mehta",
9             "department": "Data Science",
10            "skills": ["Python", "SQL", "ML"],
11            "salary": 75000,
12            "manager_id": null
13        },
14        {
15            "id": 2,
16            "name": "Vikram Singh",
17            "department": "Data Science",
18            "skills": ["Python", "Tableau"],
19            "salary": null,
20            "manager_id": 1
21        }
22    ],
23    "metadata": {
24        "total_count": 2,
25        "page": 1
26    }
27 }
```

(c) Python Function to Detect JSON Issues:

```

1 import json
2 import re
3
4 def detect_json_issues(json_string):
5     issues = []
6
7     # Check for single quotes
8     if re.search(r'(?<!\)')',', json_string):
9         issues.append("Single quotes detected - use double quotes")
10
11    # Check for trailing commas
12    if re.search(r',\s*[\}\]]', json_string):
13        issues.append("Trailing comma detected")
14
15    # Check for undefined/NaN
16    if 'undefined' in json_string:
17        issues.append("'undefined' is not valid JSON - use null")
18    if re.search(r':\s*NaN', json_string):
19        issues.append("NaN is not valid JSON - use null")
20
```

```

21 # Check for unquoted strings/dates
22 if re.search(r':\s*\d{4}-\d{2}-\d{2}T', json_string):
23     issues.append("Possible unquoted datetime detected")
24
25 # Try parsing
26 try:
27     json.loads(json_string)
28     issues.append("JSON parsed successfully")
29 except json.JSONDecodeError as e:
30     issues.append(f"Parse error: {e.msg} at line {e.lineno}")
31
32 return issues

```

---

## 5.2 Section B: SQL Normalization - ANSWERS

### (a) Data Integrity Issues:

1. **Duplicate ProductID:** P001 appears twice with different prices
2. **Multi-valued attributes:** Category has "Electronics,Mobile", OrderHistory has multiple orders
3. **Composite data:** SupplierContact contains phone AND email together
4. **Negative stock:** StockQty = -5 (invalid business rule)
5. **NULL supplier contact:** Data incompleteness

### (b) Normal Form Violations:

- **1NF Violation:** Multi-valued attributes (Category, OrderHistory, SupplierContact)
- **2NF Violation:** Supplier info depends on SupplierName, not ProductID
- **3NF Violation:** SupplierContact transitively depends on SupplierName

### (c) Normalized Schema (3NF):

```

1 Products(ProductID PK, ProductName, Price, StockQty)
2 Categories(CategoryID PK, CategoryName)
3 Product_Categories(ProductID FK, CategoryID FK) - Composite PK
4 Suppliers(SupplierID PK, SupplierName, Phone, Email)
5 Product_Suppliers(ProductID FK, SupplierID FK)
6 Orders(OrderID PK, ProductID FK, Quantity, OrderDate)

```

### (d) SQL Implementation:

```

1 CREATE TABLE Products (
2     ProductID VARCHAR(10) PRIMARY KEY,
3     ProductName VARCHAR(100) NOT NULL,
4     Price DECIMAL(10,2) NOT NULL CHECK (Price > 0),
5     StockQty INT DEFAULT 0 CHECK (StockQty >= 0)
6 );
7
8 CREATE TABLE Categories (
9     CategoryID INT PRIMARY KEY AUTO_INCREMENT,
10    CategoryName VARCHAR(50) UNIQUE NOT NULL

```

```

11 );
12
13 CREATE TABLE Product_Categories (
14     ProductID VARCHAR(10),
15     CategoryID INT,
16     PRIMARY KEY (ProductID, CategoryID),
17     FOREIGN KEY (ProductID) REFERENCES Products(ProductID),
18     FOREIGN KEY (CategoryID) REFERENCES Categories(CategoryID)
19 );
20
21 CREATE TABLE Suppliers (
22     SupplierID INT PRIMARY KEY AUTO_INCREMENT,
23     SupplierName VARCHAR(100) NOT NULL,
24     Phone VARCHAR(15),
25     Email VARCHAR(100)
26 );
27
28 -- Query for negative stock (would be prevented by CHECK constraint)
29 SELECT ProductID, ProductName, StockQty
30 FROM Products
31 WHERE StockQty < 0;

```

## 5.3 Section C: Python Script Debugging - ANSWERS

### (a) Errors:

Line	Error	Explanation
5	Missing colon	'class DataProcessor' needs ':'
8	Square brackets	'StandardScaler[]' should be 'StandardScaler()'
12	Lowercase true	Python boolean is 'True', not 'true'
15	Missing colon	'for' loop needs ':' at end
19	Missing parentheses	'median' is method, needs '()'
21	Missing colon	Function definition needs ':'
23	Wrong usage	'pd.get_dummies' creates DataFrame, not fit for column assignment
26	Missing axis	'drop()' needs 'axis=1' for columns
28	Missing comma	Between 'y' and 'test_size'
33	Square brackets	Function call should use '()' not '[]'

### (b) Corrected Script:

```

1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler, LabelEncoder
3 from sklearn.model_selection import train_test_split
4
5 class DataProcessor:
6     def __init__(self, data_path):
7         self.df = pd.read_csv(data_path)
8         self.scaler = StandardScaler()
9
10    def clean_data(self):
11        # Remove duplicates

```



```

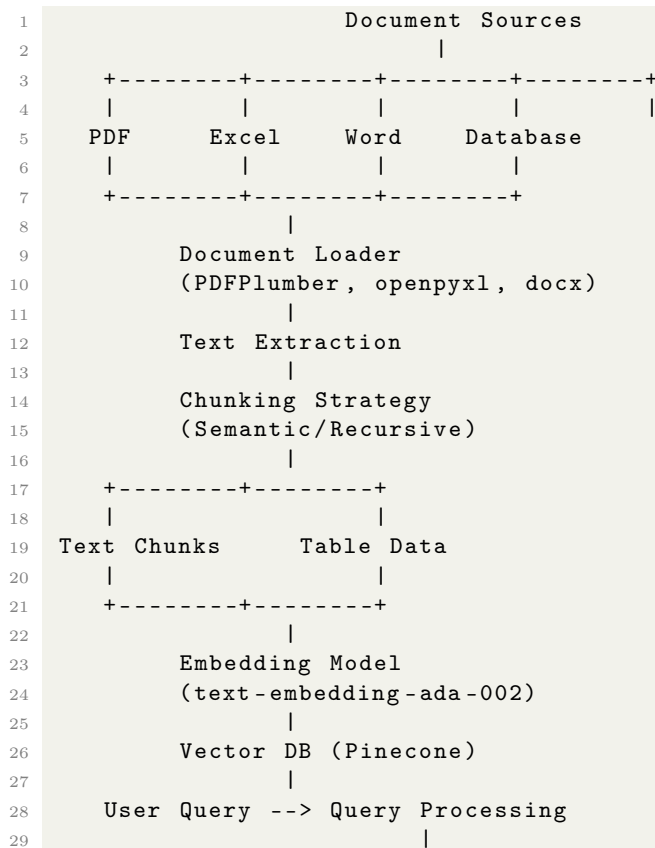
12     self.df.drop_duplicates(inplace=True)
13
14     # Handle missing values
15     for col in self.df.columns:
16         if self.df[col].dtype == 'object':
17             self.df[col].fillna('Unknown', inplace=True)
18         else:
19             self.df[col].fillna(self.df[col].median(), inplace=True)
20
21     def encode_features(self, categorical_cols):
22         le = LabelEncoder()
23         for col in categorical_cols:
24             self.df[col] = le.fit_transform(self.df[col])
25
26     def split_data(self, target_col, test_size=0.2):
27         X = self.df.drop(target_col, axis=1)
28         y = self.df[target_col]
29         return train_test_split(X, y, test_size=test_size, random_state=42)
30
31 if __name__ == '__main__':
32     processor = DataProcessor('data.csv')
33     processor.clean_data()
34     X_train, X_test, y_train, y_test = processor.split_data('target')

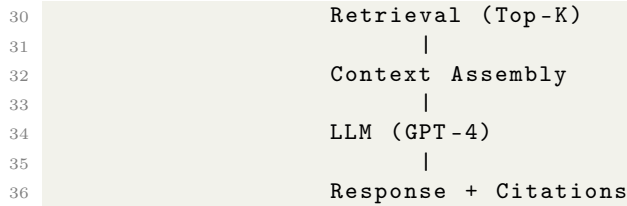
```

---

## 5.4 Section D: RAG Workflow - ANSWERS

### (a) Financial RAG Architecture:





### (b) Chunking Strategies:

1. **Fixed-size chunking:** Simple but may break semantic units
2. **Semantic chunking:** Respects paragraph/section boundaries, better context
3. **Recursive splitting:** Tries multiple separators (paragraphs → sentences → words)
4. **Document-aware chunking:** Preserves table structures and financial metrics together

**Trade-offs:** Smaller chunks = more precise retrieval but less context; Larger chunks = more context but may include irrelevant info

### (c) Handling Numerical Data:

- Extract tables separately using specialized parsers
- Store tables as structured JSON in metadata
- Use hybrid search (keyword + semantic) for numerical queries
- Consider using SQL interface for complex numerical queries

## 5.5 Section E: Model Selection - ANSWERS

### S1: Loan Default Amount Prediction

- **Model:** Gradient Boosting Regressor / XGBoost Regressor
- **Rules:**
  1. Continuous target (amount) requires regression
  2. XGBoost handles non-linear relationships in financial data
  3. Built-in handling for missing values common in loan data

### S2: Topic Identification (Topic Modeling)

- **Model:** LDA (Latent Dirichlet Allocation) or BERTopic
- **Rules:**
  1. No predefined labels - unsupervised approach needed
  2. LDA identifies latent topics from word distributions
  3. BERTopic uses transformers for better semantic grouping

### S3: Real-time Fraud Detection

- **Model:** Online Learning (SGD Classifier) or Streaming ML
- **Rules:**
  1. Model must update incrementally with new data
  2. SGD supports `partial_fit` for online learning
  3. Low latency requirements suit linear models

### S4: Time-Series Sales Prediction

- **Model:** ARIMA, Prophet, or LSTM

- **Rules:**

1. Temporal dependencies require time-series specific models
2. Prophet handles seasonality and holidays well
3. LSTM captures long-term patterns in sequential data

**S5: Spam Classification (Limited Data)**

- **Model:** Transfer Learning (Fine-tuned BERT) or Few-shot Learning

- **Rules:**

1. Limited data (100 samples) - pre-trained models essential
2. BERT's language understanding transfers to spam detection
3. Data augmentation can expand training set

**S6: Manufacturing Defect Detection**

- **Model:** CNN (Convolutional Neural Network) or Transfer Learning (ResNet)

- **Rules:**

1. Image data requires CNN architecture
2. Pre-trained models (ResNet, VGG) work well with limited manufacturing images
3. Can detect visual patterns humans might miss

## 6 PRACTICE PAPER SET 3

### 6.1 Section A: JSON Data Debugging (20 marks)

**Question 1:** Given the following JSON configuration for a machine learning pipeline. Find and fix all issues.

```
1 {
2   "pipeline_config": {
3     "model_type": "RandomForest",
4     "hyperparameters": {
5       "n_estimators": 100,
6       "max_depth": None,
7       "min_samples_split": "2",
8       "random_state": 42.0
9     },
10    "features": [
11      "age"
12      "income"
13      "credit_score",
14    ],
15    "target": 'churn',
16    "preprocessing": {
17      "scaling": true,
18      "encoding": "onehot",
19      "handle_missing": {
20        "strategy": "mean"
21        "columns": ["age", "income"]
22      }
23    }
24  },
```

```

25     "training_params": {
26         "epochs": 50,
27         "batch_size": 32,
28         "learning_rate": .001,
29         "callbacks": [EarlyStopping, ModelCheckpoint]
30     }
31 }

```

**Tasks:**

- Identify all JSON flaws and their potential impacts (10 marks)
- Provide the corrected JSON (6 marks)
- How would you implement JSON schema validation for this config? (4 marks)

## 6.2 Section B: SQL Debugging and Normalization (25 marks)

**Question 2:** Consider the following analytics data table:

LogID	UserID	UserEmail	EventType	EventData	SessionID	DeviceType
1	U001	user1@test.com	click	{"page":"home","button":"signup"}	S001	Mobile
2	U001	user1@test.com	view	home	S001	Mobile
3	NULL	guest@test.com	click	{}	S002	Desktop

**Tasks:**

- Identify all data quality issues and inconsistencies (6 marks)
- What anomalies exist in this design? (5 marks)
- Design a star schema or normalized schema for analytics (7 marks)
- Write SQL to find users with multiple devices and sessions (7 marks)

## 6.3 Section C: Python Script Debugging (20 marks)

**Question 3:** Fix the following ML embedding and prediction script:

```

1 1. import numpy as np
2 2. import pandas as pd
3 3. from sklearn.ensemble import RandomForestClassifier
4 4. from gensim.models import Word2Vec
5 5.
6 6. def create_embeddings(texts):
7 7.     # Tokenize texts
8 8.     tokenized = [text.split() for text in texts]
9 9.
10 10.    # Train Word2Vec
11 11.    model = Word2Vec(tokenized, vector_size=100, window=5)
12 12.
13 13.    # Get document embeddings
14 14.    embeddings = []

```

```

15 15.     for tokens in tokenized
16 16.         doc_vec = np.mean([model.wv[token] for token in tokens])
17 17.         embeddings.append(doc_vec)
18 18.
19 19.     return np.array(embeddings)
20 20.
21 21. def train_model(X, y):
22 22.     clf = RandomForestClassifier(n_estimators=100)
23 23.     clf.fit(X, y)
24 24.
25 25.     # Get feature importance
26 26.     importance = clf.feature_importances_
27 27.     print(f"Top features: {importance.argsort()[-5:]}")
28 28.
29 29.     return clf
30 30.
31 31. if __name__ == '__main__':
32 32.     df = pd.read_csv('data.csv')
33 33.
34 34.     # Create embeddings
35 35.     X = create_embeddings(df['text'])
36 36.     y = df['label']
37 37.
38 38.     # Train
39 39.     model = train_model(X y)
40 40.
41 41.     # Predict
42 42.     predictions = model.predict[X[:10]]
43 43.     print(f"Predictions: (predictions)")

```

#### Tasks:

- Identify all errors with explanations (12 marks)
- Rewrite the corrected script with proper error handling (8 marks)

## 6.4 Section D: Embedding Demonstration (15 marks)

**Question 4:** Explain and demonstrate word embeddings.

#### Tasks:

- Explain the difference between Word2Vec (CBOW vs Skip-gram), GloVe, and BERT embeddings (5 marks)
- Write Python code to generate embeddings for a set of product descriptions using a pre-trained model (6 marks)
- How would you use embeddings for a product similarity search? (4 marks)

## 6.5 Section E: Model Selection Scenarios (20 marks)

**Question 5:** Choose appropriate techniques for each scenario with 3 rules each.

Scenario	Description
S1	Predicting employee attrition with highly imbalanced data (5% attrition rate)
S2	Detecting fake reviews among millions of product reviews
S3	Building a chatbot that needs to understand user intent
S4	Predicting stock price direction (up/down) for next day
S5	Identifying objects in satellite imagery for urban planning
S6	Matching job candidates to job postings based on skills

## 7 PRACTICE PAPER SET 3 - ANSWERS

### 7.1 Section A: JSON Data Debugging - ANSWERS

(a) JSON Flaws Identified:

Issue	Line	Impact
'None' (Python keyword)	hyperparameters	Invalid - use 'null'
"2" as string	min_samples_split	Should be integer '2'
'42.0' float for seed	random_state	Works but typically integer
Missing commas	features array	Syntax error
Trailing comma	features array	Invalid JSON
Single quotes	target	Must use double quotes
'true' lowercase	scaling	Valid in JSON, but note case
Missing comma	handle_missing	Between strategy and columns
'0.001' without leading zero	learning_rate	Valid but inconsistent
Unquoted values	callbacks array	Variable names not valid

(b) Corrected JSON:

```

1 {
2   "pipeline_config": {
3     "model_type": "RandomForest",
4     "hyperparameters": {
5       "n_estimators": 100,
6       "max_depth": null,
7       "min_samples_split": 2,
8       "random_state": 42
9     },
10    "features": [
11      "age",
12      "income",
13      "credit_score"
14    ],
15    "target": "churn",
16    "preprocessing": {
17      "scaling": true,
18      "encoding": "onehot",
19      "handle_missing": {
20        "strategy": "mean",

```

```

21         "columns": ["age", "income"]
22     }
23 },
24 },
25 "training_params": {
26     "epochs": 50,
27     "batch_size": 32,
28     "learning_rate": 0.001,
29     "callbacks": ["EarlyStopping", "ModelCheckpoint"]
30 }
31 }

```

### (c) JSON Schema Validation:

```

1 from jsonschema import validate
2
3 schema = {
4     "type": "object",
5     "required": ["pipeline_config", "training_params"],
6     "properties": {
7         "pipeline_config": {
8             "type": "object",
9             "required": ["model_type", "hyperparameters", "features", "target"]
10        },
11        "properties": {
12            "model_type": {"type": "string", "enum": ["RandomForest", "XGBoost", "LogisticRegression"]},
13            "hyperparameters": {
14                "type": "object",
15                "properties": {
16                    "n_estimators": {"type": "integer", "minimum": 1},
17                    "max_depth": {"type": ["integer", "null"]},
18                    "min_samples_split": {"type": "integer", "minimum": 2}
19                }
20            },
21            "features": {"type": "array", "items": {"type": "string"}},
22            "target": {"type": "string"}
23        }
24    }
25 }
26
27 validate(instance=config_data, schema=schema)

```

## 7.2 Section B: SQL Normalization - ANSWERS

### (a) Data Quality Issues:

1. **Multi-valued column:** DeviceOS has "iOS,Android"
2. **Inconsistent EventData:** JSON object vs plain string
3. **NULL UserID:** Row 3 has NULL with email present
4. **Invalid timestamp:** "Invalid Timestamp" in Row 3
5. **Empty JSON:** EventData "{}" in Row 3
6. **Redundant data:** UserEmail repeated with UserID

### (b) Anomalies:

1. **Update Anomaly:** Changing user email requires multiple row updates
2. **Insertion Anomaly:** Can't add user without event
3. **Deletion Anomaly:** Deleting events loses user info
4. **Data Inconsistency:** DeviceOS tied to session but device might change

### (c) Star Schema Design:

```
1 DIMENSION TABLES:
2 - dim_users(UserID PK, UserEmail, RegistrationDate)
3 - dim_devices(DeviceID PK, DeviceType, DeviceOS)
4 - dim_events(EventTypeID PK, EventTypename)
5 - dim_time(TimeID PK, Date, Hour, DayOfWeek)
6
7 FACT TABLE:
8 - fact_events(LogID PK, UserID FK, DeviceID FK, EventTypeID FK,
9               TimeID FK, SessionID, EventData, Timestamp)
```

### (d) SQL Queries:

```
1 -- Find users with multiple devices
2 SELECT u.UserID, u.UserEmail, COUNT(DISTINCT d.DeviceID) as device_count
3 FROM fact_events f
4 JOIN dim_users u ON f.UserID = u.UserID
5 JOIN dim_devices d ON f.DeviceID = d.DeviceID
6 GROUP BY u.UserID, u.UserEmail
7 HAVING COUNT(DISTINCT d.DeviceID) > 1;
8
9 -- Find users with multiple sessions
10 SELECT UserID, COUNT(DISTINCT SessionID) as session_count
11 FROM fact_events
12 WHERE UserID IS NOT NULL
13 GROUP BY UserID
14 HAVING COUNT(DISTINCT SessionID) > 1;
```

## 7.3 Section C: Python Script Debugging - ANSWERS

### (a) Errors:

Line	Error	Explanation
15	Missing colon	'for tokens in tokenized' needs ':'
16	Missing axis in mean	'np.mean()' on list of arrays needs 'axis=0'
16	KeyError possible	Token might not exist in vocabulary
26	Missing parentheses	'feature_importances' is attribute, but needs ' _ '
39	Missing comma	Between 'X' and 'y'
42	Square brackets	'predict[]' should be 'predict()'
43	Wrong brackets	'(predictions)' should be '{predictions}'



---

(b) Corrected Script:

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.ensemble import RandomForestClassifier
4 from gensim.models import Word2Vec
5
6 def create_embeddings(texts):
7     # Tokenize texts
8     tokenized = [text.split() for text in texts]
9
10    # Train Word2Vec
11    model = Word2Vec(tokenized, vector_size=100, window=5, min_count=1)
12
13    # Get document embeddings
14    embeddings = []
15    for tokens in tokenized:
16        # Filter tokens that exist in vocabulary
17        valid_tokens = [token for token in tokens if token in model.wv]
18        if valid_tokens:
19            doc_vec = np.mean([model.wv[token] for token in valid_tokens], axis
20                               =0)
21        else:
22            doc_vec = np.zeros(100) # Zero vector if no valid tokens
23        embeddings.append(doc_vec)
24
25    return np.array(embeddings)
26
27 def train_model(X, y):
28     clf = RandomForestClassifier(n_estimators=100)
29     clf.fit(X, y)
30
31     # Get feature importance
32     importance = clf.feature_importances_
33     print(f"Top features: {importance.argsort()[-5:]}")
34
35     return clf
36
37 if __name__ == '__main__':
38     df = pd.read_csv('data.csv')
39
40     # Create embeddings
41     X = create_embeddings(df['text'])
42     y = df['label']
43
44     # Train
45     model = train_model(X, y)
46
47     # Predict
48     predictions = model.predict(X[:10])
49     print(f"Predictions: {predictions}")
```

---

## 7.4 Section D: Embedding Demonstration - ANSWERS

### (a) Embedding Comparison:

Method	CBOW	Skip-gram	GloVe	BERT
Approach	Predict target from context	Predict context from target	Global co-occurrence matrix	Bidirectional
Training	Local context window	Local context window	Global statistics	Large context
Context	Fixed window	Fixed window	Document-level	Full sentence
Speed	Faster	Slower	Fast after matrix	Slow, iterative
Best for	Frequent words	Rare words	Analogies	Contextual meaning

### (b) Product Embedding Code:

```
1 from sentence_transformers import SentenceTransformer
2 import numpy as np
3
4 def generate_product_embeddings(descriptions):
5     # Load pre-trained model
6     model = SentenceTransformer('all-MiniLM-L6-v2')
7
8     # Generate embeddings
9     embeddings = model.encode(descriptions, show_progress_bar=True)
10
11     return embeddings
12
13 # Example usage
14 product_descriptions = [
15     "Wireless bluetooth headphones with noise cancellation",
16     "Premium over-ear headset for gaming",
17     "Running shoes with cushioned sole",
18     "Athletic sneakers for marathon training"
19 ]
20
21 embeddings = generate_product_embeddings(product_descriptions)
22 print(f"Embedding shape: {embeddings.shape}") # (4, 384)
```

### (c) Product Similarity Search:

```
1 from sklearn.metrics.pairwise import cosine_similarity
2 import faiss
3
4 # Build FAISS index for fast similarity search
5 def build_similarity_index(embeddings):
6     dimension = embeddings.shape[1]
7     index = faiss.IndexFlatIP(dimension) # Inner product (cosine after
8     # normalization)
9
10    # Normalize for cosine similarity
11    faiss.normalize_L2(embeddings)
12    index.add(embeddings)
13
14    return index
15
16 def find_similar_products(query_embedding, index, k=5):
```

```
16 faiss.normalize_L2(query_embedding.reshape(1, -1))
17 distances, indices = index.search(query_embedding.reshape(1, -1), k)
18 return indices[0], distances[0]
```

---

## 7.5 Section E: Model Selection - ANSWERS

### S1: Imbalanced Employee Attrition

- **Model:** XGBoost with class weights or SMOTE + Random Forest

- **Rules:**

1. Class imbalance (5%) requires balancing techniques
2. `scale_pos_weight` parameter handles imbalance in XGBoost
3. Use F1-score/AUC-ROC instead of accuracy for evaluation

### S2: Fake Review Detection

- **Model:** BERT fine-tuned for text classification

- **Rules:**

1. NLP task requires understanding semantic patterns
2. BERT captures contextual language nuances
3. Transfer learning reduces need for large labeled dataset

### S3: Intent Classification Chatbot

- **Model:** BERT/RoBERTa or Rasa NLU

- **Rules:**

1. Intent classification is multi-class text problem
2. Pre-trained language models understand diverse phrasings
3. Can handle out-of-vocabulary words through subword tokenization

### S4: Stock Direction Prediction

- **Model:** LSTM with technical indicators or Gradient Boosting

- **Rules:**

1. Time-series data with temporal dependencies
2. Binary classification (up/down) suits classification approach
3. Feature engineering crucial (moving averages, RSI, etc.)

### S5: Satellite Imagery Object Detection

- **Model:** YOLO or Faster R-CNN

- **Rules:**

1. Object detection requires localization + classification
2. YOLO provides real-time detection capability
3. Transfer learning from ImageNet reduces data requirements

### S6: Job-Candidate Matching

- **Model:** Siamese Network or Sentence-BERT

- **Rules:**

1. Similarity matching problem between two text sources
2. Siamese networks learn to compare document pairs
3. Can rank candidates by similarity score

## 8 PRACTICE PAPER SET 4

### 8.1 Section A: JSON Data Debugging (20 marks)

**Question 1:** The following JSON represents a dataset schema for a data science project. Find and correct all issues.

```
1 {
2   "dataset_info": {
3     "name": "Customer_Analytics",
4     "version": 1.2.0,
5     "created_at": new Date(),
6     "author": "DS Team"
7   },
8   "schema": {
9     "columns": [
10      {
11        "name": "customer_id",
12        "dtype": "int64",
13        "nullable": False,
14        "constraints": {unique: true, primary_key: true}
15      },
16      {
17        "name": "revenue",
18        "dtype": "float64",
19        "nullable": True,
20        "default": 0.0,
21        "validation": {
22          "min_value": 0,
23          "max_value": Infinity
24        }
25      },
26      {
27        "name": "segment"
28        "dtype": "category",
29        "categories": ["Bronze", "Silver", "Gold", "Platinum",]
30      }
31    ],
32  },
33  "transformations": [
34    {"step": 1, "operation": "fillna", "params": {"value": None}},
35    {"step": 2, "operation": "encode", "params": {method: "label"}}
36  ]
37 }
```

#### Tasks:

- List all JSON errors with their line references and impacts (10 marks)
- Provide the fully corrected JSON (6 marks)
- Write code to validate JSON against expected schema programmatically (4 marks)

## 8.2 Section B: SQL Debugging and Normalization (25 marks)

**Question 2:** Given the following e-commerce orders table:

OrderID	CustomerID	CustomerName	CustomerPhone	ProductID	ProductName	ProductCategory
1001	C1	Amit	9999999999	P1	Laptop	Electronics
1001	C1	Amit	9999999999	P2	Mouse	Electronics
1002	C1	Amit K	8888888888	P1	Laptop	Electronics

### Tasks:

- Identify ALL data quality and structural issues (7 marks)
- Normalize to 3NF with entity-relationship diagram description (8 marks)
- Write SQL to create normalized tables with proper constraints (5 marks)
- Write SQL query to find customers with multiple shipping addresses (5 marks)

## 8.3 Section C: Python Script Debugging (20 marks)

**Question 3:** Debug this data pipeline script:

```
1 1. import pandas as pd
2 2. import numpy as np
3 3. from sklearn.pipeline import Pipeline
4 4. from sklearn.impute import SimpleImputer
5 5. from sklearn.preprocessing import StandardScaler, OneHotEncoder
6 6. from sklearn.compose import ColumnTransformer
7 7.
8 8. def build_preprocessor(numeric_cols, categorical_cols):
9 9.     numeric_transformer = Pipeline([
10 10.         ('imputer', SimpleImputer(strategy='median'))
11 11.         ('scaler', StandardScaler())
12 12.     ])
13 13.
14 14.     categorical_transformer = Pipeline([
15 15.         ('imputer', SimpleImputer(strategy='constant' fill_value='missing'))
16 16.         ('encoder', OneHotEncoder(handle_unknown='ignore'))
17 17.     ])
18 18.
19 19.     preprocessor = ColumnTransformer([
20 20.         transformers(
21 21.             ('num', numeric_transformer, numeric_cols),
22 22.             ('cat', categorical_transformer, categorical_cols)
23 23.         )
24 24.     ])
25 25.
26 26.     return preprocessor
27 27.
28 28. def prepare_data(df, target_col):
29 29.     # Identify column types
30 30.     numeric_cols = df.select_dtypes(include=[np.number]).columns.tolist()
31 31.     categorical_cols = df.select_dtypes(include=['object']).columns.tolist()
()
```

```

32 32.
33 33.     # Remove target from features
34 34.     if target_col in numeric_cols
35 35.         numeric_cols.remove(target_col)
36 36.
37 37.     # Build and apply preprocessor
38 38.     preprocessor = build_preprocessor(numeric_cols, categorical_cols)
39 39.     X = df.drop(target_col, axis=1)
40 40.     y = df[target_col]
41 41.
42 42.     X_processed = preprocessor.fit_transform(X)
43 43.
44 44.     return X_processed, y, preprocessor
45 45.
46 46. if __name__ == "__main__"
47 47.     data = pd.read_csv("train.csv")
48 48.     X, y, prep = prepare_data(data, 'target')
49 49.     print(f"Processed shape: {X.shape}")

```

#### Tasks:

- Find all errors with line numbers and explanations (12 marks)
- Provide the corrected script (8 marks)

### 8.4 Section D: RAG Workflow Design (15 marks)

**Question 4:** Design a RAG system for Kelp Global's internal knowledge base that handles financial reports, meeting notes, and investor presentations.

#### Tasks:

- Design the complete RAG pipeline with component specifications (6 marks)
- How would you handle document versioning and updates in the RAG system? (4 marks)
- Explain retrieval strategies: sparse vs dense retrieval and when to use hybrid (5 marks)

### 8.5 Section E: Model Selection Scenarios (20 marks)

**Question 5:** Select the best approach for each scenario with 3 justified rules.

Scenario	Description
S1	Forecasting inventory demand for next 30 days across 1000 products
S2	Identifying which features most impact customer lifetime value
S3	Detecting unusual server behavior from system logs
S4	Translating internal documents from English to Hindi
S5	Automatically extracting key information from contracts
S6	Personalized content ranking on a news feed

## 9 PRACTICE PAPER SET 4 - ANSWERS

### 9.1 Section A: JSON Data Debugging - ANSWERS

#### (a) JSON Errors:

Line	Error	Impact
4	'1.2.0' not quoted	Version should be string "1.2.0"
5	'new Date()'	JavaScript code not valid in JSON
11	'False' Python keyword	Should be 'false' (lowercase)
12	Unquoted keys	{'unique': true, 'primary_key': true} needs quoted keys
16	'True' Python keyword	Should be 'true' (lowercase)
21	'Infinity'	Not valid JSON, use very large number or null
24	Missing comma	Between "name": "segment" and next line
25	Trailing comma	In categories array
27	Trailing comma	After columns array
29	'None' Python keyword	Should be 'null'
30	Unquoted key	'method': "label" needs "method"

#### (b) Corrected JSON:

```
1 {
2   "dataset_info": {
3     "name": "Customer_Analytics",
4     "version": "1.2.0",
5     "created_at": "2024-12-02T10:00:00Z",
6     "author": "DS Team"
7   },
8   "schema": {
9     "columns": [
10      {
11        "name": "customer_id",
12        "dtype": "int64",
13        "nullable": false,
14        "constraints": {"unique": true, "primary_key": true}
15      },
16      {
17        "name": "revenue",
18        "dtype": "float64",
19        "nullable": true,
20        "default": 0.0,
21        "validation": {
22          "min_value": 0,
23          "max_value": 999999999
24        }
25      },
26      {
27        "name": "segment",
28        "dtype": "category",
29        "categories": ["Bronze", "Silver", "Gold", "Platinum"]
30      }
31    ]
32  },
33  "transformations": [
34    {"step": 1, "operation": "fillna", "params": {"value": null}},
```

```

35     {"step": 2, "operation": "encode", "params": {"method": "label"}}
36 ]
37 }

```

### (c) Schema Validation Code:

```

1 from jsonschema import validate, Draft7Validator
2
3 schema = {
4     "$schema": "http://json-schema.org/draft-07/schema#",
5     "type": "object",
6     "required": ["dataset_info", "schema"],
7     "properties": {
8         "dataset_info": {
9             "type": "object",
10            "required": ["name", "version"],
11            "properties": {
12                "name": {"type": "string", "pattern": "^[A-Za-z_]+$"},
13                "version": {"type": "string", "pattern": "^\\d+\\.\\d+\\.\\d+$"}
14            }
15        },
16        "schema": {
17            "type": "object",
18            "properties": {
19                "columns": {
20                    "type": "array",
21                    "items": {
22                        "type": "object",
23                        "required": ["name", "dtype"],
24                        "properties": {
25                            "name": {"type": "string"},
26                            "dtype": {"type": "string", "enum": ["int64", "float64", "category", "string"]},
27                            "nullable": {"type": "boolean"}
28                        }
29                    }
30                }
31            }
32        }
33    }
34 }
35
36 validator = Draft7Validator(schema)
37 errors = list(validator.iter_errors(config_data))
38 for error in errors:
39     print(f"Error: {error.message} at {error.path}")

```

## 9.2 Section B: SQL Normalization - ANSWERS

### (a) Data Quality Issues:

1. **Calculated column mismatch:** TotalPrice 99999  $2 \times 50000$  (should be 100000)
2. **Multi-valued column:** PaymentMethod has "COD, Card"
3. **Inconsistent customer names:** "Amit" vs "Amit K" for same customer
4. **Different phone numbers:** Same customer C1 has different phones
5. **Redundant data:** Customer info repeated in each order



6. **Same OrderID:** Orders 1001 appear twice (multi-line order vs duplication)
7. **No order line number:** Cannot distinguish products in same order

**(b) Normalized Schema (3NF):**

```

1 CUSTOMERS (CustomerID PK, CustomerName, CustomerPhone)
2   |
3   |-- 1:N
4   v
5 ADDRESSES (AddressID PK, CustomerID FK, Address, AddressType)
6   |
7   |-- N:1
8   v
9 ORDERS (OrderID PK, CustomerID FK, OrderDate, ShippingAddressID FK)
10  |
11  |-- 1:N
12  v
13 ORDER_ITEMS (OrderID FK, ProductID FK, Quantity, UnitPrice, LineTotal)
14               [Composite PK: OrderID, ProductID]
15  |
16  |-- N:1
17  v
18 PRODUCTS (ProductID PK, ProductName, CategoryID FK)
19  |
20  |-- N:1
21  v
22 CATEGORIES (CategoryID PK, CategoryName)
23
24 PAYMENTS (PaymentID PK, OrderID FK, PaymentMethod, Amount)

```

**(c) SQL Creation:**

```

1 CREATE TABLE Customers (
2     CustomerID VARCHAR(10) PRIMARY KEY,
3     CustomerName VARCHAR(100) NOT NULL,
4     CustomerPhone VARCHAR(15)
5 );
6
7 CREATE TABLE Addresses (
8     AddressID INT PRIMARY KEY AUTO_INCREMENT,
9     CustomerID VARCHAR(10) NOT NULL,
10    Address VARCHAR(200) NOT NULL,
11    AddressType VARCHAR(20) DEFAULT 'Shipping',
12    FOREIGN KEY (CustomerID) REFERENCES Customers (CustomerID)
13 );
14
15 CREATE TABLE Categories (
16     CategoryID INT PRIMARY KEY AUTO_INCREMENT,
17     CategoryName VARCHAR(50) UNIQUE NOT NULL
18 );
19
20 CREATE TABLE Products (
21     ProductID VARCHAR(10) PRIMARY KEY,
22     ProductName VARCHAR(100) NOT NULL,
23     CategoryID INT,
24     UnitPrice DECIMAL(10,2) NOT NULL,
25     FOREIGN KEY (CategoryID) REFERENCES Categories (CategoryID)
26 );

```

```

27
28 CREATE TABLE Orders (
29     OrderID INT PRIMARY KEY,
30     CustomerID VARCHAR(10) NOT NULL,
31     OrderDate DATE NOT NULL,
32     ShippingAddressID INT,
33     FOREIGN KEY (CustomerID) REFERENCES Customers(CustomerID),
34     FOREIGN KEY (ShippingAddressID) REFERENCES Addresses(AddressID)
35 );
36
37 CREATE TABLE Order_Items (
38     OrderID INT,
39     ProductID VARCHAR(10),
40     Quantity INT NOT NULL CHECK (Quantity > 0),
41     UnitPrice DECIMAL(10,2) NOT NULL,
42     LineTotal DECIMAL(12,2) GENERATED ALWAYS AS (Quantity * UnitPrice),
43     PRIMARY KEY (OrderID, ProductID),
44     FOREIGN KEY (OrderID) REFERENCES Orders(OrderID),
45     FOREIGN KEY (ProductID) REFERENCES Products(ProductID)
46 );

```

#### (d) Query for Multiple Addresses:

```

1 SELECT c.CustomerID, c.CustomerName, COUNT(DISTINCT a.Address) as address_count
2 FROM Customers c
3 JOIN Addresses a ON c.CustomerID = a.CustomerID
4 GROUP BY c.CustomerID, c.CustomerName
5 HAVING COUNT(DISTINCT a.Address) > 1;

```

## 9.3 Section C: Python Script Debugging - ANSWERS

#### (a) Errors:

Line	Error	Explanation
10-11	Missing comma	Between imputer and scaler in Pipeline list
15	Missing comma	Between 'constant' and 'fill_value'
19-24	Wrong syntax	'transformers()' is not valid; should be 'transformers=[]'
30	Wrong method	'toList()' should be 'tolist()' (lowercase)
34	Missing colon	'if' statement needs ':'
42	Missing parenthesis	Unclosed 'fit_transform(X'
46	Missing colon	'if __name__' needs ':'

#### (b) Corrected Script:

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.pipeline import Pipeline
4 from sklearn.impute import SimpleImputer
5 from sklearn.preprocessing import StandardScaler, OneHotEncoder

```

```

6 from sklearn.compose import ColumnTransformer
7
8 def build_preprocessor(numeric_cols, categorical_cols):
9     numeric_transformer = Pipeline([
10         ('imputer', SimpleImputer(strategy='median')),
11         ('scaler', StandardScaler())
12     ])
13
14     categorical_transformer = Pipeline([
15         ('imputer', SimpleImputer(strategy='constant', fill_value='missing')),
16         ('encoder', OneHotEncoder(handle_unknown='ignore'))
17     ])
18
19     preprocessor = ColumnTransformer(
20         transformers=[
21             ('num', numeric_transformer, numeric_cols),
22             ('cat', categorical_transformer, categorical_cols)
23         ]
24     )
25
26     return preprocessor
27
28 def prepare_data(df, target_col):
29     # Identify column types
30     numeric_cols = df.select_dtypes(include=[np.number]).columns.tolist()
31     categorical_cols = df.select_dtypes(include=['object']).columns.tolist()
32
33     # Remove target from features
34     if target_col in numeric_cols:
35         numeric_cols.remove(target_col)
36
37     # Build and apply preprocessor
38     preprocessor = build_preprocessor(numeric_cols, categorical_cols)
39     X = df.drop(target_col, axis=1)
40     y = df[target_col]
41
42     X_processed = preprocessor.fit_transform(X)
43
44     return X_processed, y, preprocessor
45
46 if __name__ == "__main__":
47     data = pd.read_csv("train.csv")
48     X, y, prep = prepare_data(data, 'target')
49     print(f"Processed shape: {X.shape}")

```

## 9.4 Section D: RAG Workflow - ANSWERS

### (a) RAG Pipeline Design:

```

1 Component Specifications:
2
3 1. Document Ingestion Layer:
4   - PDF Parser: PyPDF2/PDFPlumber
5   - DOCX Parser: python-docx
6   - PPT Parser: python-pptx
7   - Metadata extraction for versioning
8

```

```

9 2. Chunking Strategy:
10   - Chunk size: 512 tokens with 50 overlap
11   - Respect document sections/headers
12   - Keep tables intact
13
14 3. Embedding Layer:
15   - Model: text-embedding-3-small (OpenAI) or BGE
16   - Dimension: 1536/768
17
18 4. Vector Store:
19   - Primary: Pinecone/Chroma
20   - Metadata: document_id, version, date, type
21
22 5. Retrieval:
23   - Top-K: 5-10 chunks
24   - Hybrid: BM25 + Dense retrieval
25
26 6. LLM:
27   - GPT-4 or LLaMA-2-70B
28   - Context window: 8K tokens

```

## (b) Document Versioning:

1. **Version metadata:** Store version\_id, timestamp with each chunk
2. **Incremental updates:** Re-embed only changed documents
3. **Document lineage:** Maintain parent-child relationships
4. **Time-based retrieval:** Allow queries specific to document versions
5. **Soft delete:** Mark old versions inactive rather than deleting

## (c) Retrieval Strategies:

Strategy	How it Works	When to Use
<b>**Sparse (BM25)**</b>	Keyword matching based on TF-IDF	Exact term matches, specific terminology
<b>**Dense**</b>	Semantic similarity via embeddings	Conceptual queries, paraphrased questions
<b>**Hybrid**</b>	Combines both with weighted scoring	Production systems needing both precision and recall

Use hybrid when: queries mix specific terms with conceptual questions (common in enterprise Q&A).

## 9.5 Section E: Model Selection - ANSWERS

### S1: Multi-Product Inventory Forecasting

- **Model:** LightGBM or Global Time-Series Model (like N-BEATS)
- **Rules:**
  1. 1000 products need scalable approach
  2. LightGBM handles categorical features (product types) well
  3. Can leverage cross-product patterns for new products

### S2: Feature Importance for CLV

- **Model:** SHAP with XGBoost or Permutation Importance

- **Rules:**

1. Need interpretable feature rankings
2. SHAP provides local and global explanations
3. Business stakeholders need actionable insights

**S3: Server Log Anomaly Detection**

- **Model:** Isolation Forest or LSTM Autoencoder

- **Rules:**

1. Unsupervised - normal behavior not fully defined
2. Isolation Forest works with high-dimensional log features
3. LSTM captures temporal patterns in log sequences

**S4: English to Hindi Translation**

- **Model:** Fine-tuned mBART or IndicTrans2

- **Rules:**

1. Neural machine translation for quality
2. IndicTrans2 specifically trained for Indian languages
3. Can handle domain-specific terminology with fine-tuning

**S5: Contract Information Extraction**

- **Model:** Named Entity Recognition (NER) + Document AI

- **Rules:**

1. Structured extraction from unstructured text
2. Pre-trained legal NER models available (spaCy, Hugging Face)
3. Layout-aware models handle document structure

**S6: Personalized News Ranking**

- **Model:** Learning to Rank (LambdaMART) or Two-Tower Model

- **Rules:**

1. User-item interaction requires collaborative approach
2. Two-tower encodes user preferences and content separately
3. Real-time scoring with pre-computed embeddings

## 10 PRACTICE PAPER SET 5

### 10.1 Section A: JSON Data Debugging (20 marks)

**Question 1:** Analyze this RAG system configuration JSON and fix all issues:

```
1 {  
2   "rag_config": {  
3     "embedding_model": {  
4       "name": "text-embedding-ada-002",  
5       "dimension": 1536,  
6       "max_tokens": 8191,  
7       "batch_size": 100.0  
8     },
```

```

9      "vector_store": {
10          "provider": 'pinecone',
11          "index_name": "kelp-knowledge-base",
12          "metric": "cosine",
13          "replicas": 1
14          "pods": 1
15      },
16      "chunking": {
17          "strategy": "recursive",
18          "chunk_size": 512,
19          "chunk_overlap": 50,
20          "separators": ["\\n\\n", "\\n", " "],
21      },
22      "retrieval": {
23          "top_k": 5,
24          "score_threshold": 0.75,
25          "reranker": {
26              "enabled": TRUE,
27              "model": "cross-encoder/ms-marco-MiniLM-L-12-v2"
28          },
29          "filters": {
30              "doc_type": ["pdf", "docx"],
31              "max_age_days": Null
32          }
33      },
34      "llm_config": {
35          "model": "gpt-4-turbo",
36          "temperature": .7,
37          "max_tokens": 2000,
38          "system_prompt": "You are a helpful assistant for Kelp Global...",
39          "response_format": {"type": json}
40      }
41  },
42  "metadata": {
43      "version": "1.0",
44      "created_by": "DS Team",
45      "last_updated": 2024-12-01
46  }
47 }

```

### Tasks:

- Find all structural and syntactical errors with explanations (10 marks)
- Provide corrected JSON (6 marks)
- How would you test this configuration before deployment? (4 marks)

## 10.2 Section B: SQL Debugging and Normalization (25 marks)

**Question 2:** Given a data science project tracking table:

ProjectID	ProjectName	TeamMembers	DataSources	ModelTypes	Metrics	Status
DS001	Churn_Model	Priya,Amit,Raj	CRM,Analytics	LogReg,XGB,RF	AUC:0.85,F1:0.78	Active
DS001	Churn_Model	Priya,Amit	CRM,Analytics	LogReg,XGB	AUC:0.85	Completed
DS002	Sales_Forecast	NULL	Sales DB	Prophet	MAPE:12%	Active

### Tasks:

- (a) Identify data quality issues, duplicates, and anomalies (6 marks)
- (b) Why does this violate normal forms? (5 marks)
- (c) Design normalized schema with proper relationships (7 marks)
- (d) Write queries for: finding projects with no team, all active projects with their team count (7 marks)

## 10.3 Section C: Python Script Debugging (20 marks)

**Question 3:** Fix the following RAG implementation script:

```
1 1. from langchain.embeddings import OpenAIEmbeddings
2 2. from langchain.vectorstores import Chroma
3 3. from langchain.chains import RetrievalQA
4 4. from langchain.llms import ChatOpenAI
5 5. from langchain.text_splitter import RecursiveCharacterTextSplitter
6 6. import os
7 7.
8 8. class RAGPipeline
9 9.     def __init__(self, api_key):
10 10.         os.environ['OPENAI_API_KEY'] = api_key
11 11.         self.embeddings = OpenAIEmbeddings[]
12 12.         self.vector_store = None
13 13.         self.qa_chain = none
14 14.
15 15.     def ingest_documents(self, documents):
16 16.         # Split documents
17 17.         splitter = RecursiveCharacterTextSplitter(
18 18.             chunk_size=500,
19 19.             chunk_overlap=50
20 20.             separators=["\\n\\n", "\\n", " "]
21 21.         )
22 22.         chunks = splitter.split_documents(documents)
23 23.
24 24.         # Create vector store
25 25.         self.vector_store = Chroma.from_documents(
26 26.             documents=chunks
27 27.             embedding=self.embeddings,
28 28.             persist_directory="./chroma_db"
29 29.         )
30 30.
31 31.     def setup_qa_chain(self, model_name='gpt-3.5-turbo'):
32 32.         llm = ChatOpenAI(model_name=model_name, temperature=0)
33 33.
34 34.         self.qa_chain = RetrievalQA.from_chain_type(
35 35.             llm=llm,
36 36.             chain_type="stuff"
37 37.             retriever=self.vector_store.as_retriever(search_kwargs={"k":
5})
38 38.         )
39 39.
40 40.     def query(self, question)
41 41.         if self.qa_chain == None:
```

```

42 42.         raise ValueError("QA chain not initialized")
43 43.
44 44.         response = self.qa_chain.run[question]
45 45.         return response
46 46.
47 47. if __name__ == "__main__":
48 48.     rag = RAGPipeline("sk-xxx")
49 49.     # Load and ingest documents
50 50.     from langchain.document_loaders import DirectoryLoader
51 51.     loader = DirectoryLoader("./docs", glob="**/*.pdf")
52 52.     docs = loader.load()
53 53.     rag.ingest_documents(docs)
54 54.     rag.setup_qa_chain
55 55.     answer = rag.query("What is our revenue target?")
56 56.     print(answer)

```

#### Tasks:

- Identify all errors with line numbers and reasons (12 marks)
- Provide the corrected and improved script (8 marks)

### 10.4 Section D: Complete RAG + Embedding Task (15 marks)

**Question 4:** Design an end-to-end RAG system for private equity document analysis at Kelp Global.

#### Tasks:

- Design the complete system architecture diagram (5 marks)
- Explain embedding choices: why would you choose one embedding model over another? (5 marks)
- Write pseudocode/code for the retrieval and response generation pipeline (5 marks)

### 10.5 Section E: Model Selection Scenarios (20 marks)

**Question 5:** Comprehensive scenario-based model selection.

Scenario	Description
S1	Predicting customer segment (Premium, Standard, Basic) based on behavior
S2	Real-time recommendation of next best action for sales team
S3	Extracting and linking entities from financial news articles
S4	Predicting probability of deal closure for private equity investments
S5	Generating automated summaries of lengthy legal documents
S6	Building a semantic search system for internal company wiki

For each scenario: Model selection + 3 Rules + Key preprocessing steps



## 11 PRACTICE PAPER SET 5 - ANSWERS

### 11.1 Section A: JSON Data Debugging - ANSWERS

#### (a) JSON Errors:

Line	Error	Explanation
6	100.0	Should be integer 100 for batch_size
10	Single quotes	'pinecone' should use double quotes
13	Missing comma	After "replicas": 1
20	Trailing comma	After separators array
26	TRUE	JavaScript convention; should be true
31	Null	Wrong case; should be null
34	Leading zero missing	.7 should be 0.7
37	Unquoted value	{"type": json} should be {"type": "json"}
42	Unquoted date	2024-12-01 should be "2024-12-01"

#### (b) Corrected JSON:

```
1 {
2   "rag_config": {
3     "embedding_model": {
4       "name": "text-embedding-ada-002",
5       "dimension": 1536,
6       "max_tokens": 8191,
7       "batch_size": 100
8     },
9     "vector_store": {
10      "provider": "pinecone",
11      "index_name": "kelp-knowledge-base",
12      "metric": "cosine",
13      "replicas": 1,
14      "pods": 1
15    },
16    "chunking": {
17      "strategy": "recursive",
18      "chunk_size": 512,
19      "chunk_overlap": 50,
20      "separators": ["\\n\\n", "\\n", " "]
21    },
22    "retrieval": {
23      "top_k": 5,
24      "score_threshold": 0.75,
25      "reranker": {
26        "enabled": true,
27        "model": "cross-encoder/ms-marco-MiniLM-L-12-v2"
28      },
29      "filters": {
30        "doc_type": ["pdf", "docx"],
31        "max_age_days": null
32      }
33    },
34    "llm_config": {
35      "model": "gpt-4-turbo",
36      "temperature": 0.7,
```

```

37         "max_tokens": 2000,
38         "system_prompt": "You are a helpful assistant for Kelp Global...",
39         "response_format": {"type": "json"}
40     },
41 },
42 "metadata": {
43     "version": "1.0",
44     "created_by": "DS Team",
45     "last_updated": "2024-12-01"
46 }
47 }

```

### (c) Testing Strategy:

```

1 import json
2 from jsonschema import validate
3
4 # 1. Syntax validation
5 def test_json_syntax(config_path):
6     try:
7         with open(config_path) as f:
8             config = json.load(f)
9         return True, config
10    except json.JSONDecodeError as e:
11        return False, str(e)
12
13 # 2. Schema validation
14 def test_schema(config, schema):
15     validate(instance=config, schema=schema)
16
17 # 3. Configuration testing
18 def test_config_values(config):
19     assert config['rag_config']['chunking']['chunk_size'] > config['rag_config']
20     assert config['rag_config']['retrieval']['top_k'] > 0
21     assert 0 <= config['rag_config']['llm_config']['temperature'] <= 1
22
23 # 4. Integration testing - test with sample query before production

```

## 11.2 Section B: SQL Normalization - ANSWERS

### (a) Data Quality Issues:

1. **Duplicate project:** DS001 appears twice with different data
2. **Multi-valued columns:** TeamMembers, DataSources, ModelTypes, Metrics
3. **Inconsistent naming:** "Churn\_Model" vs "Churn Model"
4. **NULL team:** DS002 has no team members
5. **Invalid date:** "TBD" is not a proper date value
6. **Negative budget:** -50000 is invalid
7. **Inconsistent status:** Same project has Active and Completed

### (b) Normal Form Violations:

- **1NF:** Multi-valued attributes (TeamMembers, DataSources, etc.)

- **2NF**: Metrics depend on ModelType, not just ProjectID
- **3NF**: Team member details might have transitive dependencies

### (c) Normalized Schema:

```

1 PROJECTS (ProjectID PK, ProjectName, Status, StartDate, EndDate, Budget)
2     - CHECK constraint: Budget > 0
3     - CHECK constraint: Status IN ('Active', 'Completed', 'On Hold')
4
5 TEAM_MEMBERS (MemberID PK, MemberName, Role, Email)
6
7 PROJECT_TEAM (ProjectID FK, MemberID FK) - Composite PK
8
9 DATA_SOURCES (SourceID PK, SourceName, SourceType)
10
11 PROJECT_DATASOURCES (ProjectID FK, SourceID FK)
12
13 MODELS (ModelID PK, ProjectID FK, ModelType, ModelVersion)
14
15 MODEL_METRICS (MetricID PK, ModelID FK, MetricName, MetricValue, EvaluationDate
16 )

```

### (d) SQL Queries:

```

1 -- Projects with no team members
2 SELECT p.ProjectID, p.ProjectName
3 FROM Projects p
4 LEFT JOIN Project_Team pt ON p.ProjectID = pt.ProjectID
5 WHERE pt.MemberID IS NULL;
6
7 -- Active projects with team count
8 SELECT
9     p.ProjectID,
10    p.ProjectName,
11    COUNT(pt.MemberID) as team_count
12 FROM Projects p
13 LEFT JOIN Project_Team pt ON p.ProjectID = pt.ProjectID
14 WHERE p.Status = 'Active'
15 GROUP BY p.ProjectID, p.ProjectName
16 ORDER BY team_count DESC;

```

## 11.3 Section C: Python Script Debugging - ANSWERS

### (a) Errors:

Line	Error	Explanation
4	Wrong import	Should be from langchain.chat_models import ChatOpenAI
8	Missing colon	Class definition needs :
11	Square brackets	OpenAIEmbeddings[] should be OpenAIEmbeddings()
13	Wrong keyword	none should be None

Line	Error	Explanation
19-20	Missing comma	Between chunk_overlap and separators
26	Missing comma	Between chunks and embedding
36	Missing comma	After "stuff"
40	Missing colon	Function definition needs :
41	Best practice	Use is None instead of == None
44	Square brackets	run[] should be run()
54	Missing parentheses	setup_qa_chain needs () to call

## (b) Corrected Script:

```

1 from langchain.embeddings import OpenAIEmbeddings
2 from langchain.vectorstores import Chroma
3 from langchain.chains import RetrievalQA
4 from langchain.chat_models import ChatOpenAI
5 from langchain.text_splitter import RecursiveCharacterTextSplitter
6 import os
7
8 class RAGPipeline:
9     def __init__(self, api_key):
10         os.environ['OPENAI_API_KEY'] = api_key
11         self.embeddings = OpenAIEmbeddings()
12         self.vector_store = None
13         self.qa_chain = None
14
15     def ingest_documents(self, documents):
16         splitter = RecursiveCharacterTextSplitter(
17             chunk_size=500,
18             chunk_overlap=50,
19             separators=["\n\n", "\n", " "]
20         )
21         chunks = splitter.split_documents(documents)
22
23         self.vector_store = Chroma.from_documents(
24             documents=chunks,
25             embedding=self.embeddings,
26             persist_directory="./chroma_db"
27         )
28         self.vector_store.persist()
29
30     def setup_qa_chain(self, model_name='gpt-3.5-turbo'):
31         llm = ChatOpenAI(model_name=model_name, temperature=0)
32
33         self.qa_chain = RetrievalQA.from_chain_type(
34             llm=llm,
35             chain_type="stuff",
36             retriever=self.vector_store.as_retriever(search_kwargs={"k": 5})
37         )
38
39     def query(self, question):
40         if self.qa_chain is None:
41             raise ValueError("QA chain not initialized")
42
43         response = self.qa_chain.run(question)
44         return response
45
46 if __name__ == "__main__":
47     rag = RAGPipeline("sk-xxx")

```

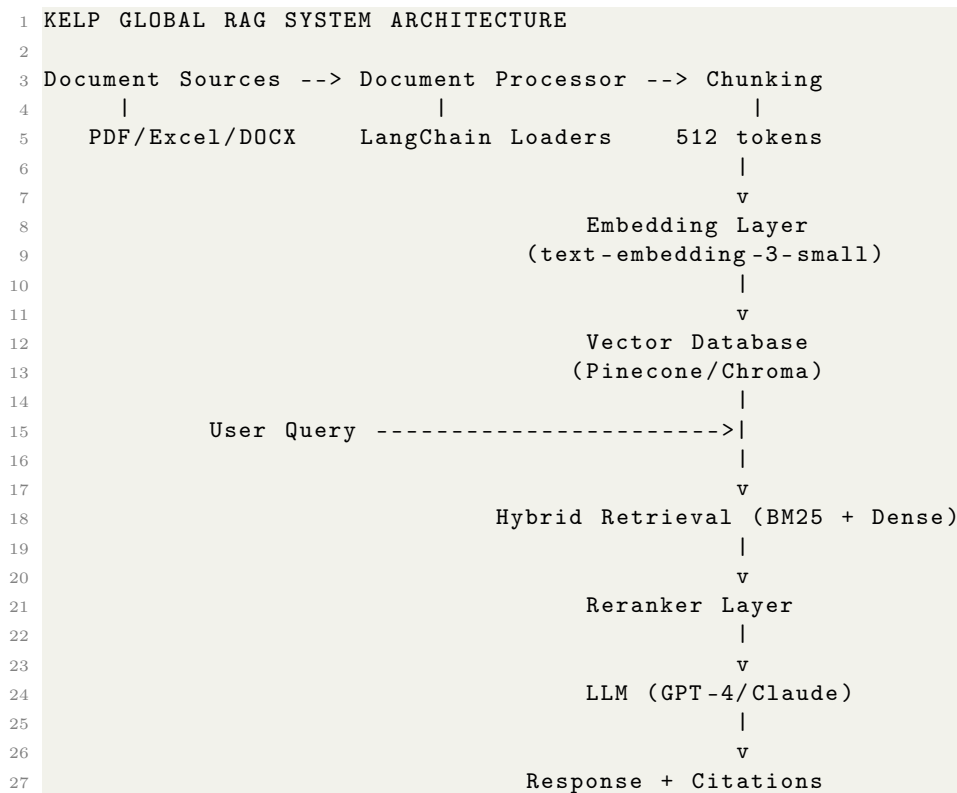
```

48 from langchain.document_loaders import DirectoryLoader
49 loader = DirectoryLoader("./docs", glob="**/*.pdf")
50 docs = loader.load()
51 rag.ingest_documents(docs)
52 rag.setup_qa_chain()
53 answer = rag.query("What is our revenue target?")
54 print(answer)

```

## 11.4 Section D: RAG + Embedding System - ANSWERS

### (a) System Architecture:



### (b) Embedding Model Selection:

Model	Dimension	Use Case	Trade-off
text-embedding-ada-002	1536	General purpose	Higher cost, larger storage
text-embedding-3-small	1536	Cost-effective	Good balance
BGE-base	768	Open source	No API costs, needs hosting
all-MiniLM-L6-v2	384	Fast, lightweight	Lower quality, faster

#### Selection Criteria:

1. Domain specificity - Financial domain may benefit from fine-tuned embeddings
2. Latency requirements - Smaller models for real-time applications

3. Cost constraints - Open-source vs commercial
4. Data privacy - Self-hosted for sensitive documents

### (c) Pipeline Pseudocode:

```

1 def rag_query_pipeline(user_query, config):
2     # 1. Embed the query
3     query_embedding = embedding_model.encode(user_query)
4
5     # 2. Retrieve relevant chunks
6     sparse_results = bm25_index.search(user_query, top_k=20)
7     dense_results = vector_db.similarity_search(query_embedding, top_k=20)
8
9     # 3. Hybrid fusion
10    combined_results = reciprocal_rank_fusion(sparse_results, dense_results)
11
12    # 4. Rerank
13    reranked = reranker_model.rerank(query=user_query, documents=
combined_results[:20])
14    top_contexts = reranked[:5]
15
16    # 5. Build prompt with context
17    context_text = "\n---\n".join([doc.page_content for doc in top_contexts])
18
19    # 6. Generate response
20    response = llm.generate(prompt_with_context, temperature=0.3)
21
22    # 7. Add citations
23    return add_source_citations(response, top_contexts)

```

---

## 11.5 Section E: Model Selection - ANSWERS

### S1: Customer Segment Classification

- **Model:** XGBoost / LightGBM Multi-class Classifier
- **Rules:** Multi-class classification; Tree-based models handle mixed features; Feature importance explains segments
- **Preprocessing:** Handle class imbalance, encode categoricals, scale numerics

### S2: Real-time Next Best Action

- **Model:** Multi-Armed Bandit or Reinforcement Learning
- **Rules:** Balance exploration vs exploitation; Real-time with continuous learning; Thompson Sampling handles uncertainty
- **Preprocessing:** Feature engineering for user context, action encoding

### S3: Financial Entity Extraction

- **Model:** SpaCy NER + Entity Linking (fine-tuned on financial data)
- **Rules:** NER for extraction; Entity linking to knowledge base; Pre-trained financial NER available
- **Preprocessing:** Text normalization, sentence segmentation

### S4: Deal Closure Probability

- **Model:** Logistic Regression or Gradient Boosting with calibration

- **Rules:** Binary classification with probability; Calibrated probabilities for decisions; Feature importance for insights
- **Preprocessing:** Handle missing data, engineer timeline features

#### **S5: Legal Document Summarization**

- **Model:** Fine-tuned T5 or LongT5 / GPT-4 with RAG
- **Rules:** Long document handling; Abstractive summarization; Legal terminology preservation
- **Preprocessing:** Document chunking, section identification

#### **S6: Semantic Search for Wiki**

- **Model:** Dense Retrieval (Sentence-BERT + FAISS)
- **Rules:** Semantic similarity over keywords; Fast retrieval with ANN; Easy to update
- **Preprocessing:** Document chunking, metadata extraction