

**UNIVERSITY CHATBOT SYSTEM*****Using Custom TF-IDF Implementation Without AI Libraries***

<i>Submitted By</i>	Zara Imran
<i>SAP ID</i>	70149236
<i>Section</i>	BSIET-5I
<i>Submitted To</i>	Sir Hamedoon
<i>Course</i>	Artificial Intelligence
<i>Date</i>	January 2026

# PROJECT REPORT

---

## 1. PROJECT DESCRIPTION

### 1.1 Introduction

The *University Chatbot System* is an intelligent web-based application designed to automatically answer common university-related questions. It helps students, parents, and visitors obtain information about admissions, courses, fees, facilities, and other services without requiring manual staff support.

This project is implemented entirely from scratch without using any AI or machine learning libraries such as scikit-learn, TensorFlow, or spaCy. Instead, the system manually implements the **TF-IDF (Term Frequency–Inverse Document Frequency)** algorithm and **cosine similarity** to match user queries with the most relevant answers.

---

### 1.2 Problem Statement

Universities receive hundreds of repetitive inquiries daily such as:

- How to apply for admission?
- What courses are available?
- What are the fees?
- Where is the library?
- When does the semester start?

Handling these queries manually is time-consuming and resource-intensive. Students often experience delays and frustration. An automated chatbot can significantly reduce this workload while providing instant responses.

---

### 1.3 Project Objectives

The key objectives of this project are:

1. To develop a functional chatbot without using any AI libraries.
2. To manually implement TF-IDF for text vectorization.
3. To compute cosine similarity using mathematical formulas.
4. To provide accurate responses based on similarity matching.

5. To create a user-friendly web interface.
  6. To maintain conversation logs and performance statistics.
  7. To allow easy expansion of the knowledge base.
- 

## 1.4 Key Features

### Technical Features:

- Custom-built TF-IDF vectorizer
- Manual cosine similarity computation
- Text preprocessing (tokenization, stop-word removal)
- Flask-based RESTful API
- JSON-based knowledge base
- Real-time responses
- Conversation history tracking

### User Features:

- Simple chat interface
  - Instant automated replies
  - Similarity score for confidence
  - Topic browsing
  - Conversation history viewing
  - Web-based accessibility
- 

## 1.5 Technologies Used

- **Programming Language:** Python 3.8+
  - **Framework:** Flask, Flask-CORS
  - **Frontend:** HTML, CSS, JavaScript
  - **Data Storage:** JSON
  - **Libraries:** Only Flask-related (no AI/ML libraries)
-

## 1.6 Importance of the Project

### Educational Value:

- Demonstrates fundamental NLP concepts
- Explains TF-IDF and vector mathematics
- Builds algorithmic problem-solving skills

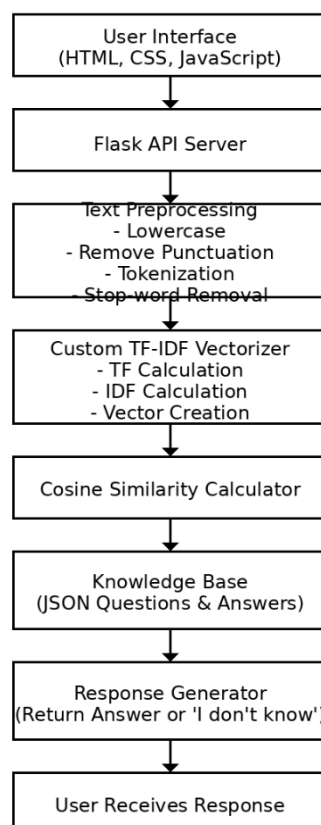
### Practical Applications:

- Reduces staff workload
- Provides 24/7 student support
- Can be adapted for hospitals, government offices, etc.

---

## 2. PROJECT FLOW DIAGRAM

**University Chatbot System Flow Diagram**

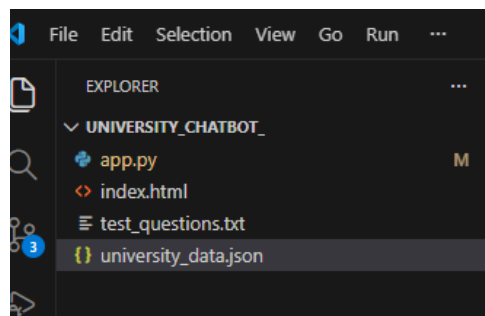


## 2.2 Step-by-Step Workflow

1. User enters a question in the chat interface.
  2. JavaScript sends the request to Flask API.
  3. Server receives input via `/api/chat`.
  4. Text preprocessing is applied.
  5. Query is converted into TF-IDF vector.
  6. Similarity is calculated against all stored questions.
  7. Best match is selected based on highest score.
  8. If score > threshold, corresponding answer is returned.
  9. Conversation is logged with timestamp and similarity score.
  10. Response is sent back to frontend and displayed.
- 

## 3. PROJECT CODE (SCREENSHOT REFERENCES)

### 3.1 File Structure



### 3.2 Major Components

1. **CustomVectorizer Class**
    - Handles tokenization, TF, IDF, and vector creation.
  2. **UniversityChatbot Class**
    - Controls chatbot logic, training, and matching.
  3. **Flask API**
    - Manages HTTP requests and responses.
-

### 3.3 Core Algorithm Summary

**TF (Term Frequency):**

$$TF = (\text{Word count in document}) / (\text{Total words})$$

**IDF (Inverse Document Frequency):**

$$IDF = \log(\text{Total documents} / \text{Documents containing the word})$$

**TF-IDF:**

$$TF-IDF = TF \times IDF$$

**Cosine Similarity:**

$$\text{Similarity} = (A \cdot B) / (|A| \times |B|)$$

---

### 3.4 Code Screenshots

#### ➤ CustomVectorizer Class (TF-IDF Implementation)

```

12
13 class CustomVectorizer:
14     """Custom TF-IDF implementation without using sklearn"""
15
16     def __init__(self):
17         self.vocabulary = {} # word -> index mapping
18         self.idf_values = {} # word -> IDF score
19         self.documents = [] # store original documents
20
21     def _tokenize(self, text):
22         """Convert text to lowercase and split into words"""
23         text = text.lower()
24         text = re.sub(r'[^\w\s]', '', text) # remove punctuation
25         words = text.split()
26         # Remove common stop words manually
27         stop_words = {'the', 'a', 'an', 'and', 'or', 'but', 'in', 'on', 'at',
28                       'to', 'for', 'of', 'with', 'is', 'are', 'was', 'were',
29                       'be', 'been', 'being', 'have', 'has', 'had', 'do', 'does',
30                       'did', 'will', 'would', 'could', 'should', 'may', 'might',
31                       'can', 'this', 'that', 'these', 'those', 'i', 'you', 'he',
32                       'she', 'it', 'we', 'they', 'what', 'which', 'who', 'when',
33                       'where', 'why', 'how'}
34         return [w for w in words if w not in stop_words and len(w) > 1]
35
36     def _calculate_tf(self, words):
37         """Calculate Term Frequency for a document"""
38         word_count = len(words)
39         word_freq = Counter(words)
40         tf_scores = {}
41         for word, count in word_freq.items():
42             tf_scores[word] = count / word_count
43         return tf_scores
44
45     def _calculate_idf(self):
46         """Calculate Inverse Document Frequency for all words"""
47         num_docs = len(self.documents)
48         word_doc_count = {}
49
50         # Count how many documents contain each word
51         for doc in self.documents:
52             unique_words = set(doc)
53             for word in unique_words:
54                 word_doc_count[word] = word_doc_count.get(word, 0) + 1
55
56         # Calculate IDF: log(total_docs / docs_containing_word)
57         for word, doc_count in word_doc_count.items():
58             self.idf_values[word] = math.log(num_docs / doc_count)
59
60     def fit_transform(self, texts):
61         """Train on documents and return TF-IDF vectors"""
62         # Tokenize all documents
63         self.documents = [self._tokenize(text) for text in texts]
64
65         # Build vocabulary
66         all_words = set()
67         for doc in self.documents:
68             all_words.update(doc)
69
70         self.vocabulary = {word: idx for idx, word in enumerate(sorted(all_words))}
71
72         # Calculate IDF values
73         self._calculate_idf()
74
75         # Create TF-IDF vectors for all documents
76         vectors = []
77         for doc in self.documents:
78             vectors.append(self._create_vector(doc))
79
80         return vectors
81
82     def transform(self, texts):
83         """Transform new texts to TF-IDF vectors"""
84         vectors = []
85         for text in texts:
86             words = self._tokenize(text)
87
88             vectors.append(self._create_vector(words))
89         return vectors
90
91     def _create_vector(self, words):
92         """Create TF-IDF vector for a document"""
93         # Initialize vector with zeros
94         vector = [0.0] * len(self.vocabulary)
95
96         # Calculate TF for this document
97         tf_scores = self._calculate_tf(words)
98
99         # Calculate TF-IDF for each word
100         for word, tf in tf_scores.items():
101             if word in self.vocabulary:
102                 idx = self.vocabulary[word]
103                 idf = self.idf_values.get(word, 0)
104                 vector[idx] = tf * idf
105
106         return vector

```

## ➤ Manual Cosine Similarity Function

```

108 def cosine_similarity_manual(vec1, vec2):
109     """Calculate cosine similarity between two vectors manually"""
110     # Dot product
111     dot_product = sum(a * b for a, b in zip(vec1, vec2))
112
113     # Magnitude of vec1
114     magnitude1 = math.sqrt(sum(a * a for a in vec1))
115
116     # Magnitude of vec2
117     magnitude2 = math.sqrt(sum(b * b for b in vec2))
118
119     # Avoid division by zero
120     if magnitude1 == 0 or magnitude2 == 0:
121         return 0.0
122
123     # Cosine similarity
124     return dot_product / (magnitude1 * magnitude2)

```

## ➤ University Chatbot Matching Logic (Best Question Selection)

```

170 def find_best_match(self, user_question):
171     """Find the best matching question manually"""
172     if not self.is_trained:
173         return {
174             "answer": "Chatbot is not trained yet.",
175             "similarity_score": 0.0,
176             "matched_question": None
177         }
178
179     # Convert user question to vector using our custom vectorizer
180     user_vector = self.vectorizer.transform([user_question])[0]
181
182     # Manually calculate similarity with each stored question
183     similarities = []
184     for stored_vector in self.question_vectors:
185         sim = cosine_similarity_manual(user_vector, stored_vector)
186         similarities.append(sim)
187
188     # Find the best match manually
189     best_match_idx = 0
190     best_score = similarities[0]
191     for i, score in enumerate(similarities):
192         if score > best_score:
193             best_score = score
194             best_match_idx = i

```

## ➤ Flask API Routes

```

267 # API Routes
268 @app.route('/api/chat', methods=['POST'])
269 def chat():
270     """Handle chat messages"""
271     try:
272         data = request.get_json()
273         user_question = data.get('question', '')
274
275         if not user_question:
276             return jsonify({"error": "No question provided"}), 400
277
278         response = chatbot.find_best_match(user_question)
279         return jsonify(response), 200
280
281     except Exception as e:
282         return jsonify({"error": str(e)}), 500
283
284
285 @app.route('/api/stats', methods=['GET'])
286 def get_stats():
287     """Get chatbot statistics"""
288     try:
289         stats = chatbot.get_stats()
290         return jsonify(stats), 200
291     except Exception as e:
292         return jsonify({"error": str(e)}), 500
293
294
295 @app.route('/api/topics', methods=['GET'])
296 def get_topics():
297     """Get available topics"""
298     try:
299         topics = chatbot.get_topics()
300         return jsonify(topics), 200
301     except Exception as e:
302         return jsonify({"error": str(e)}), 500

```

```

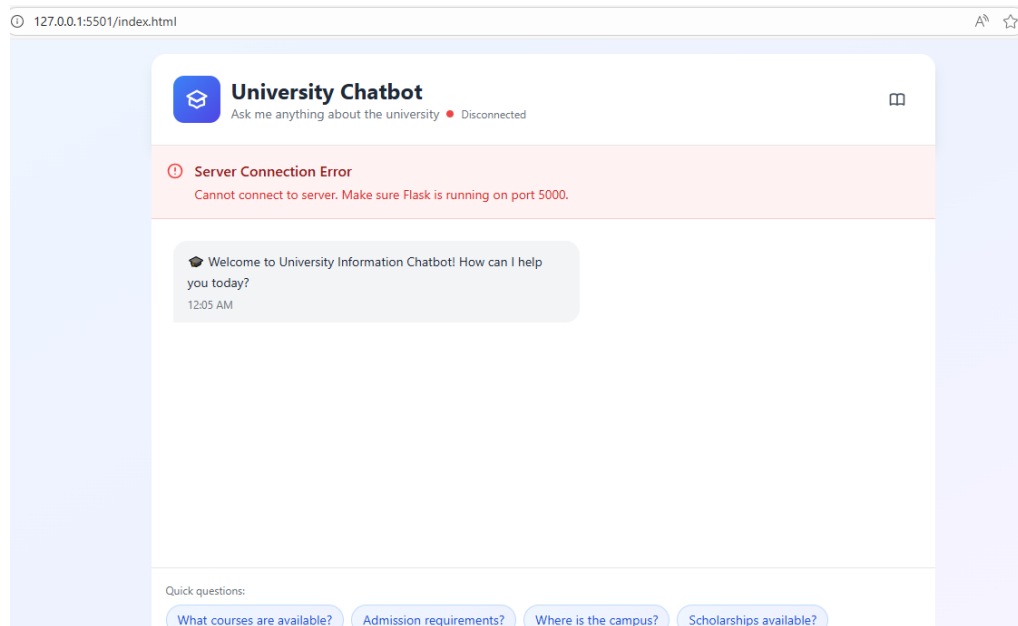
304
305 @app.route('/api/history', methods=['GET'])
306 def get_history():
307     """Get conversation history"""
308     try:
309         return jsonify(chatbot.conversation_history), 200
310     except Exception as e:
311         return jsonify({"error": str(e)}), 500
312
313
314 @app.route('/api/health', methods=['GET'])
315 def health_check():
316     """Health check endpoint"""
317     return jsonify({
318         "status": "healthy",
319         "is_trained": chatbot.is_trained,
320         "knowledge_base_size": len(chatbot.knowledge_base),
321         "implementation": "Custom (No AI Libraries)"
322     }), 200
323
324
325 @app.route('/api/add_data', methods=['POST'])
326 def add_data():
327     """Add new data to knowledge base"""
328     try:
329         data = request.get_json()
330         chatbot.add_data(data)
331         chatbot.train()
332         return jsonify({"message": "Data added successfully"}), 200
333     except Exception as e:
334         return jsonify({"error": str(e)}), 500
335

```



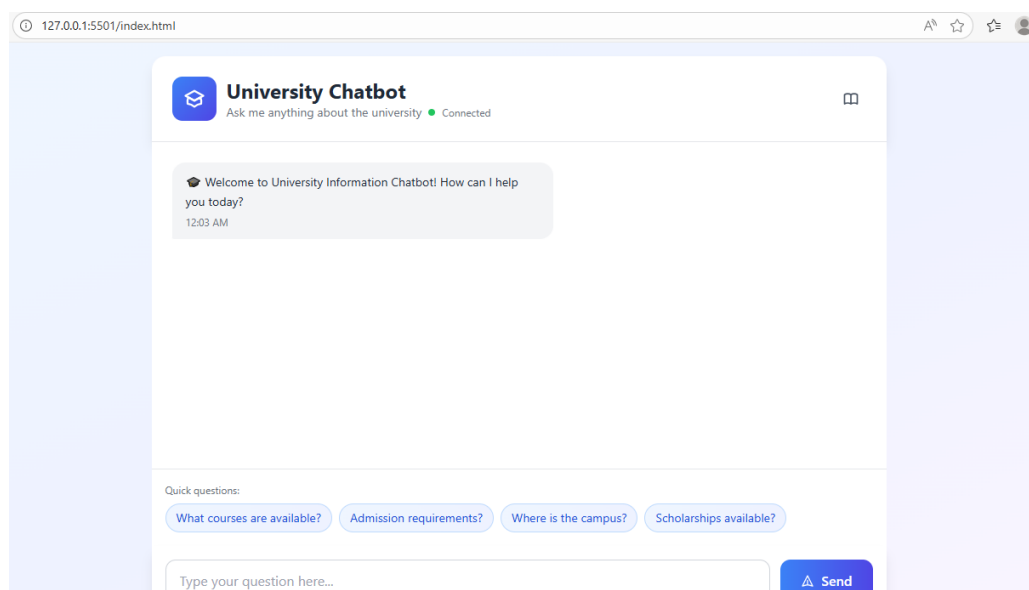
## 4. OUTPUT (SCREENSHOTS)

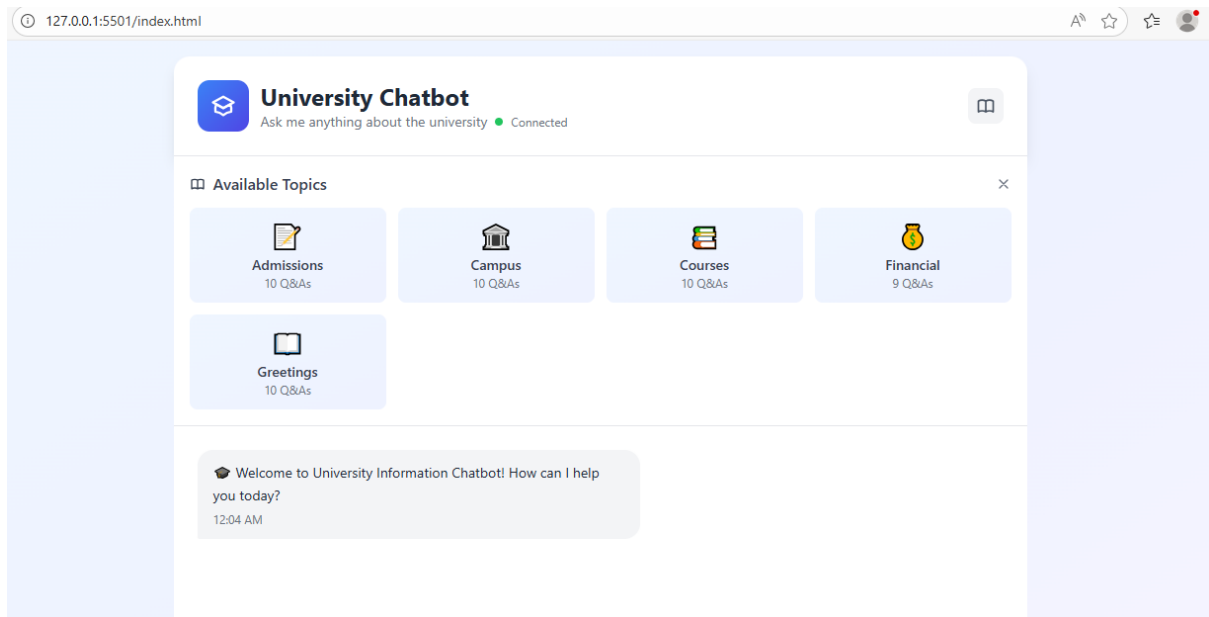
### 4.1 Chat Interface



```
📦 Loading university data...
* Running on http://127.0.0.1:5000
Press CTRL+C to quit
* Restarting with stat
📦 Loading university data...
✅ Loaded 49 Q&A pairs
🤖 Training chatbot (custom implementation)...
✅ Chatbot trained on 49 Q&A pairs

=====
🚀 University Chatbot API Server
(Custom Implementation - No AI Libraries)
=====
✅ Chatbot trained with 49 Q&A pairs
✅ Vocabulary size: 69 words
📢 Server running on http://localhost:5000
=====
```





---

## 5. SCOPE OF WORK

### 5.1 Backend Development

- Implemented custom TF-IDF without AI libraries
- Built preprocessing system
- Developed cosine similarity algorithm
- Created JSON-based knowledge base
- Designed Flask API with 6 endpoints
- Implemented conversation logging and statistics

---

### 5.2 Frontend Development

- Chat interface design
  - Statistics dashboard
  - Topics browser
  - History viewer
  - AJAX-based API communication
-

### 5.3 Algorithms Implemented

- Term Frequency (TF)
  - Inverse Document Frequency (IDF)
  - TF-IDF weighting
  - Vector Space Model
  - Cosine Similarity
- 

### 5.4 Skills Demonstrated

- Python programming
  - Object-Oriented Programming
  - Flask framework
  - RESTful API design
  - Web development
  - NLP fundamentals
  - Algorithm implementation
  - Data handling with JSON
- 

### 5.5 Challenges and Solutions

- **Understanding TF-IDF:** Solved by step-by-step testing.
  - **Slow Vector Comparisons:** Optimized using efficient loops.
  - **Different Question Phrasing:** Improved preprocessing and stop-word list.
  - **Low Similarity Scores:** Adjusted threshold and debugging vectors.
- 

### 5.6 Testing and Validation

- Unit testing of algorithms
- Integration testing of APIs
- User testing and feedback

- Performance testing for scalability
- 

## 5.7 Future Enhancements

- Spell checking
  - Voice interaction
  - Multi-language support
  - Mobile app version
  - Database integration
  - Authentication system
  - Admin dashboard
  - Cloud deployment
- 

## 6. GITHUB REPOSITORY LINK

### Repository URL:

[https://github.com/zaraimran03/AI\\_PROJECT.git](https://github.com/zaraimran03/AI_PROJECT.git)

### Installation Instructions:

1. Clone repository  
`git clone https://github.com/[username]/university-chatbot.git`
  2. Install dependencies  
`pip install flask flask-cors`
  3. Run application  
`python app.py`
  4. Open browser  
`http://localhost:5000`
- 

## CONCLUSION

This project successfully demonstrates the implementation of a chatbot system using fundamental NLP techniques without relying on AI libraries. By manually implementing TF-IDF and cosine similarity, the system provides accurate, fast, and scalable responses to user queries.

**Key Achievements:**

- Built custom TF-IDF algorithm
- Implemented manual cosine similarity
- Developed REST API using Flask
- Designed user-friendly interface
- Created comprehensive documentation

This project reflects strong understanding of AI fundamentals, software engineering principles, and full-stack development.

---