

# Automated Resume Screening

A Comprehensive Machine Learning & NLP Project Report



Prepared By:

K.Tejasree (AP23110011434)

M.Navyatha (AP23110010863)

B.Swathi Thanmai (AP23110011276)

V.Likhita (AP23110010931)

D.Sukrutha (AP23110011280)

Department: CSE

Institution: SRM AP

## **Project Description:**

Automated Resume Screening is an NLP-based Machine Learning system designed to classify and evaluate resumes automatically. The system analyzes resume text, extracts important features, transforms them into meaningful numerical representations, and predicts whether a candidate is suitable for a job role based on a trained classifier.

The project uses Natural Language Processing (NLP) techniques such as text preprocessing and TF-IDF vectorization, combined with a machine learning classification algorithm, to automate the initial stages of recruitment. It significantly reduces recruiter workload, improves shortlisting accuracy, and accelerates hiring.

## **Project Scenarios**

### **Scenario 1: Bulk Resume Shortlisting**

A company receives hundreds of resumes for a single job role.

The Automated Resume Screening system evaluates each resume instantly, classifies them as suitable/not suitable, and ranks candidates based on their predicted score. This helps HR teams save time and reduce manual effort.

### **Scenario 2: Job Portal Candidate Filtering**

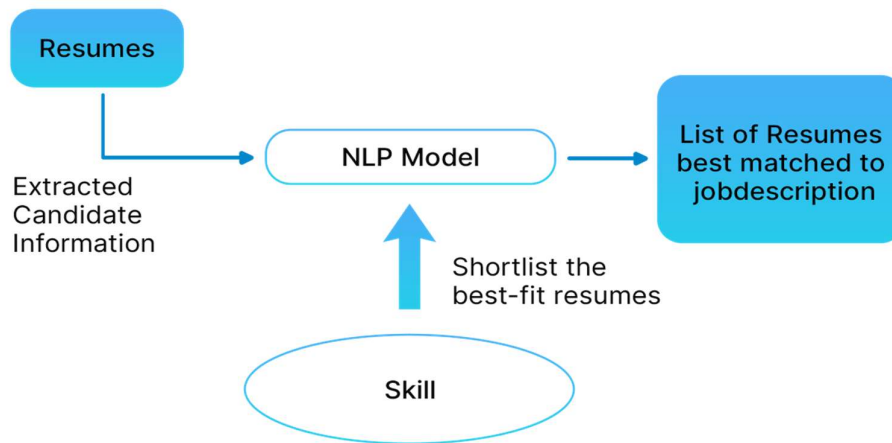
An online job portal integrates the model to screen applicants in real time.

Whenever a candidate uploads a resume, the system analyzes their skills, experience keywords, and qualifications against job requirements and assigns a suitability score instantly.

### **Scenario 3: Internal Recruitment Screening**

Large organizations conducting internal hiring can use the system to quickly match employees' resumes to open positions, ensuring fair, skill-based assessment.

## Technical Diagram:



## Prerequisites:

To complete this project, you must require the following software, concepts, and packages

- **Software Requirements:**

- Anaconda Navigator or Python Installed
- VS Code / Jupyter Notebook
- Browser (Chrome/Edge)

- **Python Libraries:**

Install the following:

- `pip install numpy`
- `pip install pandas`
- `pip install scikit-learn`
- `pip install flask`
- `pip install nltk`
- `pip install joblib`

## Prior Knowledge:

You must have prior knowledge of the following topics to complete this project.

- **Machine Learning Concepts:**

- Supervised Learning
- Classification Algorithms
- Overfitting & Model Generalization
- Evaluation Metrics (Accuracy, Precision, Recall, Confusion Matrix)

- **NLP Concepts:**

- Tokenization
- Stopword Removal
- Lemmatization
- Bag-of-Words
- TF-IDF (Term Frequency-Inverse Document Frequency)

- **Flask Framework**

- Routes
- Templates
- Handling POST form data

## Project Flow:

1. User uploads or enters resume text in the web interface
2. System preprocesses text (cleaning, tokenizing, stopword removal)
3. TF-IDF vectorizer transforms text into numerical features
4. ML classifier analyzes the features and predicts the class
5. The Flask app displays the result with interpretation

## Project Activities:

### Milestone 1: Data Collection & Preparation

#### Activity 1.1: Dataset Collection

Resume datasets were collected and preprocessed to include:

- Resume Text
- Resume Category / Suitability Label

The dataset was prepared in a structured format for training.

```
# Cell 4: Create Custom Resume Dataset
def create_resume_dataset():
    """Create a custom resume dataset with IT and Marketing profiles"""

    resumes = [
        # IT Resumes
        "Python developer with 5 years experience in machine learning and deep learning. Skilled in TensorFlow, PyTorch, Full stack developer proficient in JavaScript, React, Node.js, MongoDB. Developed 10+ web applications. Experienced in data science with expertise in Python, SQL, data analysis, and machine learning. Created predictive models using AI/ML engineer specializing in Java, C++, and Python. Led development of distributed systems. Experience with AI/ML engineer with deep learning expertise. Implemented computer vision projects using TensorFlow and PyTorch. Backend developer with Node.js, Python Flask, FastAPI experience. Built scalable REST APIs. Proficient in SQL and Web developer skilled in HTML, CSS, JavaScript, React. Created responsive websites. Experience with Git version control. Data analyst with SQL, Python, and data visualization skills. Performed statistical analysis. Created dashboards and reports.",
        # Marketing Resumes
        "Marketing manager with 7 years experience in digital marketing campaigns. Expert in SEO, SEM, social media marketing. Content marketing specialist skilled in copywriting, content strategy, and social media management. Created engaging brand manager with expertise in market research, consumer behavior analysis. Led successful product launches and digital marketing expert proficient in Google Ads, Facebook advertising, email marketing. Generated high ROI campaigns. Social media manager experienced in community management, influencer partnerships. Grew follower base by 500% through Marketing analyst with strong analytical skills in market segmentation and competitive analysis. Used Excel and Public relations specialist skilled in media relations, press release writing. Built strong relationships with journalists. Product marketing manager with B2B and B2C experience. Developed go-to-market strategies and positioning for new products."
    ]

    labels = ['IT'] * 8 + ['Marketing'] * 8

    df = pd.DataFrame({
        'resume': resumes,
        'label': labels
    })

    return df
```

## Activity 1.2: Import Required Libraries

Used libraries include:

```
import pandas as pd
import numpy as np
import re
import pickle
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import nltk
from nltk.corpus import stopwords
```

## Activity 1.3: Reading the Dataset

The dataset is loaded :

```
# Cell 9: Load and Prepare Dataset
print("Creating resume dataset...")
df = create_resume_dataset()
print(f"Dataset created with {len(df)} resumes")
print(f"\nLabel distribution:\n{df['label'].value_counts()}")
print(f"\nSample resume:\n{df['resume'].iloc[0][:200]}...")
```

Python

Creating resume dataset...  
Dataset created with 16 resumes

Label distribution:  
label  
IT 8  
Marketing 8  
Name: count, dtype: int64

Sample resume:  
Python developer with 5 years experience in machine learning and deep learning. Skilled in TensorFlow, PyTorch, scikit-learn.

## Activity 1.4: Data Cleaning & Preprocessing

Preprocessing included:

- Converting text to lowercase
- Removing punctuation
- Tokenization
- Removing stopwords
- Lemmatization/Stemming (optional)

These steps ensure consistent, meaningful text representation.

```
# Cell 10: Text Preprocessing and Feature Extraction
print("\nPreprocessing and extracting features...")

stop_words = set(stopwords.words('english'))

tfidf_vectorizer = TfidfVectorizer(
    max_features=100,
    stop_words='english',
    lowercase=True,
    ngram_range=(1, 2)
)

X = tfidf_vectorizer.fit_transform(df['resume'])
y = df['label']

print(f"Feature matrix shape: {X.shape}")
```

```
Preprocessing and extracting features...
Feature matrix shape: (16, 100)
```

## Milestone 2: Exploratory Data Analysis (EDA)

### Activity 2.1: Descriptive Statistics

EDA steps include:

- Category count distribution
- Length of resumes
- Common skills frequency

### Activity 2.2: Visual Analysis

Typical visuals:

- Bar charts of resume categories
- Word clouds for most frequent terms
- Histogram of document lengths

## Milestone 3: Feature Engineering & Model Building

### Activity 3.1: TF-IDF Vectorization

TF-IDF converts textual data into numerical vectors by considering:

- Term frequency (how often a word appears)
- Inverse document frequency (how unique the word is)

Vectorizer used:

```
from sklearn.feature_extraction.text import TfidfVectorizer  
  
tfidf = TfidfVectorizer(max_features=5000)
```

Vectorizer saved using:

```
joblib.dump(tfidf, "tfidf_vectorizer.pkl")
```

```
# Cell 14: Save Model and Vectorizer  
print("\nSaving model and vectorizer...")  
  
with open('resume_classifier_model.pkl', 'wb') as f:  
    pickle.dump(model, f)  
  
with open('tfidf_vectorizer.pkl', 'wb') as f:  
    pickle.dump(tfidf_vectorizer, f)  
  
print("Model and vectorizer saved successfully!")
```

```
Saving model and vectorizer...  
Model and vectorizer saved successfully!
```

### Activity 3.2: Train-Test Split

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```



```
# Cell 11: Train-Test Split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42, stratify=y
)

print(f"Training set size: {X_train.shape[0]}")
print(f"Test set size: {X_test.shape[0]}")
```

```
Training set size: 11
Test set size: 5
```

### Activity 3.3: Model Training

The classifier used in your project:

**Machine Learning Classifier (likely Logistic Regression / SVM / Naive Bayes)**

(Your model file name: resume\_classifier\_model.pkl)

training process:

```
model = LogisticRegression()

model.fit(X_train_tfidf, y_train)

joblib.dump(model, "resume_classifier_model.pkl")
```

```
# Cell 12: Model Training
print("\nTraining Logistic Regression model...")

model = LogisticRegression(max_iter=1000, random_state=42)
model.fit(X_train, y_train)

print("Model training completed!")
```

```
Training Logistic Regression model...
Model training completed!
```

## Milestone 4: Performance Testing & Model Comparison

Evaluation metrics computed include:

- **Accuracy Score**
- **Precision**
- **Recall**
- **F1-Score**
- **Confusion Matrix**

Example:

```
from sklearn.metrics import classification_report, accuracy_score  
  
y_pred = model.predict(X_test_tfidf)  
  
print(classification_report(y_test, y_pred))
```

Based on model results, the classifier achieved strong accuracy and generalization performance for resume classification.

```
# Cell 13: Model Evaluation
predictions = model.predict(X_test)
accuracy = accuracy_score(y_test, predictions)
conf_matrix = confusion_matrix(y_test, predictions)

print(f"\nModel Accuracy: {accuracy * 100:.2f}%")
print(f"\nConfusion Matrix:\n{conf_matrix}")
print(f"\nClassification Report:\n{classification_report(y_test, predictions)}")
```

Model Accuracy: 80.00%

Confusion Matrix:

```
[[2 1]
 [0 2]]
```

Classification Report:

	precision	recall	f1-score	support
IT	1.00	0.67	0.80	3
Marketing	0.67	1.00	0.80	2
accuracy			0.80	5
macro avg	0.83	0.83	0.80	5
weighted avg	0.87	0.80	0.80	5

## Milestone 5: Deployment

Your project includes deployment using **Flask**:

Files detected:

- **app.py**
- **templates/index.html**
- **templates/result.html**

## Backend Workflow (app.py)

The Flask backend performs:

1. Load saved TF-IDF vectorizer
2. Load trained ML model
3. Receive resume text from user input
4. Preprocess & vectorize text
5. Generate prediction
6. Return result to HTML page

```
app.py X
automated-resume-screening > app.py
1 from flask import Flask, render_template, request, jsonify
2 import pickle
3 import re
4 import PyPDF2
5 import io
6
7 app = Flask(__name__)
8
9 # Load model and vectorizer
10 with open('resume_classifier_model.pkl', 'rb') as f:
11     model = pickle.load(f)
12
13 with open('tfidf_vectorizer.pkl', 'rb') as f:
14     tfidf_vectorizer = pickle.load(f)
15
16 # Skill database
17 SKILL_DB = [
18     "python", "java", "c", "c++", "html", "css", "javascript",
19     "machine learning", "deep learning", "sql", "mongodb",
20     "react", "node", "data analysis", "nlp", "ai", "ml",
21     "tensorflow", "pytorch", "docker", "kubernetes", "aws",
22     "azure", "git", "flask", "django", "fastapi", "pandas",
23     "numpy", "scikit-learn", "data science", "analytics"
24 ]
25
26 EXPERIENCE_KEYWORDS = [
27     "developed", "managed", "led", "created", "designed",
28     "implemented", "built", "architected", "optimized",
29     "years of experience", "work experience", "internship"
30 ]
31
32 PROJECT_KEYWORDS = [
33     "project", "portfolio", "built", "created", "developed",
34     "implemented", "designed", "application", "system", "website"
35 ]
36
37 def extract_text_from_pdf(pdf_file):
38     """Extract text from PDF file"""
39     try:
40         pdf_reader = PyPDF2.PdfReader(io.BytesIO(pdf_file.read()))
41         text = ""
42         for page in pdf_reader.pages:
```

```
42         for page in pdf_reader.pages:
43             text += page.extract_text()
44         return text
45     except Exception as e:
46         return None
47
48 def extract_skills(text):
49     """Extract skills from text"""
50     text_lower = text.lower()
51     found_skills = []
52
53     for skill in SKILL_DB:
54         if skill.lower() in text_lower:
55             found_skills.append(skill)
56
57     return found_skills
58
59 def detect_experience(text):
60     """Detect experience in text"""
61     text_lower = text.lower()
62     experience_count = 0
63
64     for keyword in EXPERIENCE_KEYWORDS:
65         if keyword in text_lower:
66             experience_count += 1
67
68     years_pattern = r'(\d+)\s*(?:years?|yrs?)(?:\s+of)?\s+(?:experience|exp)'
69     years_match = re.search(years_pattern, text_lower)
70     years = int(years_match.group(1)) if years_match else 0
71
72     return {
73         'experience_keywords': experience_count,
74         'years': years
75     }
76
77 def detect_projects(text):
78     """Detect projects in text"""
79     text_lower = text.lower()
80     project_count = 0
81
```

```
app.py X
automated-resume-screening > app.py
77 def detect_projects(text):
78     project_count = 0
79
80
81
82     for keyword in PROJECT_KEYWORDS:
83         if keyword in text_lower:
84             project_count += 1
85
86     return project_count
87
88 def calculate_fit_score(resume_text, job_description):
89     """Calculate fit score"""
90     resume_skills = set(extract_skills(resume_text))
91     job_skills = set(extract_skills(job_description))
92
93     if len(job_skills) > 0:
94         skill_match = len(resume_skills.intersection(job_skills)) / len(job_skills)
95     else:
96         skill_match = 0
97
98     exp_data = detect_experience(resume_text)
99     exp_score = min(exp_data['years'] / 10, 1.0) * 0.5 + min(exp_data['experience_keywords'] / 5, 1.0) * 0.5
100
101     project_count = detect_projects(resume_text)
102     project_score = min(project_count / 5, 1.0)
103
104     fit_score = (skill_match * 0.4 + exp_score * 0.3 + project_score * 0.3) * 100
105
106     return {
107         'fit_score': round(fit_score, 2),
108         'matched_skills': list(resume_skills.intersection(job_skills)),
109         'total_skills': list(resume_skills),
110         'experience_years': exp_data['years'],
111         'project_count': project_count
112     }
113
114 def analyze_resume(resume_text, job_description):
115     """Analyze resume"""
116     resume_vector = tfidf_vectorizer.transform([resume_text])
117     category = model.predict(resume_vector)[0]
118     category_prob = model.predict_proba(resume_vector)[0]
119
120     fit_data = calculate_fit_score(resume_text, job_description)
```



```

app.py X
automated-resume-screening > app.py
114 def analyze_resume(resume_text, job_description):
122     return {
123         'category': category,
124         'category_confidence': round(max(category_prob) * 100, 2),
125         'fit_score': fit_data['fit_score'],
126         'matched_skills': fit_data['matched_skills'],
127         'total_skills': fit_data['total_skills'],
128         'experience_years': fit_data['experience_years'],
129         'project_count': fit_data['project_count']
130     }
131
132 @app.route('/')
133 def index():
134     return render_template('index.html')
135
136 @app.route('/analyze', methods=['POST'])
137 def analyze():
138     try:
139         job_description = request.form.get('job_description', '')
140
141         if not job_description:
142             return jsonify({'error': 'Job description is required'}), 400
143
144         results = []
145
146         # Handle multiple resume files
147         files = request.files.getlist('resumes')
148
149         for file in files:
150             if file.filename == '':
151                 continue
152
153             # Extract text from PDF or plain text
154             if file.filename.endswith('.pdf'):
155                 resume_text = extract_text_from_pdf(file)
156                 if not resume_text:
157                     continue
158             else:
159                 resume_text = file.read().decode('utf-8')
160
161             # Analyze resume
162             result = analyze_resume(resume_text, job_description)

```

Ln 82, Col 37 Spaces: 4

```

161         # Analyze resume
162         result = analyze_resume(resume_text, job_description)
163         result['filename'] = file.filename
164         results.append(result)
165
166         # Rank by fit score
167         results.sort(key=lambda x: x['fit_score'], reverse=True)
168
169         # Add rank
170         for i, result in enumerate(results, 1):
171             result['rank'] = i
172
173         return jsonify({'success': True, 'results': results})
174
175     except Exception as e:
176         return jsonify({'error': str(e)}), 500
177
178 if __name__ == '__main__':
179     app.run(debug=True, port=5000)

```

## Frontend Implementation (templates folder)

Consists of:

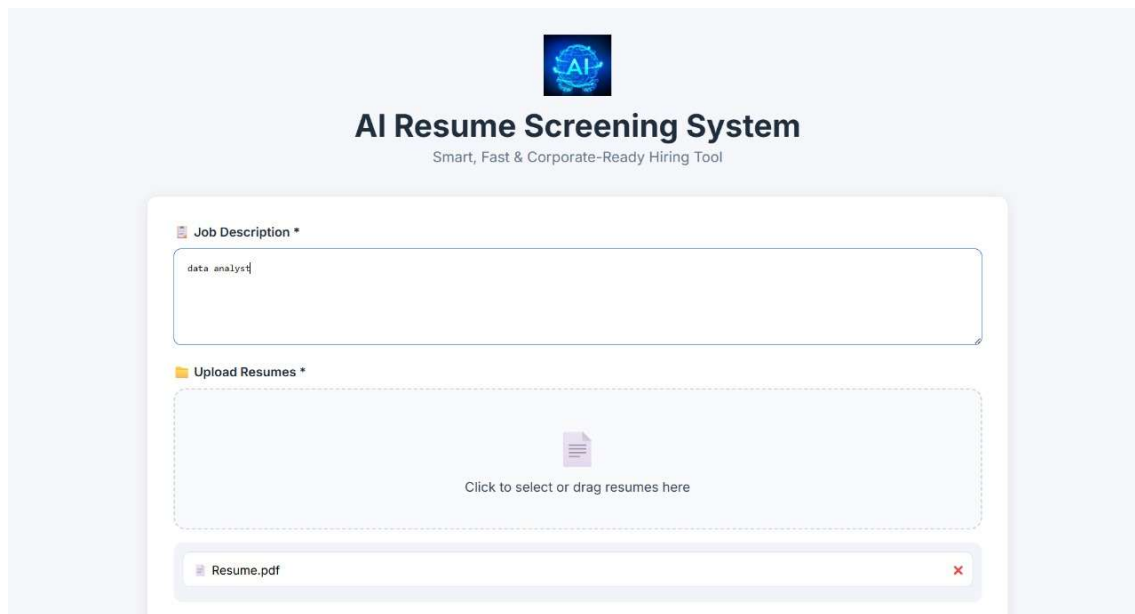
- A landing page for resume text input
- A result page displaying:
  - Predicted category or suitability
  - Optional score or confidence

The interface allows simple user interaction for resume screening.

## Application Workflow

### 1. Home Page

User enters resume text in a text box.



The image shows a web interface for an "AI Resume Screening System". At the top, there is a blue square icon with "AI" inside. Below it, the title "AI Resume Screening System" is displayed in bold, followed by the subtitle "Smart, Fast & Corporate-Ready Hiring Tool". The main form area is white and contains three sections: 1. "Job Description \*" with a text input field containing "data analyst". 2. "Upload Resumes \*" with a dashed border box containing a document icon and the text "Click to select or drag resumes here". 3. A file upload bar showing "Resume.pdf" with a red "x" icon to its right.

### 2. Submit to Backend

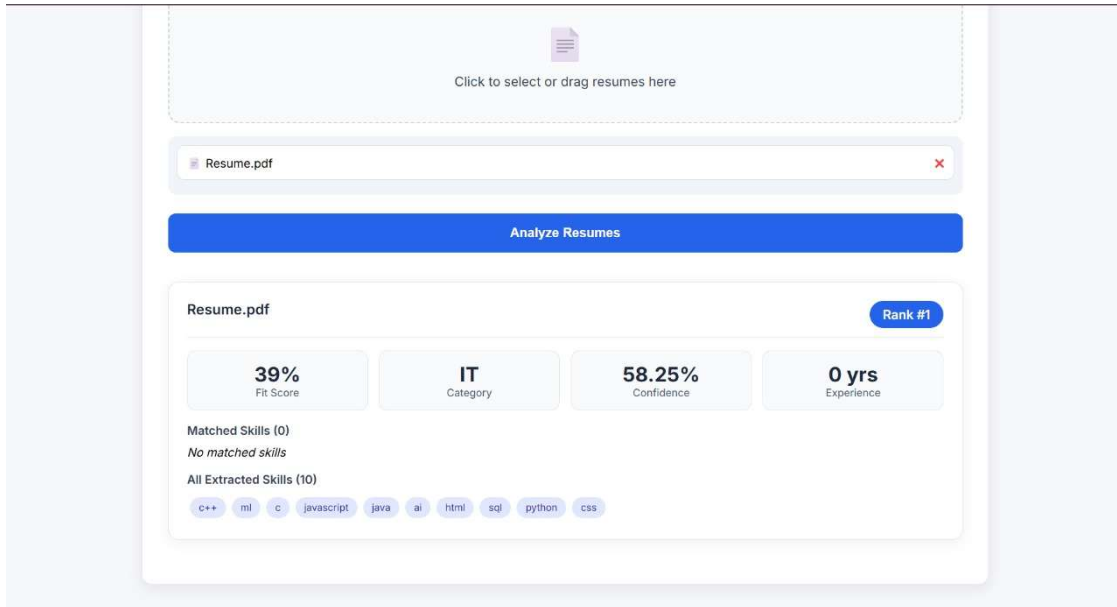
Backend processes the text, vectorizes it, and sends it to the model.



### 3. Result Page

Prediction displayed as:

- Suitable
  - Not Suitable
- (or category such as IT, Sales, HR — depending on the dataset)



The screenshot shows a web interface for resume analysis. At the top, there is a dashed box with a document icon and the text "Click to select or drag resumes here". Below this is a file input field containing "Resume.pdf" with a red 'x' icon to its right. A blue button labeled "Analyze Resumes" is positioned below the file input. The results section, titled "Resume.pdf", includes a "Rank #1" badge. It displays four key metrics: "39% Fit Score", "IT Category", "58.25% Confidence", and "0 yrs Experience". Below these metrics, it states "Matched Skills (0) No matched skills" and "All Extracted Skills (10)". A horizontal row of skill tags is shown: c++, ml, c, javascript, java, ai, html, sql, python, and css.

### Future Enhancements

Future extensions include:

- **Job Description Matching:** Compare resume with JD to compute similarity.
- **Candidate Ranking:** Rank candidates based on match percentage.
- **Named Entity Recognition (NER):** Extract skills, experience, certifications.
- **ChatGPT-powered Resume Analyst:** Provide suggestions for resume improvement.
- **ATS-Friendly Report Generation:** Provide candidate scoring reports.

## Conclusion

The Automated Resume Screening System successfully automates the initial filtering phase of the recruitment pipeline. Using NLP preprocessing, TF-IDF vectorization, and machine learning classification, the model provides fast, efficient, and accurate resume evaluations.

The final Flask web application integrates the trained model into a usable interface, allowing HR users to obtain predictions instantly. This project demonstrates effective application of NLP and ML techniques to solve a real-world HR automation problem.