



Smart Internship **Recommendation System**

A Content-Based Filtering Approach using Python & Machine Learning

Project By: [Your Name]

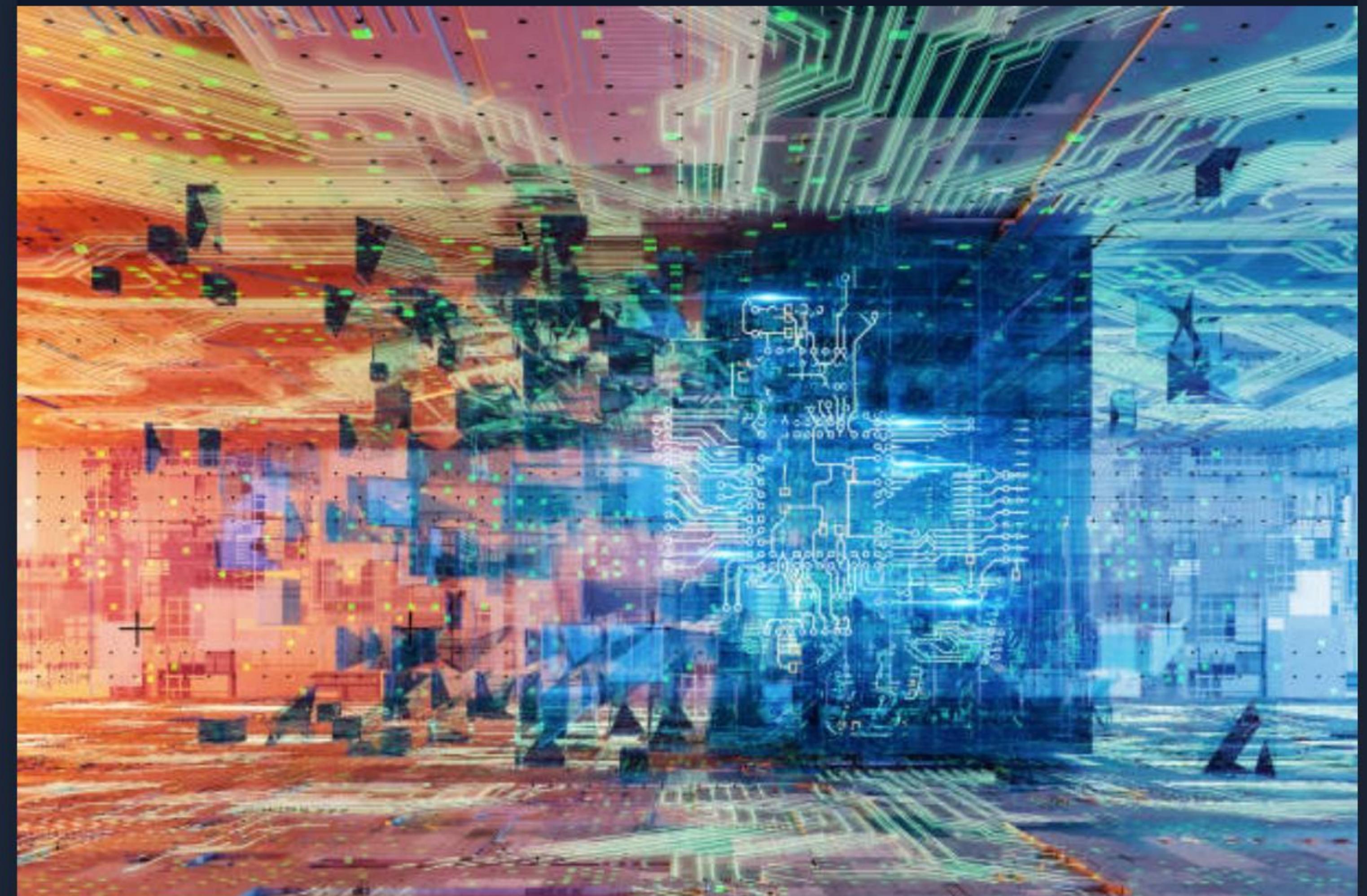
The Challenge

Information Overload

Students today are bombarded with thousands of internship listings. Filtering through irrelevant roles to find one that matches their specific skills, location, and preferred sector is time-consuming and inefficient.

The Mismatch

Manual searching often leads to missed opportunities where a student's qualifications perfectly match a company's requirements, but the listing is buried under noise.



Our Solution: Intelligent Matching



Automated & Personalized

We built an AI-driven engine that moves beyond simple keyword search. By analyzing the "vectorized" signature of a candidate's profile, we calculate the mathematical probability of a fit.

- ✓ **Skill-Centric:** Uses Machine Learning to match technical abilities.
- ✓ **Context-Aware:** Factors in location, sector, and social category.
- ✓ **Ranked Output:** Provides a prioritized list of top opportunities.

Technology Stack



Python 3.x

The core programming language used for data processing and logic implementation due to its rich ecosystem of AI libraries.



Pandas

Used for structured data manipulation, creating the internship database, and handling user inputs efficiently.



Scikit-Learn

The machine learning powerhouse. Used for **MultiLabelBinarizer** (Encoding) and **Cosine Similarity** (Matching).

Process Architecture



1. Input

User provides Skills, Location, Sector, and Category.



2. Vectorize

Convert text skills into binary vectors (0s and 1s).



3. Compute

Calculate Cosine Similarity & apply weighted formula.



4. Rank

Sort results by Final Score and display to user.

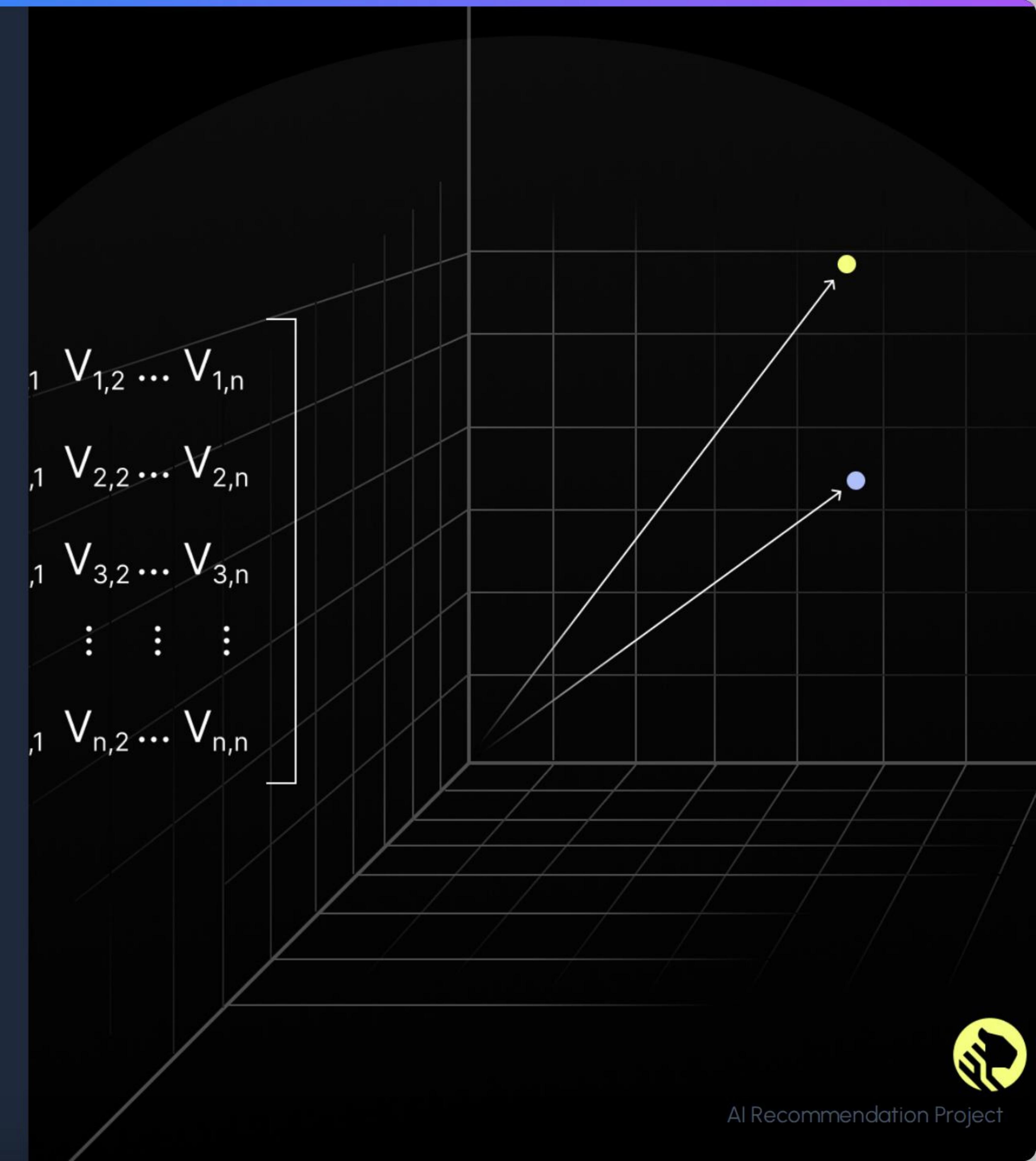
Core Algo: Cosine Similarity

To the computer, "Python" and "Java" are just math.

We represent every student and every internship as a **Vector** in multi-dimensional space.

The Logic:

- ✓ If the angle (θ) between vectors is 0° , they are identical (Score = 1.0).
- ✓ If the angle is 90° , they are unrelated (Score = 0.0).
- ✓ This allows us to mathematically quantify "Closeness."



The "Hybrid" Scoring Logic

We don't rely on AI alone. We use a **Weighted Hybrid Approach** that combines the semantic skill match with hard business constraints (Location, Sector).

$$\text{FinalScore} = (\text{Skill}) + 0.2(\text{Location}) + 0.2(\text{Sector}) + 0.1(\text{Category})$$

50% AI Match

The Cosine Similarity of the skills vector drives half the decision.

40% Preferences

Strict matching for Location and Industry Sector.

10% Diversity

Bonus points for matching priority social categories.

Data Pre-Processing

MultiLabelBinarizer

Computers cannot read comma-separated lists like "Python, SQL". We transform these lists into a **One-Hot Encoded Matrix**.

This creates a master list of all possible skills. If an internship requires a skill, it gets a '1', otherwise '0'.

```
mlb = MultiLabelBinarizer()  
matrix = mlb.fit_transform(internships['Skills'])
```

1	1	0	1	1	1	1	1	1	0	0	0	0	1	0	1	1	1	0	0	1	1	0						
1	1	1	1	0	0	0	1	0	0	1	1	0	1	0	1	1	0	0	0	0	1	1	0					
0	0	0	0	1	1	0	1	1	1	0	1	0	1	0	0	0	0	1	1	1	1	0	1	0				
0	0	1	0	0	0	0	1	0	0	1	0	1	0	0	0	1	1	1	0	0	1	0	0	1	1			
1	1	0	0	1	0	0	0	0	1	1	0	1	0	1	1	1	0	1	1	0	1	1	0	0	0			
0	0	1	1	0	1	1	1	0	0	0	0	0	0	0	0	0	1	1	0	1	0	0	1	1	0			
1	1	0	0	1	0	0	1	1	1	1	1	1	1	1	1	1	0	1	1	0	1	1	0	0	1	1		
1	1	0	0	1	0	0	1	0	1	1	1	1	1	1	1	1	0	1	1	0	1	1	0	0	1	1		
0	0	1	1	0	0	0	0	1	0	0	0	0	0	1	0	0	1	0	0	1	1	0	0	0	0	1	1	
0	0	0	0	1	0	1	1	0	0	0	1	1	1	1	1	0	0	0	0	0	0	1	1	1	1	1	1	
1	1	0	1	0	1	1	0	0	0	1	1	1	1	1	0	0	1	1	1	0	1	1	1	0	0	0	1	1
1	1	0	1	0	1	0	0	0	1	0	1	1	1	1	1	0	0	1	1	1	0	1	1	0	0	0	1	1
1	0	1	0	1	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	
0	1	0	1	0	0	0	0	1	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	1	1	
0	1	0	1	0	0	0	1	1	0	0	0	0	0	1	1	1	1	0	0	0	0	0	0	0	0	1	1	
1	1	0	1	0	1	1	0	0	0	1	1	1	1	1	1	0	0	1	1	1	0	1	1	0	0	0	1	1
1	1	0	1	0	1	0	0	0	1	0	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	1	1	
0	1	0	1	0	0	0	1	1	0	0	0	0	0	1	0	0	0	1	1	1	0	0	0	0	0	1	1	
0	0	0	0	0	0	0	1	1	0	0	0	0	0	1	0	0	0	1	1	1	0	0	0	0	0	1	1	
1	1	1	1	1	1	0	0	1	1	1	1	1	1	1	1	0	0	1	1	1	1	1	0	0	0	0	1	1
1	1	1	1	1	1	0	0	1	1	1	1	1	1	1	1	0	0	1	1	1	1	1	1	0	0	0	1	1
0	0	0	0	1	0	0	1	1	0	0	1	0	1	0	1	0	0	1	1	1	0	0	0	1	0	1	1	
0	0	0	1	1	1	1	0	0	0	1	0	1	0	1	0	0	0	0	0	0	1	1	1	1	0	1	0	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	1	1	1	1	1	0	0	1	1
1	1	0	1	0	0	0	1	1	0	1	0	1	1	0	1	0	0	1	1	1	0	0	0	1	0	1	0	1
0	1	0	1	0	0	0	0	1	0	1	0	0	1	1	0	0	0	0	1	1	0	0	0	1	0	0	0	1
0	1	1	0	1	0	0	1	0	0	0	1	1	1	0	0	0	0	0	1	1	0	0	0	1	0	0	0	1
1	0	0	0	0	0	0	0	1	0	0	0	1	1	1	1	1	0	0	0	0	0	0	0	0	1	0	1	1
0	1	1	0	1	0	0	1	0	0	0	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	1	0	0	1	0	1	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	1	1	1	1	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	0	1	0	0	1	0	1	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	1	1	1	1	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	0	1	0	0	1	0	1	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	1	1	1	1	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	0	1	0	0	1	0	1	0	0	1	0	1	0													

The Matching Engine

The core iteration loop that drives the recommendations:

```
for i, row in internships.iterrows(): # 1. Calculate Cosine Similarity for Skills skill_score = cosine_similarity([candidate_vec], [internship_matrix[i]])[0][0] # 2. Check Constraints (Binary 1 or 0) loc_match = int(user_loc == row['Location']) sec_match = int(user_sec == row['Sector']) # 3. Apply Weighted Formula final_score = (0.5 * skill_score) + (0.2 * loc_match) + ... results.append({'Company': row['Name'], 'Score': final_score})
```

System Output

If a user inputs: **Skills: "Python", Location: "Delhi", Sector: "IT"**

The system generates a ranked dataframe:

Rank	Internship ID	Course Name	Company	Final Score
1	101	Data Analytics	TechCorp	0.95 (High Match)
2	103	Machine Learning	MLLabs	0.82 (Good Match)
3	102	Finance Internship	FinanceInc	0.10 (Low Match)

Future Scope



Web Integration

Deploying the Python logic as a Flask/Django API to create a user-friendly web interface.



Resume Parsing

Integrating OCR to automatically extract skills from uploaded PDF resumes instead of manual typing.



Collaborative Filtering

Adding a "Users like you applied for..." feature based on historical application data.



Questions & Answers

Thank you for your attention.

Image Sources



https://media.istockphoto.com/id/1485546343/photo/abstract-cloud-computing-image.jpg?s=612x612&w=0&k=20&c=zIMwzGpAQCnzbORGR-OQxRBX_XXWBNav6TLDiwTg=

Source: www.istockphoto.com



<https://static.vecteezy.com/system/resources/thumbnails/067/727/242/small/concept-of-seo-content-marketing-and-blog-article-optimization-includes-rating-stars-web-layout-coding-and-target-icons-perfect-for-web-digital-marketing-and-editorial-design-vector.jpg>

Source: www.vecteezy.com



https://www.tigerdata.com/_next/image?url=https%3A%2F%2Ftimescale.ghost.io%2Fblog%2Fcontent%2Fimages%2F2024%2F04%2FA-Beginner-s-Guide-to-Vector-Embeddings-1.png&w=3840&q=100

Source: www.tigerdata.com



https://static.vecteezy.com/system/resources/previews/006/060/283/non_2x/binary-code-zero-one-matrix-white-background-technology-connection-digital-data-abstract-background-vector.jpg

Source: www.vecteezy.com