

In [1]:

```

1 import pandas as pd
2
3 # Sample data
4 data = pd.DataFrame({
5     'Gender': ['Male', 'Female', 'Female', 'Male', 'Male', 'Female'],
6     'Department': ['HR', 'IT', 'Sales', 'HR', 'IT', 'Sales'],
7     'Education': ['Bachelor', 'Master', 'PhD', 'Master', 'PhD', 'Bachelor'],
8     'Salary': [50000, 60000, 55000, 52000, 58000, 53000],
9     'Experience': [1.5, 3.0, 4.2, 2.3, 3.7, 2.9]
10 })
11
12 data.head()
13

```

c:\users\vamsi2001\appdata\local\programs\python\python39\lib\site-packages\numpy_distributor_init.py:30: UserWarning: loaded more than 1 DLL from .libs:
c:\users\vamsi2001\appdata\local\programs\python\python39\lib\site-packages\numpy\.libs\libopenblas.EL2C6PLE4ZYW3ECEVIV30XXGRN2NRFM2.gfortran-win_amd64.dll
c:\users\vamsi2001\appdata\local\programs\python\python39\lib\site-packages\numpy\.libs\libopenblas.XWYDX2IKJW2NMTWSFYNGFUWKQU3LYTCZ.gfortran-win_amd64.dll
warnings.warn("loaded more than 1 DLL from .libs:")

Out[1]:

	Gender	Department	Education	Salary	Experience
0	Male	HR	Bachelor	50000	1.5
1	Female	IT	Master	60000	3.0
2	Female	Sales	PhD	55000	4.2
3	Male	HR	Master	52000	2.3
4	Male	IT	PhD	58000	3.7

In [2]:

```

1 data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Gender          6 non-null     object
1   Department      6 non-null     object
2   Education       6 non-null     object
3   Salary          6 non-null     int64
4   Experience       6 non-null     float64
dtypes: float64(1), int64(1), object(3)
memory usage: 368.0+ bytes

```

```
In [12]: 1 #Use LabelEncoder to convert categories into integers
2 from sklearn.preprocessing import LabelEncoder
3
4 # Make a copy to preserve original
5 encoded_data = data.copy()
6
7 le = LabelEncoder()
8 encoded_data['Gender'] = le.fit_transform(data['Gender'])
9 encoded_data['Department'] = le.fit_transform(data['Department'])
10 encoded_data['Education'] = le.fit_transform(data['Education'])
11 encoded_data.head()
12 encoded_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Gender      6 non-null      int32
1   Department  6 non-null      int32
2   Education   6 non-null      int32
3   Salary      6 non-null      int64
4   Experience  6 non-null      float64
dtypes: float64(1), int32(3), int64(1)
memory usage: 296.0 bytes
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Gender      6 non-null      category
1   Department  6 non-null      int32
2   Education   6 non-null      int32
3   Salary      6 non-null      int64
4   Experience  6 non-null      float64
dtypes: category(1), float64(1), int32(2), int64(1)
memory usage: 402.0 bytes
```

```
In [13]: 1 encoded_data['Gender']=encoded_data['Gender'].astype('category')
2 encoded_data['Department']=encoded_data['Department'].astype('category')
3 encoded_data['Education']=encoded_data['Education'].astype('category')
4 encoded_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Gender      6 non-null      category
1   Department  6 non-null      category
2   Education   6 non-null      category
3   Salary      6 non-null      int64
4   Experience  6 non-null      float64
dtypes: category(3), float64(1), int64(1)
memory usage: 630.0 bytes
```

```
In [20]: 1 #One-Hot Encoding Use pd.get_dummies() to create binary columns.
2
3 onehot_encoded_data = pd.get_dummies(data, columns=['Gender', 'Department',
4
5 onehot_encoded_data
6
```

```
Out[20]:
```

	Salary	Experience	Gender_Female	Gender_Male	Department_HR	Department_IT	Department_S
0	50000	1.5	0	1	1	0	
1	60000	3.0	1	0	0	1	
2	55000	4.2	1	0	0	0	
3	52000	2.3	0	1	1	0	
4	58000	3.7	0	1	0	1	
5	53000	2.9	1	0	0	0	

```
In [8]: 1 from sklearn.preprocessing import LabelBinarizer
2
3 lb = LabelBinarizer()
4 en_data=data.copy()
5 # Gender (Binary)
6 en_data['Gender'] = lb.fit_transform(data['Gender'])
7 en_data['Department'] = lb.fit_transform(data['Department'])
8 en_data['Education'] = lb.fit_transform(data['Education'])
9
10 en_data
11
```

```
Out[8]:
```

	Gender	Department	Education	Salary	Experience
0	1	1	1	50000	1.5
1	0	0	0	60000	3.0
2	0	0	0	55000	4.2
3	1	1	0	52000	2.3
4	1	0	0	58000	3.7
5	0	0	1	53000	2.9

```
In [9]: 1 # Department (Multiclass)
2 dept_binarized = pd.DataFrame(LabelBinarizer().fit_transform(data['Department',
3                                     columns=['Dept_HR', 'Dept_IT', 'Dept_Sales']))
4 dept_binarized
```

```
Out[9]:
```

	Dept_HR	Dept_IT	Dept_Sales
0	1	0	0
1	0	1	0
2	0	0	1
3	1	0	0
4	0	1	0
5	0	0	1

```
In [10]: 1 # Education (Multiclass)
2 edu_binarized = pd.DataFrame(LabelBinarizer().fit_transform(data['Education',
3                                     columns=['Edu_Bachelor', 'Edu_Master', 'Edu_PhD']))
4
5 edu_binarized
```

```
Out[10]:
```

	Edu_Bachelor	Edu_Master	Edu_PhD
0	1	0	0
1	0	1	0
2	0	0	1
3	0	1	0
4	0	0	1
5	1	0	0

Similarity Measures

```
In [17]: 1 import numpy as np
2
3 # Two sample vectors
4 A = np.array([3.7, 2.90, 9.87])
5 B = np.array([5.0, 6.5, 8.6])
6 c=np.array([3.7, 0, 0])
7
```

1. Euclidean Distance (L2 norm)

Straight-line distance between points.

used: Good for continuous features, geometric closeness.

```
In [19]: 1 from scipy.spatial.distance import euclidean
2
3 print("Euclidean Distance:", euclidean(A, B,c))
4
```

Euclidean Distance: 2.5005999280172744

2. Manhattan Distance (L1 norm)

Sum of absolute differences.

Used in urban planning or anytime "right-angle" travel is needed.

```
In [21]: 1 from scipy.spatial.distance import cityblock
2
3 print("Manhattan Distance:", cityblock(A, B))
4
```

Manhattan Distance: 6.17

3. Cosine Similarity

Measures angle between vectors (not magnitude).

Ideal for text similarity, TF-IDF, embeddings.

```
In [23]: 1 from sklearn.metrics.pairwise import cosine_similarity
2
3 cos_sim = cosine_similarity([A], [B])
4 print("Cosine Similarity:", cos_sim[0][0])
5
```

Cosine Similarity: 0.9408871372284504

4. Jaccard Similarity

Works with sets or binary data.

Great for recommendations (e.g., common likes), binary vectors.

```
In [24]: 1 A_set = set([1, 2, 3])
2 B_set = set([2, 3, 4])
3
4 jaccard = len(A_set & B_set) / len(A_set | B_set)
5 print("Jaccard Similarity:", jaccard)
6
```

Jaccard Similarity: 0.5

5. Hamming Distance

Number of positions with different values (binary or strings).

Best for comparing binary vectors or DNA sequences.

```
In [34]: 1 from scipy.spatial.distance import hamming
          2
          3 A_bin = np.array([1, 0, 1, 1])
          4 B_bin = np.array([1, 1, 0, 1])
          5
          6 print("Hamming Distance:", hamming(A_bin, B_bin))
          7
```

Hamming Distance: 0.5

1. Cosine Similarity on Text

Cosine Similarity works best when we vectorize the text (e.g., using TF-IDF or CountVectorizer).

Using TfidfVectorizer

```
In [28]: 1 text1 = "I love machine learning"
          2 text2 = "I enjoy learning machines"
          3
```

```
In [29]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.metrics.pairwise import cosine_similarity
3
4 # Sample texts
5 texts = [text1, text2]
6
7 # Vectorize
8 vectorizer = TfidfVectorizer()
9 tfidf_matrix = vectorizer.fit_transform(texts)
10
11 # Show feature names (unique words)
12 print(vectorizer.get_feature_names_out())
13
14 # Show the TF-IDF matrix as an array
15 tf=pd.DataFrame(tfidf_matrix.toarray(),columns=vectorizer.get_feature_names_
16 tf
```

```
['enjoy' 'learning' 'love' 'machine' 'machines']
```

```
Out[29]:
```

	enjoy	learning	love	machine	machines
0	0.000000	0.449436	0.631667	0.631667	0.000000
1	0.631667	0.449436	0.000000	0.000000	0.631667

```
In [30]: 1 # Compute Cosine Similarity
2 cos_sim = cosine_similarity(tfidf_matrix[0:1], tfidf_matrix[1:2])
3
4 print("Cosine Similarity:", cos_sim[0][0])
5
```

```
Cosine Similarity: 0.20199309249791833
```

```
In [ ]: 1
```

```
In [ ]: 1
```

2. Hamming Distance on Text

Hamming Distance only works on equal-length binary strings or character sequences.

So, we'll:

Convert each sentence to a fixed-length binary or character sequence

Pad if needed

Convert to character sequence (basic example)

```
In [31]: 1 from scipy.spatial.distance import hamming
2
3 # Shortened equal-length strings
4 text1 = "machine"
5 text2 = "mashing"
6
7 # Ensure equal length
8 min_len = min(len(text1), len(text2))
9 text1 = text1[:min_len]
10 text2 = text2[:min_len]
11
12 # Convert to list of characters
13 list1 = list(text1)
14 list2 = list(text2)
15
16 # Compute Hamming Distance
17 hamm_dist = hamming(list1, list2)
18 print("Hamming Distance:", hamm_dist)
19
```

Hamming Distance: 0.2857142857142857

```
In [ ]: 1 # # i.e., 2 out of 7 characters are different
2 # Value is normalized (0 = same, 1 = completely different)
```