# Week-2:Plagiarism Detection using A* Algorithm

Chetankumar kamani(20251603005), Sakariya Devraj (20251603006), Divyesh Dodiya (20251603007),
M.Tech Sem-1 IIIT Vadodara

## I. PROBLEM STATEMENT

Given two text documents, the task is to align their sentences and detect plagiarism using the A* search algorithm. The alignment should minimize the edit distance (or maximize similarity) between corresponding sentences.

## II. PYTHON CODE

```python
import sys, re
from heapq import heappush, heappop

# --- Removing Punctutation Marks ---

def normalize_the_text(s: str) -> str:
    s = s.lower()
    s = re.sub(r"[^\w\s\.!\?]", " ", s)
    return re.sub(r"\s+", " ", s).strip()

def tokenize(s: str):
    s = re.sub(r"\s*([\.!\?])\s*", r"\1 ", s)
    parts = re.split(r"[\.!\?]\s+", s)
    return [t.strip() for t in parts if t.strip()]

def words(s: str):
    return [t for t in re.split(r"\W+", s) if t]

def compute_edit_distance(a: str, b: str) -> int:
    A = words(a)
    B = words(b)
    m = len(A)
    n = len(B)
    if m == 0:
        return n
    if n == 0:
        return m
    dp = list(range(n + 1))
    #lavenshtein distance calculation
    for i in range(1, m + 1):
        prev = dp[0]
        dp[0] = i
        for j in range(1, n + 1):
            tmp = dp[j]
            if A[i-1] == B[j-1]:
                cost = 0
            else:
                cost = 1
            # The following two lines were
                incorrectly indented.
            # They must be inside the inner loop
                to work correctly.
            dp[j] = min(dp[j] + 1, dp[j-1] + 1,
                prev + cost)
            prev = tmp
    return dp[-1]


def ned(text_a: str, text_b: str) -> float:
    num_words_a = len(words(text_a))
    num_words_b = len(words(text_b))
    max_length = max(num_words_a, num_words_b)
    if max_length == 0:
        return 0.0
    raw_distance = compute_edit_distance(text_a,
        text_b)
    return raw_distance / max_length


# --- A* function
def a_star_function(SA, SB, skip_penalty=3.0):
    m, n = len(SA), len(SB)
    def h(i, j):
        return abs((m - i) - (n - j)) *
            skip_penalty

    openq = []
    heappush(openq, (h(0,0), 0.0, 0, 0))
    best = {(0,0): 0.0}
    prev = {}

    while openq:
        f, g, i, j = heappop(openq)
        if (i, j) == (m, n):
            steps = []
            cur = (i, j)
            while cur != (0,0):
                p, op = prev[cur]
                steps.append(op)
                cur = p
            return g, list(reversed(steps))

        # SKIP_A
        if i < m:
            ng, s = g + skip_penalty, (i+1, j)
            if ng < best.get(s, 1e18):
                best[s] = ng; prev[s] = ((i,j), ("
                    SKIP_A", i, -1, skip_penalty))
                heappush(openq, (ng + h(*s), ng, *
                    s))
        # SKIP_B
        if j < n:
            ng, s = g + skip_penalty, (i, j+1)
            if ng < best.get(s, 1e18):
                best[s] = ng; prev[s] = ((i,j), ("
                    SKIP_B", -1, j, skip_penalty))
                heappush(openq, (ng + h(*s), ng, *
                    s))
        # ALIGN
        if i < m and j < n:
            c = compute_edit_distance(SA[i], SB[j
                ])
            ng, s = g + c, (i+1, j+1)
            if ng < best.get(s, 1e18):
                best[s] = ng; prev[s] = ((i,j), ("
                    ALIGN", i, j, float(c)))
                heappush(openq, (ng + h(*s), ng, *
                    s))

    return float("inf"), []

def compare_document(fileA, fileB, t_value,
    sp_value):

    SA = tokenize(normalize_the_text(fileA))
    SB = tokenize(normalize_the_text(fileB))
```

```python
    aligned=[]
    total, steps = a_star_function(SA, SB,
        sp_value)

    for s in steps:
        if s[0]=="ALIGN":
            aligned.append(s)

    plag = []

    print("------------Result--------------")
    print("number of sentence in the file A:", len
        (SA))
    print("number of sentence in the file B:", len
        (SB))
    print("total_path_cost:", total)

    print("\nALIGNMENT")
    for op, i, j, c in steps:
        if op == "ALIGN":
            ne = round(ned(SA[i], SB[j]), 3)
            flag = (ne <= t_value)
            if flag: plag.append(1)
            status = "PLAGIARIZED" if flag else "
                ORIGINAL"
            print(f"[ALIGN] A[{i}]<->B[{j}] cost={
                int(c)} NED={ne} Status: {status}"
                )
            print("    A:", SA[i])
            print("    B:", SB[j])
        elif op == "SKIP_A":
            print(f"[SKIP_A] A[{i}] {SA[i]}")
        else:
            print(f"[SKIP_B] B[{j}] {SB[j]}")

    pr = (sum(plag)/len(aligned)) if aligned else
        0.0
    plagiarism_percentage = round(pr * 100, 2)

    # The documents are "identical" at the
        percentage of aligned content that is not
        plagiarized.
    originality_percentage = round((1.0 - pr) *
        100, 2)

    print("\n----------------SUMMARY
        ------------------")
    print(f"Total number of plagiarized pairs
        found: {sum(plag)} out of {len(aligned)}
        aligned sentences.")
    print(f"The documents are {
        plagiarism_percentage}% plagiarized.")
    print(f"The documents are {
        originality_percentage}% original.")
    print(f"Based on the analysis, the two
        documents are {originality_percentage}%
        identical.")


#Input two files while executing the program at
    the command line.

if len(sys.argv) < 3:
    print("Expecting Files. You have not given
        files.")
    sys.exit(1)

#path of file 1
A_path = sys.argv[1]

#path of file 2
B_path = sys.argv[2]
```

```python
#Getting the threshold value from the user.
t_value = float(input("Enter threshold value: "))

#Getting the skip-panelty value from the user
sp_value  = float(input("Enter skip-panelty value:
        ")) # Corrected typo in the prompt

#initially checking the length of the file. If the
     file length

#reading the content from the file A
with open(A_path, "r") as f:
    A = f.read()

#reading the content from the file B
with open(B_path, "r") as f:
    B = f.read()

compare_document(A, B, t_value, sp_value)
```

## III. INPUT FILES

### A. Test Case 1: Identical Documents

#### T1docA.txt

I like Artificial Intelligence. I like IIIT
    Vadodara
I like probability & Statistics. I like IIIT
    Vadodara.

#### T1docB.txt

I like Artificial Intelligence.
I like IIIT Vadodara.
I like probability and statistics.

### B. Test Case 2: Slightly Modified Documents

#### T2docA.txt

I like Artificial Intelligence.
I like IIIT Vadodara.
I like probability and statistics.

#### T2docB.txt

I like Artificial Intelligence.
I like IIIT Vadodara.
I like probability and statistics.

### C. Test Case 3: Completely Different Documents

#### T3docA.txt

My name is Chetan Kamani.
I live in Jamnagar.
I work as a lecturer in Government Polytechnic.

#### T3docB.txt

Sorting algorithms arrange data in ascending or
    descending order.
QuickSort uses partitioning and recursion.
Heaps are used for priority queues.

## D. Test Case 4: Partial Overlap
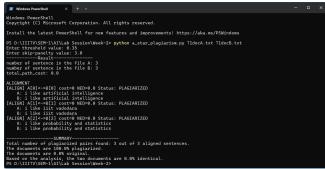
### T4docA.txt

My name is Chetan Kamani.
I like IIIT Vadodara.
I enjoy probability and statistics.
Artificial Intelligence is my favorite subject.
I live in Jamnagar.

### T4docB.txt

I like Artificial Intelligence.
I like IIIT Vadodara.
I study probability and statistics.
I live in Jamnagar.

# IV. SCREENSHOTS OF RESULTS

## A. Test Case 1



## B. Test Case 2



## C. Test Case 3



## D. Test Case 4



# V. CODE AVAILABILITY

The complete source code is available at: GitHub Repository (CS659 – AI Laboratory).

**GitHub Repository:**

https://github.com/ChetanKamani/CS659-AI-Lab