# PLAGIARISM CHECKER REPORT

## Ahmet Hüzeyfe Demir

HTTPS://GİTHUB.COM/KOYUBOY?TAB=REPOSİTORİES

https://www.linkedin.com/in/ahmet-h%C3%BCzeyfe-demir-80515b157/

# INTRODUCTION

A Line or word based algorithms, deep learning based methods etc. detects and reveals a plagiarized content, whether intentional or not.

There are many approaches and methods to detect plagiarized contents. Line or word based algorithms, deep learning based methods etc. Usually, plagiarism checkers work online, which means they search and detect content online.

We developed an algorithm based on row analysis in our project and we made use of the scipy library while doing this. We used Cosine Similarity Function from Scipy library.

Cosine similarity is a metric used to determine how similar the documents are irrespective of their size. Another definition is  Cosine similarity measures the similarity between two vectors of an inner product space.

The function is -> cosine(u, v)

u and v are one dimensional arrays. Function computes the Cosine distance between two arrays. After returns a double value.

All program based on row_similarity(row1, row2) funciton. So we will understand how this works. This function takes two arguments. They are row1 and row2 and  they are in string type. Function take this rows and splitis. After that uses two dictionary to reveal frequency of splitted elements. After that, three things are important to us. Two lists vector1 , vector2 and a set of all_words_set. We used a set to eliminate duplicate elements. We have filled vector1 and vector2 with this set in a loop with a rule. After that cosine similarity function called with vector1 and vector2.

There is an example of this operation:

row1 = "I have an apple"

row2 = "She has an apple and apple juice"

all_words_set = ["I", "have", "an", "apple", "She", "has", "and", "juice"]

vector1 = [1, 1, 1, 1, 0, 0, 0, 0]

vector2 = [0, 0, 1, 2, 1, 1, 1, 1]

1 – spatial.distance.cosine(vector1, vector2)

This operation returns a similarity ratio. We can manually choose a threshold for accepting these two contents as duplicates. We finished explaining what this function does. Let's see what the whole program does.

Other jobs of the program is quite simple. The program and the files to be compared must be in the same directory. It reads files with Python extension, splits them into lines and puts them into some preprocessing. Lastly it sends the lines in pairs to the row_similarity function. Each line in a file is compared with lines in other files.

The number of times it has been detected in other files is written next to each line in the files. At the end, the similarity ratios of the files are printed on the screen and the program ends.

<h1 style="text-align:center; color:red">SAMPLE RUN -1</h1>

**We have 4 python files here:** simple_program_1.py, simple_program_2.py, simple_program_3.py, simple_program_4.py

```python
# simple_program_1.py > BOARD_SIZE
1    BOARD_SIZE = 8
2
3    def under_attack(col, queens):
4        left = right = col
5
6        for r, c in reversed(queens):
7            left, right = left - 1, right + 1
8
9            if c in (left, col, right):
10                return True
11        return False
12
13    def solve(n):
14        if n == 0:
15            return [[]]
16
17        smaller_solutions = solve(n - 1)
18
19        return [solution+[(n,i+1)]
20            for i in range(BOARD_SIZE)
21                for solution in smaller_solutions
22                    if not under_attack(i+1, solution)]
23    for answer in solve(BOARD_SIZE):
24        print (answer)
```

```python
# simple_program_2.py > BOARD_SIZE
1    BOARD_SIZE = 8
2
3    class BailOut(Exception):
4        pass
5
6    def validate(queens):
7        left = right = col = queens[-1]
8        for r in reversed(queens[:-1]):
9            left, right = left-1, right+1
10            if r in (left, col, right):
11                raise BailOut
12
13    def add_queen(queens):
14        for i in range(BOARD_SIZE):
15            test_queens = queens + [i]
16            try:
17                validate(test_queens)
18                if len(test_queens) == BOARD_SIZE:
19                    return test_queens
20                else:
21                    return add_queen(test_queens)
22            except BailOut:
23                pass
24        raise BailOut
25
26    queens = add_queen([])
27    print (queens)
28    print ("\n".join(". "*q + "Q " + ". "*(BOARD_SIZE-q-1) for q in queens))
```

```python
simple_program_3.py > ...
1    import random
2
3    guesses_made = 0
4
5    name = input('Hello! What is your name?\n')
6
7    number = random.randint(1, 20)
8    print ('Well, {0}, I am thinking of a number between 1 and 20.'.format(name))
9
10   while guesses_made < 6:
11
12       guess = int(input('Take a guess: '))
13
14       guesses_made += 1
15
16       if guess < number:
17           print ('Your guess is too low.')
18
19       if guess > number:
20           print ('Your guess is too high.')
21
22       if guess == number:
23           break
24
25   if guess == number:
26       print ('Good job, {0}! You guessed my number in {1} guesses!'.format(name, guesses_made))
27   else:
28       print ('Nope. The number I was thinking of was {0}'.format(number))
```

```python
simple_program_4.py > ...
1    import itertools
2
3    def iter_primes():
4        # an iterator of all numbers between 2 and +infinity
5        numbers = itertools.count(2)
6
7        # generate primes forever
8        while True:
9            # get the first number from the iterator (always a prime)
10           prime = next(numbers)
11           yield prime
12
13           # this code iteratively builds up a chain of
14           # filters...slightly tricky, but ponder it a bit
15           numbers = filter(prime.__rmod__, numbers)
16
17   for p in iter_primes():
18       if p > 1000:
19           break
20       print (p)
```

**Output of the program:**

```
Similarity of simple_program_1.py to simple_program_2.py is 22.22222222222222 %
Similarity of simple_program_1.py to simple_program_3.py is 0.0 %
Similarity of simple_program_1.py to simple_program_4.py is 0.0 %
Similarity of simple_program_2.py to simple_program_1.py is 16.666666666666664 %
Similarity of simple_program_2.py to simple_program_3.py is 4.166666666666666 %
Similarity of simple_program_2.py to simple_program_4.py is 0.0 %
Similarity of simple_program_3.py to simple_program_1.py is 0.0 %
Similarity of simple_program_3.py to simple_program_2.py is 0.0 %
Similarity of simple_program_3.py to simple_program_4.py is 5.555555555555555 %
Similarity of simple_program_4.py to simple_program_1.py is 0.0 %
Similarity of simple_program_4.py to simple_program_2.py is 0.0 %
Similarity of simple_program_4.py to simple_program_3.py is 6.25 %
```

**Python files after operation:**

simple_program_1.py > ...

```python
1   *-- BOARD_SIZE = 8!! Similar with 1 lines. --*
2   def under_attack(col, queens):
3   *--     left = right = col!! Similar with 1 lines. --*
4       for r, c in reversed(queens):
5           left, right = left - 1, right + 1
6   *--         if c in (left, col, right):!! Similar with 1 lines. --*
7               return True
8       return False
9   def solve(n):
10      if n == 0:
11          return [[]]
12      smaller_solutions = solve(n - 1)
13      return [solution+[(n,i+1)]
14  *--         for i in range(BOARD_SIZE)!! Similar with 1 lines. --*
15              for solution in smaller_solutions
16                  if not under_attack(i+1, solution)]
17  for answer in solve(BOARD_SIZE):
18      print (answer)
```

simple_program_2.py > ...

```python
1   *-- BOARD_SIZE = 8!! Similar with 1 lines. --*
2   class BailOut(Exception):
3       pass
4   def validate(queens):
5   *--     left = right = col = queens[-1]!! Similar with 1 lines. --*
6       for r in reversed(queens[:-1]):
7           left, right = left-1, right+1
8   *--         if r in (left, col, right):!! Similar with 1 lines. --*
9               raise BailOut
10  def add_queen(queens):
11  *--     for i in range(BOARD_SIZE):!! Similar with 1 lines. --*
12          test_queens = queens + [i]
13          try:
14              validate(test_queens)
15              if len(test_queens) == BOARD_SIZE:
16                  return test_queens
17  *--             else:!! Similar with 1 lines. --*
18                  return add_queen(test_queens)
19          except BailOut:
20              pass
21      raise BailOut
22  queens = add_queen([])
23  print (queens)
24  print ("\n".join(". "*q + "Q " + ". "*(BOARD_SIZE-q-1) for q in queens))
```

```python
# simple_program_3.py > ...
1   import random
2   guesses_made = 0
3   name = input('Hello! What is your name?\n')
4   number = random.randint(1, 20)
5   print ('Well, {0}, I am thinking of a number between 1 and 20.'.format(name))
6   while guesses_made < 6:
7       guess = int(input('Take a guess: '))
8       guesses_made += 1
9       if guess < number:
10          print ('Your guess is too low.')
11      if guess > number:
12          print ('Your guess is too high.')
13      if guess == number:
14 *--         break!! Similar with 1 lines. --*
15  if guess == number:
16      print ('Good job, {0}! You guessed my number in {1} guesses!'.format(name, guesses_made))
17  else:
18      print ('Nope. The number I was thinking of was {0}'.format(number))
```

```python
# simple_program_4.py > ...
1    import itertools
2    def iter_primes():
3        # an iterator of all numbers between 2 and +infinity
4        numbers = itertools.count(2)
5        # generate primes forever
6        while True:
7            # get the first number from the iterator (always a prime)
8            prime = next(numbers)
9            yield prime
10           # this code iteratively builds up a chain of
11           # filters...slightly tricky, but ponder it a bit
12           numbers = filter(prime.__rmod__, numbers)
13   for p in iter_primes():
14       if p > 1000:
15 *--             break!! Similar with 1 lines. --*
16       print (p)
```

# SAMPLE RUN -2

**We have 5 python files here:** test_1.py, test_2.py, test_3.py, test_4.py, test_5.py
These are randomly generated python codes from http://www.4geeks.de/cgi-bin/webgen.py

## Output of the program:

```
Similarity of test_1.py to test_2.py is 26.751592356687897 %
Similarity of test_1.py to test_3.py is 33.12101910828025 %
Similarity of test_1.py to test_4.py is 31.210191082802545 %
Similarity of test_1.py to test_5.py is 30.573248407643312 %
Similarity of test_2.py to test_1.py is 37.93103448275862 %
Similarity of test_2.py to test_3.py is 33.62068965517241 %
Similarity of test_2.py to test_4.py is 32.758620689655174 %
Similarity of test_2.py to test_5.py is 29.310344827586203 %
Similarity of test_3.py to test_1.py is 35.55555555555556 %
Similarity of test_3.py to test_2.py is 30.0 %
Similarity of test_3.py to test_4.py is 34.44444444444444 %
Similarity of test_3.py to test_5.py is 28.888888888888886 %
Similarity of test_4.py to test_1.py is 46.05263157894737 %
Similarity of test_4.py to test_2.py is 31.57894736842105 %
Similarity of test_4.py to test_3.py is 39.473684210526315 %
Similarity of test_4.py to test_5.py is 35.526315789473685 %
Similarity of test_5.py to test_1.py is 34.66666666666667 %
Similarity of test_5.py to test_2.py is 26.666666666666668 %
Similarity of test_5.py to test_3.py is 30.666666666666664 %
Similarity of test_5.py to test_4.py is 33.33333333333333 %
The program is over!!
```

## A part of test_1.py files after operation:

```
71    def func8():
72 v *--     closure = [6]!! Similar with 3 lines. --*
73 v     def func7(arg46, arg47):
74           closure[0] += func9(arg46, arg47)
75    *--        return closure[0]!! Similar with 2 lines. --*
76    *--     func = func7!! Similar with 2 lines. --*
77    *--     return func!! Similar with 2 lines. --*
78    var48 = func8()
79 v def func6(arg26, arg27):
80        var28 = 414368693 | (-1236223149 ^ -1634524315) | arg27
81        var29 = arg26 - (arg26 + 1573171526) | var28
82 v *--     var30 = var29 + var29!! Similar with 2 lines. --*
83        var31 = ((-163403563 ^ arg27) & -1207192254) ^ -1703484484
84 v *--     var32 = ((var28 | var30) ^ var28) | var29!! Similar with 1 lines. --*
85        var33 = var32 - arg27
86        var34 = var29 & (var31 ^ var32) & var30
87        var35 = -65 - 872895955
```

## CONCLUSION

This program looks at every line and then every word. This is a little bad approach. Algorithm can be developed with recognize the structs like functions, loops, switch-cases etc. After that program can compares that structs with each other. This will make better accuracy.  Because some code lines can be in all the programs and this lines increase our plagiarism ratio. Therefore program should look code files by part based not only line based.

Also in our program threshold is setting manually by programmer. This threshold ratio should be setted nicely. A deep learning model or AI based programs can be used to determine the threshold.

Lastly if this algorithm improves it can be detailed checker. Because we can say that this plagiarism checker algorithm is variable-free. It means is that revealing plagiarism is not affected from changing function names, if program improved to part based algorithm.