

## Implement LSH

### Group members : ¶

1. Mintra Sojiphan 6220422057 DS6
2. Kavin Singhakhet 6310422040 DS7
3. Porawoot Buranadilok 6220422064 DS6

We implement LSH using two packages with two different datasets

1. Package: SnaPy // Datasets: A Corpus of Plagiarised Short Answers

References: <https://github.com/justinnbt/SnaPy> (<https://github.com/justinnbt/SnaPy>)

[https://ir.shef.ac.uk/cloughie/resources/plagiarism\\_corpus.html](https://ir.shef.ac.uk/cloughie/resources/plagiarism_corpus.html)  
([https://ir.shef.ac.uk/cloughie/resources/plagiarism\\_corpus.html](https://ir.shef.ac.uk/cloughie/resources/plagiarism_corpus.html))

2. Package: datasketch // Dataset: News headlines

References: <http://ekzhu.com/datasketch/lshforest.html>  
(<http://ekzhu.com/datasketch/lshforest.html>)

1. Package: SnaPy // Datasets: A Corpus of Plagiarised Short Answers

**We will implement LSH on plagiarised answers of five different questions namely question A, B, C, D, and E to identify near duplicate answers of each question using Jaccard similarity threshold (s) = 0.5**

### Step 1: Install and import packages

In [ ]:

```
pip install snapy
```

Collecting snapy

Downloading <https://files.pythonhosted.org/packages/68/59/cdf153f7a391593060a3fcd06f13f8127dc79f675d99ef576b69f49365a0/snapy-1.0.2-py3-none-any.whl>

Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from snapy) (1.19.5)

Installing collected packages: snapy

Successfully installed snapy-1.0.2

In [ ]:

```
pip install mmh3
```

Collecting mmh3

Downloading <https://files.pythonhosted.org/packages/fa/7e/3ddcab0a9fcea034212c02eb411433db9330e34d626360b97333368b4052/mmh3-2.5.1.tar.gz>

Building wheels for collected packages: mmh3

Building wheel for mmh3 (setup.py) ... done

Created wheel for mmh3: filename=mmh3-2.5.1-cp36-cp36m-linux\_x86\_64.whl  
size=37849 sha256=3f28dc26ac98dc67078e2b1835f83f8d0ceb68783fbb8283be366793d5b16299

Stored in directory: /root/.cache/pip/wheels/38/b4/ea/6e4e321c625d3320c0c496bf4088371546d8fce5f1dd71b219

Successfully built mmh3

Installing collected packages: mmh3

Successfully installed mmh3-2.5.1

In [ ]:

```
from snapy import MinHash, LSH
from google.colab import drive
import matplotlib.pyplot as plt
import networkx as nx
```

## Step 2: Connect Google Colab with Google drive

In [ ]:

```
drive.mount("/content/drive")
```

Mounted at /content/drive

## Step 3: Create function to perform LSH

In [ ]:

```

def LSH_task(task):

    #import data
    files=[]
    x='/content/drive/MyDrive/corpus/'+task+'.txt'
    files.append(x)
    for i in range(1,20):
        z='/content/drive/MyDrive/corpus/'+task+ ' ('+str(i)+').txt'
        files.append(z)

    docs=[]
    for file in files:
        file = open(file)
        text = file.read()
        docs.append(text)
    file.close()

    print(docs)
    print()
    print('#task: ',len(docs))
    labels=[]
    labels.append(task)
    for i in range(1,20):
        a=task+' ('+str(i)+')'
        labels.append(a)

    seed =3

    #Create MinHash object.
    minhash = MinHash(docs, n_gram=9, permutations=100, hash_bits=64, seed=3)
    print('Signatures metric: ',len(minhash.signatures),'x',len(minhash.signatures[0]))
    print('#permutations used to create signatures:',minhash.permutations)
    #print('Minhash Signatures for each text:')
    #for i in minhash.signatures:
        #print(i)

    # Create LSH model.
    lsh = LSH(minhash, labels, no_of_bands=50)

    # Query to find near duplicates for each doc.
    # print()
    #for i in labels:
        #print('Near duplicate for answer',i,':',lsh.query(i, min_jaccard=0.1))

    # Check contents of documents.
    #print(lsh.contains())

    # Return adjacency list for all similar texts.
    adjacency_list = lsh.adjacency_list(min_jaccard=0.1)
    #print()
    #print('adjacency_lists: ',adjacency_list)
    #print()
    # Returns edge list for use creating a weighted graph.
    edge_list = lsh.edge_list(min_jaccard=0.1, jaccard_weighted=True)
    #print('edge lists: ',edge_list)
    #print()

    # Create Undirected weighted graph.
    fig=plt.figure(figsize =(10,6))

```

```

fig.set_facecolor("#181818")
title="Near duplicate answer of question "+task
fig.suptitle(title,color= '#cccccc',fontsize=18)
G = nx.Graph()
for i in edge_list:
    G.add_edge(i[0],i[1],weight=i[2])
e1=[(u,v) for (u,v,d) in G.edges(data=True) if d['weight'] > 0.5]
e2=[(u,v) for (u,v,d) in G.edges(data=True) if d['weight'] <= 0.5]
pos=nx.spring_layout(G)

# nodes
nx.draw_networkx_nodes(G,pos,node_size=900,node_color='yellow')

# edges
edge1=nx.draw_networkx_edges(G,pos,edgelist=e1,width=4,edge_color='red')
edge2=nx.draw_networkx_edges(G,pos,edgelist=e2,width=1,edge_color='red',style='dashed')

# Labels
nx.draw_networkx_labels(G,pos,font_family='sans-serif',font_size=10,font_color='#000000',font_weight=150)
keys=[(i[0],i[1]) for i in edge_list]
values= [(i[2]) for i in edge_list]
edge_labels = dict(zip(keys, values))
nx.draw_networkx_edge_labels(G,pos,edge_labels= edge_labels,font_color='red')
fig.set_facecolor("#181818")
plt.axis('off')
fig.legend((edge1, edge2), ('sim(i,j) > 0.5', 'sim(i,j) <= 0.5'),loc=1,fontsize=12)
plt.savefig("LSH_graph.png")
plt.show()

```

**Step 4: Perform LSH on set of 20 plagiarised answers of the following questions to identify near duplicate answers using Jaccard similarity threshold (s) = 0.5**

**Question A: The result illustrates that there are 3 pairs of answers which can be identified as near duplicate answers as shown below:**

- 1. Answer A(1) and A(18) with Jaccard similarity of 0.9**
- 2. Answer A(1) and A(6) with Jaccard similarity of 0.74**
- 3. Answer A(6) and A(18) with Jaccard similarity of 0.7**

In [ ]:

```
LSH_task('A')
```

['Object oriented programming is a style of programming that supports encapsulation, inheritance, and polymorphism. Inheritance means derived a new class from the base class. We can also say there are parents class and child classes in inheritance. Inheritance was firstly derived in 1967. The child class has all the features of parents class or we can say the base class more over it may also include some additional features. Inheritance is used for modification and implementation new features in computer programming language. It is possible that child class has all the attributes of parents class but it is not possible that all the attributes of child class must have in base class or parent class. In categorization in computer language also inheritance is a useful tool. categorization define as a powerful feature. it has been also used in generalisation and in human learning. In some areas less information need to be stored. Generalisation also some time known as inheritance. The main reason behind this is a hierarchy structure of objects and classes. We can understand this mechanism by some examples: like fruit is aq main class and mangoes apple ,orange is child classs of the main class. So obviously inherit all the properties of fruit class. In object-oriented programming, inheritance is a way to form new classes (instances of which are called objects) using classes that have already been defined. The inheritance concept was invented in 1967 for Simula. The new classes, known as derived classes, take over (or inherit) attributes and behavior of the pre-existing classes, which are referred to as base classes (or ancestor classes). It is intended to help reuse existing code with little or no modification. Inheritance provides the support for representation by categorization in computer languages. Categorization is a powerful mechanism number of information processing, crucial to human learning by means of generalization (what is known about specific entities is applied to a wider group given a belongs relation can be established) and cognitive economy (less information needs to be stored about each specific entity, only its particularities). Inheritance is also sometimes called generalization, because the is-a relationships represent a hierarchy between classes of objects. For instance, a "fruit" is a generalization of "apple", "orange", "mango" and many others. One can consider fruit to be an abstraction of apple, orange, etc. Conversely, since apples are fruit (i.e., an apple is-a fruit), apples may naturally inherit all the properties common to all fruit, such as being a fleshy container for the seed of a plant. An advantage of inheritance is that modules with sufficiently similar interfaces can share a lot of code, reducing the complexity of the program. Inheritance therefore has another view, a dual, called polymorphism, which describes many pieces of code being controlled by shared control code. Inheritance is typically accomplished either by overriding (replacing) one or more methods exposed by ancestor, or by adding new methods to those exposed by an ancestor. Complex inheritance, or inheritance used within a design that is not sufficiently mature, may lead to the Yoyo problem. Inheritance is a basic concept of Object-Oriented Programming where the basic idea is to create new classes that add extra detail to existing classes. This is done by allowing the new classes to reuse the methods and variables of the existing classes and new methods and classes are added to specialise the new class. Inheritance models the "is-kind-of" relationship between entities (or objects), for example, postgraduates and undergraduates are both kinds of student. This kind of relationship can be visualised as a tree structure, where 'student' would be the more general root node and both 'postgraduate' and 'undergraduate' would be more specialised extensions of the 'student' node (or the child nodes). In this relationship 'student' would be known as the superclass or parent class whereas, 'postgraduate' would be known as the subclass or child class because the 'postgraduate' class extends the 'student' class. Inheritance can occur on several layers, where if visualised would display a larger tree structure. For example, we could further extend the 'postgraduate' node by adding two extra extended classes to it called, 'MSc Student' and 'PhD Student' as both these types of student are kinds of postgrad

uate student. This would mean that both the 'MSc\nStudent' and 'PhD Student' classes would inherit methods and variables\nfrom both the 'postgraduate' and 'student classes'.\n', "Inheritance is a basic concept in object oriented programming. It models the reuse of existing class code in new classes - the "is a kind of" relationship.\n\nFor example, a house is a kind of building; similarly, an office block is a kind of building. Both house and office block will inherit certain characteristics from buildings, but also have their own personal characteristics - a house may have a number of occupants, whereas an office block will have a number of offices. However, these personal characteristics don't apply to all types of buildings.\n\nIn this example, the building would be considered the superclass - it contains general characteristics for other objects to inherit - and the house and office block are both subclasses - they are specific types and specialise the characteristics of the superclass.\n\nJava allows object inheritance. When one class inherits from another class, all the public variables and methods are available to the subclass.\n\npublic class Shape {\n\n private Color colour;\n\n public void setColour(Color newColour){\n\n colour = newColour;\n\n }\n\n}\n\npublic class Circle extends Shape {\n\n private int radius;\n\n public void setRadius(int newRadius){\n\n radius = newRadius;\n\n }\n\n}\n\nIn this example, the Circle class is a subclass of the Shape class. The Shape class provides a public setColour method, which will be available to the Circle class and other subclasses of Shape. However, the private variable colour (as defined in the Shape class) will not be available for direct manipulation by the Circle class because it is not inherited. The Circle class specialises the Shape class, which means that setRadius is available to the Circle class and all subclasses of Circle, but it isn't available to the superclass Shape. \n", 'inheritance in object oriented programming is where a new class is formed using classes which have already been defined. These classes have some of the behavior and attributes which were existent in the classes that it inherited from. The purpose of inheritance in object oriented programming is to minimize the reuse of existing code without modification.\n\nInheritance allows classes to be categorized, similar to the way humans categorize. It also provides a way to generalize due to the "is a" relationship between classes. For example a "cow" is a generalization of "animal" similarly so are "pigs" & "cheaters". Defining classes in this way, allows us to define attributes and behaviours which are common to all animals in one class, so cheaters would naturally inherit properties common to all animals.\n\nThe advantage of inheritance is that classes which would otherwise have a lot of similar code, can instead share the same code, thus reducing the complexity of the program. Inheritance, therefore, can also be referred to as polymorphism which is where many pieces of code are controlled by shared control code.\n\nInheritance can be accomplished by overriding methods in its ancestor, or by adding new methods. \n', 'Inheritance in object oriented programming is a way to form new classes using classes that have already been defined. The new classes, known as derived classes, inherit attributes and behaviour of the existing classes, which are referred to as base classes. With little or no modification, it is intended to help reuse existing code. It is typically accomplished either by overriding one or more methods exposed by ancestor, or by adding new methods to those exposed by an ancestor.\n\nInheritance is also sometimes called generalization, because there is a relationship that represents a hierarchy between classes of objects. A 'fruit', for instance, is a generalization of "orange", "mango", "apples" and many others. One can consider fruit to be an abstraction of apple, orange, etc. Since apples are fruit (i.e., an apple is-a fruit), conversely apples may naturally inherit all the properties common to all fruit, such as being a fleshy container for the seed of a plant.\n\nAn advantage of inheritance is that modules with sufficiently similar interfaces can share a lot of code reducing the complexity of the program. \n', 'In object-oriented programming, inheritance is a way to form new classes (instances of which are called objects) using classes that have already been defined. The inheritance

concept was invented in 1967 for Simula. The new classes, known as derived classes, take over (or inherit) attribute and behaviour of the pre-existing classes, which are referred to as base classes (or ancestor classes). It is intended to help reuse existing code with little or no modification. Inheritance provides the support for representation by categorization in computer languages. Categorization is a powerful mechanism number of information processing, crucial to human learning by means of generalization (what is known about specific entities is applied to a wider group given a belongs relation can be established) and cognitive economy (less information needs to be stored about each specific entity, only its particularities). Inheritance is also sometimes called generalization, because the is-a relationships represent a hierarchy between classes of objects. For instance, a "fruit" is a generalization of "apple", "orange", "mango" and many others. One can consider fruit to be an abstraction of apple, orange, etc. Conversely, since apples are fruit (i.e., an apple is-a fruit), apples may naturally inherit all the properties common to all fruit, such as being a fleshy container for the seed of a plant. An advantage of inheritance is that modules with sufficiently similar interfaces can share a lot of code, reducing the complexity of the program. Inheritance therefore has another view, a dual, called polymorphism, which describes many pieces of code being controlled by shared control code. Inheritance is typically accomplished either by overriding (replacing) one or more methods exposed by ancestor, or by adding new methods to those exposed by an ancestor.

\n', 'In object oriented programming, objects are grouped together into classes according to their type, structure and the functions that can be performed on them. Inheritance is a process in object oriented programming in which objects acquire (or inherit) the properties of objects of another class. It is therefore used to create relationships between one object and another. Each class groups together objects of a similar type, with similar properties. New classes can be formed by this process whose objects will have properties of both the classes from which this new class is formed. A superclass has all of the properties of the subclasses below it. At the same time subclasses are each distinctive from each other but related via the superclass. Subclasses are said to 'extend' superclasses. Due to these relationships, object oriented programmes tend to be easier to modify since they do not need to be changed when a new object, with different properties is added. Instead, a new object is made to inherit properties of objects which already exist. Inheritance can be divided into two main processes: single inheritance and multiple inheritance. Single inheritance means that the class can only inherit from one other class, whereas multiple inheritance allows for inheritance from several classes.

\n', 'Inheritance is one of the basic concepts of Object Oriented Programming. It's objective is to add more detail to pre-existing classes whilst still allowing the methods and variables of these classes to be reused. The easiest way to look at inheritance is as an "...is a kind of" relationship. For example, a guitar is a kind of string instrument, electric, acoustic and steel stringed are all types of guitar.

\n

The further down an inheritance tree you get, the more specific the classes become. An example here would be books. Books generally fall into two categories, fiction and non-fiction. Each of these can then be subdivided into more groups. Fiction for example can be split into fantasy, horror, romance and many more. Non-fiction splits the same way into other topics such as history, geography, cooking etc. History of course can be subdivided into time periods like the Romans, the Elizabethans, the World Wars and so on.

\n', 'Inheritance is a method of forming new classes using predefined classes. The new classes are called derived classes and they inherit the behaviours and attributes of the base classes. It was intended to allow existing code to be used again with minimal or no alteration. It also offers support for representation by categorization in computer languages; this is a powerful mechanism of information processing, vital to human learning by means of generalization and cognitive economy. Inheritance is occasionally referred to as generalization due to the fact that is-a rel



relationships represent a hierarchy between classes of objects. Inheritance has the advantage of reducing the complexity of a program since modules with very similar interfaces can share lots of code. Due to this, inheritance has another view called polymorphism, where many sections of code are being controlled by some shared control code. Inheritance is normally achieved by overriding one or more methods exposed by ancestor, or by creating new methods on top of those exposed by an ancestor. Inheritance has a variety of uses. Each different use focuses on different properties, for example the external behaviour of objects, internal structure of an object, inheritance hierarchy structure, or software engineering properties of inheritance. Occasionally it is advantageous to differentiate between these uses, as it is not necessarily noticeable from context. \n', 'Inheritance allows programs developed in an Object Orientated language to reuse code without having it replicated unnecessarily elsewhere within the program.\n\nTo achieve this, the programmer has to note generalisations and similarities about various aspects of the program. \n\nFor example, a program could exist to model different forms of transport. At first glance, a car and a train may not have much in common. But abstractly, both will have a speed at which they are travelling, a direction, and a current position.\nMethods utilising this data can be specified high up in the inheritance hierarchy, for example in a 'Transport' class. For example you could have a method which works out the new position of a train after travelling x minutes in direction y. Likewise, you might want to be able to find out the same information for an object of the type car. \nInheritance means that if such a method was defined in the superclass of the train and car classes, any car or train object can utilise it. \n\nThe train and car subclasses are said to 'extend' the Transport class, as they will have additional characteristics which they don't share. E.g. passenger capacity would be a class variable of both car and train (but have different values), and a train may have methods along the lines of 'is toilet engaged'. \nIf you then wanted to add additional forms of transport, such as an aeroplane, you may wish for that also to have a 'toilet engaged' function. Then you could have an extended hierarchy, where a Mass Transport class extends the Transport class. Under which you'd have a train and aeroplane, which would inherit characteristics from both super classes.', 'Inheritance is an important feature in object orientated programming. This is because it allows new classes to be made that extend previous classes and to go into more detail. \n\nThis is carried out by allowing the new class to reuse the existing class methods and variables, whilst also creating class specific methods and variables. This means that the new class, the subclass, is a more specialised version of the original, or superclass.\n\nBecause of this it means that the subclass can use all the public methods and variables from the superclass; however any private methods or variables are still private. \n\nAlso it should be noted that a class can only extend one class, e.g. can only be a subclass to one superclass. However a superclass can have more than one subclass and a class can both be a subclass and a superclass. If this occurs then all of the non-private methods and variables can be used by the most specialised class.\n\nThis means that inheritance is used when types have common factors and these would be put into the superclass. Then the subclass/es then extend these to add more detail. An example of this could be using a superclass of employee and then to have two subclasses called fulltime and part time. As employee could have name, address and other details whilst full time could just have salary and part time could work out the salary from part time hours worked, as the full time members of staff wouldn't need these.\n', 'Inheritance is a way to form new classes (instances of which are called objects) using classes that have already been defined. The new classes, known as derived classes, take over (or inherit) attributes and behavior of the pre-existing classes, which are referred to as base classes (or ancestor classes). It is intended to help reuse existing code with little or no modification.\n\nAn advantage of inheritance is that modules with sufficiently similar interfaces can share a lot of code, reducing the complexity of

the program. Inheritance therefore has another view, a dual, called polymorphism, which describes many pieces of code being controlled by shared control code.

Inheritance is typically accomplished either by overriding (replacing) one or more methods exposed by ancestor, or by adding new methods to those exposed by an ancestor.

In defining this inheritance hierarchy we have already defined certain restrictions, not all of which are desirable. Singleness: using single inheritance, a subclass can inherit from only one superclass. Visibility: whenever client code has access to an object, it generally has access to all the object's superclass data. Static: the inheritance hierarchy of an object is fixed at instantiation when the object's type is selected and does not change with time.

"When we talk about inheritance in object-oriented programming languages, which is a concept that was invented in 1967 for Simula, we are usually talking about a way to form new classes and classes are instances of which are called objects and involve using classes that have already been defined. Derived classes are intended to help reuse existing code with little or no modification and are the new classes that take over (or inherit) attributes and behavior of the pre-existing classes, usually referred to as base classes (or ancestor classes). Categorization in computer languages is a powerful way number of processing information and inheritance provides the support for representation by categorization. Furthermore, it is fundamental for helping humans learn by means of generalization in what is known about specific entities is applied to a wider group given a belongs relation can be established and cognitive processing which involves less information being acquired to be stored about each specific entity, but in actual fact only its particularities.

An instance of a "fruit" is a generalization of "apple", "orange", "mango" and many others. Inheritance can also sometimes be referred to as generalization, because is-a relationships represent a hierarchy amongst classes of objects. It can be considered that fruit is an abstraction of apple, orange, etc. Conversely, since apples are fruit, they may naturally inherit all the properties common to all fruit, such as being a fleshy container for the seed of a plant.

Modules with sufficiently similarities in interfaces would be able to share a lot of code and therefore reducing the complexity of the program. This can be known as one of the advantages of inheritance. Therefore inheritance can be known to have a further view, a dual, which describes many parts of code that are under control of shared control code, named as polymorphism.

On the other hand, inheritance is normally accomplished either by replacing one or more methods exposed by ancestor, or by adding new methods to those exposed by an ancestor. A well known term used for this replacing act is called overriding.

"In object-oriented programming, inheritance is the ability to specify one class to be a subclass of another; this leads to a hierarchy of classes, with the child classes inheriting and specialising - and sometimes adding to - the functionality and data structures of the parent classes. The hierarchy that is formed is also useful for the organisation of classes and objects, as it defines a relationship between the child and the parent (the child class is a "kind of" the parent class). Inheritance is useful for situations where several classes share common features, such as needed functions or data variables. In addition to this, child classes can be referenced in terms of their parent classes, which can be useful when storing large data structures of objects of several classes, which can all be referenced as one base class. Inheritance is a core aspect of object-oriented programming, and is available in some form or another in most, if not all, object oriented languages available today. Most of these languages provide an "extend" keyword, which is used to subclass another. Also, the functions and data variables that are inherited by the subclasses can be controlled through the use of visibility modifiers."

Inheritance is a concept in Object Oriented programming where a child- or sub-class inherits characteristics from a parent- or super-class. The concept takes its name from genetic inheritance where a child can inherit genetic characteristics from its parents.

Inheritance, at its simplest, allows programmers to model a

relationship where one object is a kind of another. For instance two classes, one representing an undergraduate student and another representing a post-graduate student could both be said to belong to a more generalised class representing all students. Similarly, we could say that dogs and cats are two kinds of animal, or that bridges and skyscrapers are two types of man-made structure.

Subclasses are said to extend or specialise their superclasses. Attributes (variables) and behaviours (functions) that are common between classes can be included in the definition of the superclass, leaving the subclass definitions containing only the attributes and behaviours that are unique to that class.

Inheritance can be used to create a multiple level architecture of classes. In such an architecture even the bottom-most subclasses inherit all of the attributes and behaviours that are defined in the very top-most superclasses. This can save the programmer time because it renders unnecessary a lot of code duplication.

In object-oriented programming, inheritance is a way to form new classes (instances of which are called objects) using classes that have already been defined.

Inheritance is also sometimes called generalization, because the is-a relationships represent a hierarchy between classes of objects. For instance, a "fruit" is a generalization of "apple", "orange", "mango" and many others. One can consider fruit to be an abstraction of apple, orange, etc. Conversely, since apples are fruit (i.e., an apple is-a fruit), apples may naturally inherit all the properties common to all fruit, such as being a fleshy container for the seed of a plant.

Inheritance is typically accomplished either by overriding (replacing) none or more methods exposed by ancestor, or by adding new methods to those exposed by an ancestor.

Inheritance is the ability of a subclass to inherit default, protected and public attributes and methods from its superclasses. Each object (except java.lang.Object) can be cast to an object of one of its superclasses. However an object cannot be cast to a class which is no relative of it. Here is an example of inheritance:

We have the class of all living things which have attributes like weight and age. We have the classes of animals, plants, viruses and fungi that are subclasses of the class of all living things. The animals have their unique attributes (organs, hair, etc.) and methods (walking, mating, etc.). They also inherit the attributes and methods of its superclass. Animals can be treated (cast) to living things. However, animals cannot be treated as fungi.

In object oriented programming inheritance is also dependant on access level modifiers. For example private attributes and methods cannot be inherited. Virtual attributes and methods can be shadowed/overridden. In Java all attributes and methods are implicitly virtual. Object variable can store a reference to the same class or a subclass (i.e. this or more specialised version). However, object variables cannot store references to a superclass (i.e. less specialised version) of the original class.

In object-oriented programming, inheritance is a way to form new classes (instances of which are called objects) using classes that have already been defined. The inheritance concept was invented in 1967 for Simula.

Inheritance provides the support for representation by categorization in computer languages. Categorization is a powerful mechanism number of information processing, crucial to human learning by means of generalization and cognitive economy (less information needs to be stored about each specific entity, only its particularities).

The new classes, known as derived classes, take over (or inherit) attributes and behavior of the pre-existing classes, which are referred to as base classes (or ancestor classes). It is intended to help reuse existing code with little or no modification.

Inheritance is also sometimes called generalization, because the is-a relationships represent a hierarchy between classes of objects. For instance, a "fruit" is a generalization of "apple", "orange", "mango" and many others. One can consider fruit to be an abstraction of apple, orange, etc. Conversely, since apples are fruit (i.e., an apple is-a fruit), apples may naturally inherit all the properties common to all fruit, such as being a fleshy container for the seed of a plant.

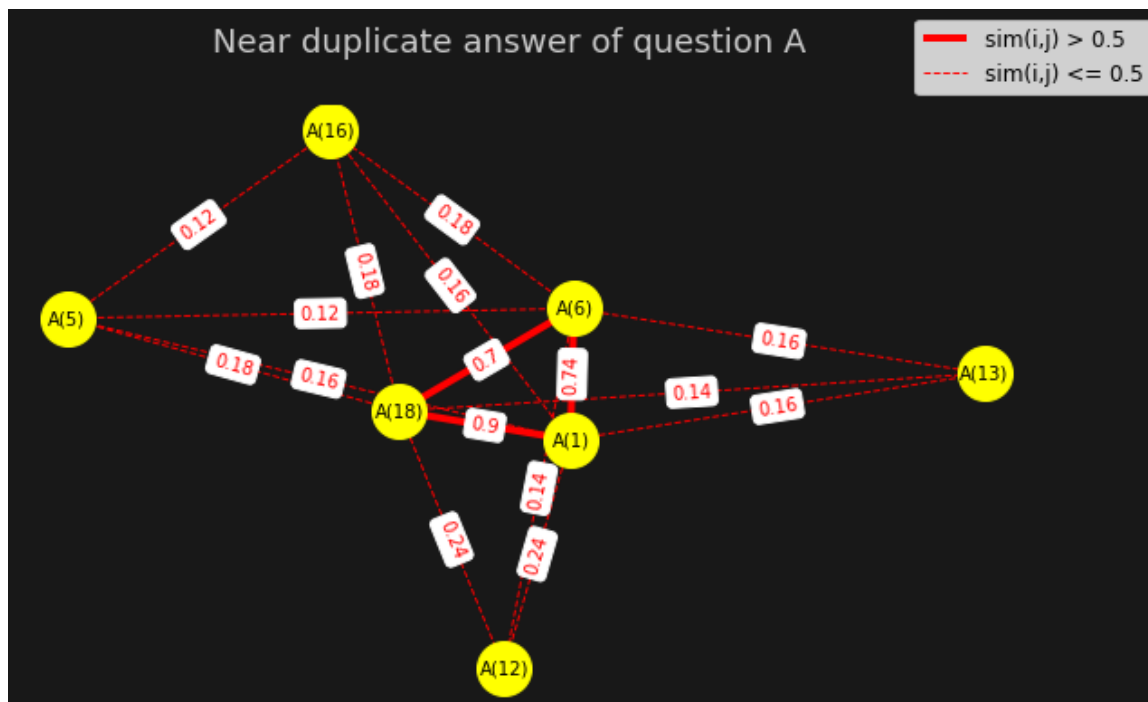
An advantage of inheritance is that modules with sufficiently similar interfa

ces can share a lot of code, reducing the complexity of the program. Inheritance therefore has another view, a dual, called polymorphism, which describes many pieces of code being controlled by shared control code. Inheritance is typically accomplished either by overriding (replacing) one or more methods exposed by ancestor, or by adding new methods to those exposed by an ancestor. Complex inheritance, or inheritance used within a design that is not sufficiently mature, may lead to the Yo-yo problem. The idea of inheritance in OOP refers to the formation of new classes with the already existing classes. The concept of inheritance was basically formulated for Simula in 1967. As a result, the newly created inherited or derived classes inherit the properties and behavior of the classes from which they are derived. These original classes are either called base classes or sometimes referred to as ancestor classes. The idea of inheritance is to reuse the existing code with little or no modification at all. The basic support provided by inheritance is that it represents by categorization in computer languages. The power mechanism number of information processing that is crucial to human learning by the means of generalization and cognitive economy is called categorization. Where generalization if the knowledge of specific entities and is applied to a wider group provided that belongs relation can be created. On the other hand cognitive economy is where less information needs to be stored about each specific entity except for some particularities. There are examples where we can have modules with similar interfaces. The advantage that inheritance provides is that it makes such modules share a lot of code which consequently reduces the complexity of the program.

#task: 20

Signatures metric: 20 x 100

#permutations used to create signatures: 100



**Question B:** The result illustrates that there are no pairs of answers which can be identified as near duplicate answers

In [ ]:

```
LSH_task('B')
```

['PageRank is a link analysis algorithm used by the Google Internet search engine that assigns a numerical weighting to each element of a hyperlinked set of documents, such as the World Wide Web, with the purpose of "measuring" its relative importance within the set. The algorithm may be applied to any collection of entities with reciprocal quotations and references. The numerical weight that it assigns to any given element E is also called the PageRank of E and denoted by  $PR(E)$ . The name "PageRank" is a trademark of Google, and the PageRank process has been patented (U.S. Patent 6,285,999). However, the patent is assigned to Stanford University and not to Google. Google has exclusive license rights on the patent from Stanford University. The university received 1.8 million shares in Google in exchange for use of the patent; the shares were sold in 2005 for \$336 million. Google describes PageRank: "PageRank relies on the uniquely democratic nature of the web by using its vast link structure as an indicator of an individual page's value. In essence, Google interprets a link from page A to page B as a vote, by page A, for page B. But, Google looks at more than the sheer volume of votes, or links a page receives; it also analyzes the page that casts the vote. Votes cast by pages that are themselves "important" weigh more heavily and help to make other pages "important". In other words, a PageRank results from a "ballot" among all the other pages on the World Wide Web about how important a page is. A hyperlink to a page counts as a vote of support. The PageRank of a page is defined recursively and depends on the number and PageRank metric of all pages that link to it ("incoming links"). A page that is linked to by many pages with high PageRank receives a high rank itself. If there are no links to a web page there is no support for that page. Google assigns a numeric weighting from 0-10 for each webpage on the Internet; this PageRank denotes a site's importance in the eyes of Google. The PageRank is derived from a theoretical probability value on a logarithmic scale like the Richter Scale. The PageRank of a particular page is roughly based upon the quantity of inbound links as well as the PageRank of the pages providing the links. It is known that other factors, e.g. relevance of search words on the page and actual visits to the page reported by the Google toolbar also influence the PageRank. In order to prevent manipulation, spoofing and Spamdexing, Google provides no specific details about how other factors influence PageRank. Numerous academic papers concerning PageRank have been published since Page and Brin's original paper. In practice, the PageRank concept has proven to be vulnerable to manipulation, and extensive research has been devoted to identifying falsely inflated PageRank and ways to ignore links from documents with falsely inflated PageRank. Other link-based ranking algorithms for Web pages include the HITS algorithm invented by Jon Kleinberg (used by Teoma and now Ask.com), the IBM CLEVER project, and the TrustRank algorithm. 'PageRank algorithm is also known as link analysis algorithm. It has been used by google. The algorithm may be applied to any collection of entities with reciprocal quotations and hyperlinked set of documents, such as the World Wide Web, with the purpose of "measuring references. The name "PageRank" is a trademark of Google, and the PageRank process has been patented (U.S. Patent 6,285,999). The numerical weight that it assigns to any given element E is also called the PageRank of E and denoted by  $PR(E)$ . The name "PageRank" is a trademark of Google, and the PageRank process has been patented (U.S. Patent 6,285,999). However, the patent is assigned to Stanford University and not to Google. Google has exclusive license rights on the patent from Stanford University. In other words, a PageRank results from a "ballot" among all the other pages on the World Wide Web about how important a page is. A hyperlink to a page counts as a vote of support. The PageRank of a page is defined recursively and depends on the number and PageRank metric of all pages that link to it ("incoming links"). Numerous academic papers concerning PageRank have been published since Page and Brin's original paper.[4] In practice, the PageRank concept has proven to be vulnerable to manipulation, and extensive research has been devoted to identifying falsely inflated PageRank and ways to ignore links from documents

nts with falsely inflated PageRank\n\t”\n\n', 'PageRank is a link analysis algorithm used by the Google Internet search engine that assigns a numerical weighting to each element of a hyperlinked set of documents, such as the World Wide Web, with the purpose of "measuring" its relative importance within the set. Google assigns a numeric weighting from 0-10 for each webpage on the Internet; this PageRank? denotes a site's importance in the eyes of Google.\n\nThe PageRank? is derived from a theoretical probability value on a logarithmic scale like the Richter Scale. The PageRank? of a particular page is roughly based upon the quantity of inbound links as well as the PageRank? of the pages providing the links. The algorithm may be applied to any collection of entities with reciprocal quotations and references. The numerical weight that it assigns to any given element E is also called the PageRank? of E and denoted by  $PR(E)$ .\n\nIt is known that other factors, e.g. relevance of search words on the page and actual visits to the page reported by the Google toolbar also influence the PageRank?. Other link-based ranking algorithms for Web pages include the HITS algorithm invented by Jon Kleinberg (used by Teoma and now Ask.com), the IBM CLEVER project, and the TrustRank? algorithm. \n', "PageRank (PR) refers to both the concept and the Google system used\nfor ranking the importance of pages on the web. The "PageRank" of a\nsite refers to its importance or value on the web in relation to the\nrest of the sites that have been "PageRank"ed.\n\nThe algorithm basically works like a popularity contest - if your site\nis linked to by popular websites, then your site is considered more\npopular. However, the PR doesn't just apply to the website as a whole\n- different pages within a website get given different PRs dependent\non a number of factors:\n\n\* Inbound links (backlinks) - how many pages (other than the ones on your website) link to this particular page\n\n\* Outbound links (forward links) - how many external pages the particular page links to\n\n\* Dangling links - how many pages with no external links are linked to from a particular page\n\n\* Deep links - how many links that are not the homepage are linked to from a particular page\n\nPR tries to emulate a "random surfer". The algorithm includes a\ndampening factor, which is the probability that a random surfer will\nget bored and go and visit a new page - by default, this is 0.85. A\nvariation on this is the "intentional surfer", where the importance of\na page is based on the actual visits to sites by users. This method is\nused in the Google Toolbar, which reports back actual site visits to\nGoogle.\n', 'There are many attributes which influence the ranking of a page in google, The main too are the content, key words, and links. The content of a webpage generally gives a good idea about what the page is about, however, there are some flaws in this, for example, for along time ibm web page didnt contain the word computer despite it being strongly associated with them. To solve this problem, web pages can assign itself key words, which contribute to its ranking in searches.\n\nThe second method is the use of links. the more sights which links to your web page and the higher the rank of those sights, the higher the rank of your site will be. This method is used as links are seen as an adoursment of a sight.\n\nWith both these methods of ranking web pages, there are issues. key words can be compromised by sparming, google solves this problem by penolizing such activity. Useing links to rank a page also has its problems, for example, link farms which have recursive links, for the sole perpos of raising there ranking, google takels this by useing a dampaning algorithem.\n\n', 'PageRank algorithm is patented by Stanford University. It is a link analysis algorithm employed by the Google Internet search engine that assigns a value used to measure the importance to each element of a hyperlinked set of documents, such as the WWW, with the purpose of "measuring" its relative significance within the set.\n\nGoogle owns exclusive license rights on the patent from Stanford University. The University received 1.8 million shares in Google in return for use of the patent. \n', 'PageRank is a link analysis algorithm used by the Google Internet search engine that assigns a numerical weighting to each element of a hyperlinked set of documents, such as the World Wide Web, with the purpose of "measuring" its rela

tive importance within the set. The algorithm may be applied to any collection of entities with reciprocal quotations and references. PageRank Uses in google toolbar: Measures popularity of a site ,Marketing value,Updated periodically, in google directory: PageRank: sort links within categories; Volunteers evaluate, classify, annotate;Open Directory project using PageRank.

\n', 'The PageRank algorithm is used to designate every aspect of a set of hyperlinked documents with a numerical weighting. It is used by the Google search engine to estimate the relative importance of a web page according to this weighting. The system uses probability distribution to determine the odds that a person randomly clicking on links will arrive at any given page. Following this, each web page is given a ranking of 0-10 according to its relevance to a search. The PageRank is calculated by taking into consideration the number of inbound links, and the PageRank of the pages supplying these links. This means therefore that if a webpage is linked to others that have a high ranking, then it too will receive a high rank.

\n\nDue to the nature of the PageRank system, it is susceptible to manipulation and has been exploited so that certain pages are given a false, exaggerated ranking. In these cases, only Google has access to the genuine PageRank. However, much research has been conducted into methods of avoiding links from documents with a false PageRank to try and iron out the bugs in this system and from 2007 Google has actively penalized schemes which try to increase rankings artificially.

\n', 'A website's page rank, is how 'important' it is on the web. It is essentially a popularity meter. Popularity or importance is determined by the amount of links relating to the page there are, there are four different types. Inbound, links from other pages to yours. Outbound, links from your page to others. Dangling, links to a page which has no links to others. Deep, links to a specific page, usually bypassing the homepage. The page rank algorithm takes the probability of a 'random surfer' becoming bored and requesting another 'random page' (otherwise known as the dampening factor) away from 1 and divides this number by the number of pages in the system, adding it to the dampening factor multiplied by the page rank of a linked page divided by the number of outbound links on that linked page. Adding on this last section for every other page linked to from the original page. Google uses this algorithm to assist intentional surfers in finding the best websites to suit their needs. One of the problems with this popularity algorithm is that it is easily manipulated and can give false values, hence the frequent recalculating of page ranks.

\n', 'PageRank is a link analysis algorithm used by the Google Internet search engine that assigns a numerical weighting to each element of a hyperlinked set of documents, such as the World Wide Web, with the purpose of "measuring" its relative importance within the set. Google assigns a numeric weighting from 0-10 for each webpage on the Internet; this PageRank denotes a site's importance in the eyes of Google. The PageRank is derived from a theoretical probability value on a logarithmic scale like the Richter Scale. PageRank is a probability distribution used to represent the likelihood that a person randomly clicking on links will arrive at any particular page. PageRank can be calculated for collections of documents of any size. It is assumed in several research papers that the distribution is evenly divided between all documents in the collection at the beginning of the computational process. The PageRank computations require several passes, called "iterations", through the collection to adjust approximate PageRank values to more closely reflect the theoretical true value. A probability is expressed as a numeric value between 0 and 1. A 0.5 probability is commonly expressed as a "50% chance" of something happening. Hence, a PageRank of 0.5 means there is a 50% chance that a person clicking on a random link will be directed to the document with the 0.5 PageRank. The PageRank theory holds that even an imaginary surfer who is randomly clicking on links will eventually stop clicking. The probability, at any step, that the person will continue is a damping factor  $d$ . Various studies have tested different damping factors, but it is generally assumed that the damping factor will be set around 0.85.

\n', 'The algorithm that Google uses to ass



ign a weighting to each element of a linked set of documents, with the purpose of "measuring" its relative importance within the set. \nA particular websites PageRank results from a "vote" from other pages on the Internet about how important that website actually is. A link to a page is seen as a vote of support. The PageRank depends on the PageRank rating and number of all pages that have links to it. Additionally, if a page is linked to by pages with a high PageRank rating, this increases the rating of the original page.\n\nThe PageRank scale ranges from 0-10. The rating of a certain page is generally based upon the quantity of inbound links as well as the perceived quality of the pages providing the links. \n\nPageRank could be described as a probability distribution representing the chance that someone randomly clicking on links will reach a certain page. The PageRank calculations require iterations through the collection of web pages to alter approximate PageRank values to accurately reflect the actual rank.\n\nIn order to prevent spamming, Google releases little information on the way in which a PageRank is calculated. The PageRank algorithm has led to many sites being spammed with links in an attempt to artificially inflate the PageRank of the linked page, notably in blog comments and message boards. In 2005 a 'no follow' tag was added as an attribute of a HTML link to be used where Google shouldn't change the PageRank of the linked page as a result of the link. ', 'The first thing to consider when talking about Google's PageRank algorithm, is that a PageRank is essentially how important that web page is to the internet. So in essence it is a popularity contest between WebPages.\n\nOriginally search engines used highest keyword density, however this could be abused if keyword spamming was implemented. Instead Google uses a system that is based on sites linking to each other, e.g. the more important a site is that is linked to yours the higher your site will become. \n\nThe algorithm Google actually uses is based on 4 factors, total number of pages, dampening factor, PageRank of a single page and the number of outbound links. A dampening factor is used to counter random surfers, who get bored and then switch to other pages. This formula is then re-used until the results seem to converge together, to find the PageRank, so it is calculated iteratively.\n\nPageRank is used by Google to measure a popularity of the site and a number between 0-10 is assigned to each webpage depending on their PageRank. This allows Google to calculate a marketing value for different WebPages.\n\nAlso it should be noted that the PageRank is periodically updated every 3 to 6 months, this is counter hackers influence on different PageRanks.\n', "The PageRank is a recursive algorithm used by Google to determine which webpages are more important than others. The algorithm considers the importance of a webpage to be reflected by how many other webpages link to that page, and the importance of those pages.\n\nFor each page that links to a page A, the PageRank between zero and one is calculated iteratively according to the following two key factors: The probability of a user navigating away from a page randomly; the PageRank of any page that links to A, divided by the total number of outbound links from that page. This assumes that a link among many outbound links is less valuable than a link among fewer outbound links. A variation of the PageRank method bases the importance of a webpage on how many visits the page gets. \n\nThe method can be abused when people deliberately link to sites in order to raise a site's PageRank. However, it is still a good indicator for search engines to use as a variable in deciding on the most appropriate results according to a query. \n", 'PageRank is a link analysis algorithm that is used by search engine such as Google Internet that assigns a numerical weighting to every element of a hyperlinked set of documents, like the World Wide Web, with the hope of "measuring" the relative importance held in the set. The algorithm may be applied to any number of entities with reciprocal quotations and references. The weight taking a numerical value which assigns to any given element E is also known as the PageRank of E and is denoted by  $PR(E)$ . \nA trademark of Google has the name "PageRank" and this process has been patented (U.S. Patent 6,285,999\xa0). Nevertheless, the patent is assigned to the University of Stanford and not to Google. Google h

as exclusive license rights on the patent from the University of Stanford and the university received 1.8 million shares in Google in exchange for use of the patent; the in the year 2005, shares were sold for \$336 million.

"The Google search engine uses a link analysis algorithm called PageRank to assign a relative numerical importance to a set of hyperlinked documents, such as the World Wide Web.

For a given page, its importance (the PageRank value) results from a ballot among all the other pages in the set. For a page to give a vote to another, it must link to it, and so the PageRank depends on the number of incoming links, and the PageRank of those pages that provide the links. Pages that are linked to by many high ranking pages will themselves obtain a high rank. If a page has no incoming links, there is no support for that page.

The PageRank is a numeric weighting of 0 to 10, and denotes how important a site is in Google's eyes. Like the Richter Scale, the PageRank is a value on a logarithmic scale that is derived from a probability. In addition to the quantity and quality of inbound links, other factors affect the PageRank, such as the number of visits to the page and the search words that are used on the page.

To prevent sites from manipulating or spoofing PageRank, very little details are provided by Google as to what factors actually affect it.

'PageRank is an algorithm that was developed by Google to provide the most relevant search results to its users queries. PageRank, along with similar algorithms developed by Google's competitors for their search engines, is part of the second generation of technologies designed to rate the importance of web pages: the first, which was solely based on keywords in the page content and meta-data, could easily be influenced by those wishing to obtain a higher ranking for their less-relevant pages.

The difference with PageRank is that it tries to determine a web page's relevance to users by attempting to determine its importance. It does this by assigning it a value of importance that is dependant upon the number of web sites that link to that page, taking into account the importance value, or PageRank, of those pages. The PageRank is computed iteratively, and it is found that the PageRank values converge fairly rapidly.

Although it is much better than simple keyword-based ranking algorithms, PageRank is not infallible: we have an internet where advertising revenue can make up most - and quite frequently all - of a web site's income and the people that run these web sites will always be trying to trick the system into giving their pages a higher PageRank. One of Google's attempts to counter this is their Google Toolbar browser plugin.

Google Toolbar is a free tool which provides a number of useful functions in a convenient location: the user's web browser window. Google's payoff is that it gets to track the behaviour of actual users. This allows them to see whether their PageRank algorithm is accurate in assigning high PageRank values to the most relevant web pages and, just as importantly, low values to those that are irrelevant and try to fool the system.

'The PageRank algorithm used by Google harnesses the implicit collective intelligence present in the structure of the world wide web. Any page on the Internet will generally link to at least one other, by modelling this link structure as a graph, we can build up a symbolic representation of the world wide web.

As the basic level, the nodes with the highest degrees can be considered the most "popular" and by inference the most important - which can be used to rank the pages when returning search results.

Expanding on this theory, we can then say that the links from an important page are themselves more important. Using this idea we can adjust the rankings of our pages so that pages linked to be the most important pages are considered more relevant.

The actual Google PageRank algorithm is much more complex than this, but follows the same underlying principles. It incorporates some more advanced reasoning to avoid website creators exploiting their knowledge of the algorithm to artificially increase their PageRank through use of web-rings and other similar reciprocal hyperlinking schemes.'

'Page rank algorithm is used to determine a webpage's importance or relevance in the web dependant on certain criteria. The criteria may include numbers of word matches with the search terms, number of

other webpages that link this one and/or cite it as a source, number of unique visits for certain amount of time etc. There are some techniques that try to fool the search engines like link farms, keyword spamming and a lot of meta tags. The last two are somewhat easier to be dealt with (simply by being ignored most of the time). Link farms are groups of sites that are producing links between each other pursuing higher link counts. The reason for such manipulations is the pursuit of higher page rank so even higher number of users will see the page which will lead to higher income. Link farms can be exploited by joining to them and get inbound linkage but refuse to add links for one's own site to the sites from the link farm. Google's toolbar tries to follow the intentional user model by counting the visits from actual users (i.e. not computer bots) to a website. Page ranks can be calculated either recursively or iteratively. One of the most important uses of page rank is its meaning to advertising.

'Since the development of the Web 2.0, Google as one of the most popular search engine in the world, there are many algorithms in the web search. Accordingly, implementations of link analysis algorithms will typically discount such "internal" links. The word computer can be exploited by web search engines such as Google. Thus, the web is just like a graph, and the PageRank, which is our first technique for analysing the link which is assigns to every node in the web graph a numerical score between 0 and 1. Since the PageRank is the most important algorithms which is used in the Google engine. For example, there are four pages group: A, B, C and D. If every page link to A, then A's PageRank value should be the total value of B, C and D.

$$PR(A) = PR(B) + PR(C) + PR(D)$$

Moreover, there is a  $q = 0.15$  which is be use in the web page, like the general algorithm below:

However, the disadvantage is of PageRank algorithm is that the renew system is too slow.

'PageRank is a probability distribution used to represent the likelihood that a person randomly clicking on links will arrive at any particular page. . It is assumed in several research papers that the distribution is evenly divided between all documents in the collection at the beginning of the computational process. PageRank can be calculated for collections of documents of any size The PageRank computations require several passes, called "iterations", through the collection to adjust approximate PageRank values to more closely reflect the theoretical true value.

A probability is expressed as a numeric value between 0 and 1. A 0.5 probability is commonly expressed as a "50% chance" of something happening. Hence, a PageRank of 0.5 means there is a 50% chance that a person clicking on a random link will be directed to the document with the 0.5 PageRank.

Simplified algorithm

How PageRank Works

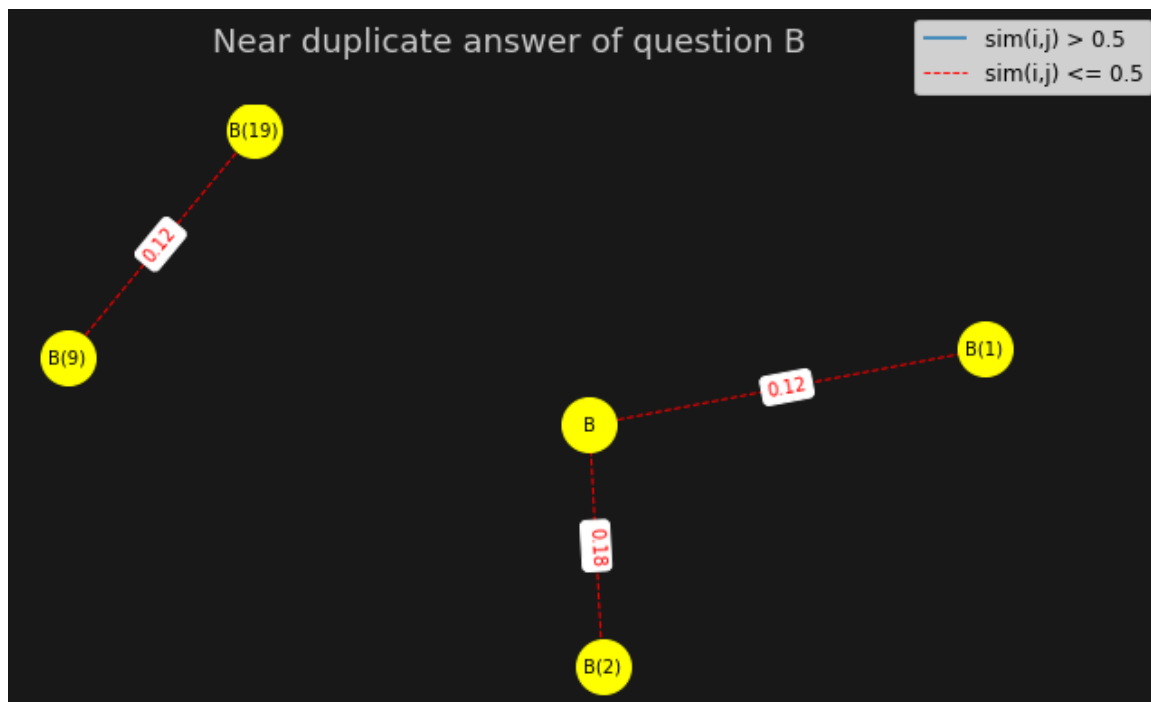
Assume a small universe of four web pages: A, B, C and D. The initial approximation of PageRank would be evenly divided between these four documents. Hence, each document would begin with an estimated PageRank of 0.25.

In the original form of PageRank initial values were simply 1. This meant that the sum of all pages was the total number of pages on the web. Later versions of PageRank (see the below formulas) would assume a probability distribution between 0 and 1. Here we're going to simply use a probability distribution hence the initial value of 0.25.]

#task: 20

Signatures metric: 20 x 100

#permutations used to create signatures: 100



**Question C:** The result illustrates that there are no pairs of answers which can be identified as near duplicate answers

In [ ]:

```
LSH_task('C')
```

['Vector space model (or term vector model) is an algebraic model for representing text documents (and any objects, in general) as vectors of identifiers, such as, for example, index terms. It is used in information filtering, information retrieval, indexing and relevancy rankings. Its first use was in the SMART Information Retrieval System.\nA document is represented as a vector. Each dimension corresponds to a separate term. If a term occurs in the document, its value in the vector is non-zero. Several different ways of computing these values, also known as (term) weights, have been developed. One of the best known schemes is tf-idf weighting (see the example below).\nThe definition of term depends on the application. Typically terms are single words, keywords, or longer phrases. If the words are chosen to be the terms, the dimensionality of the vector is the number of words in the vocabulary (the number of distinct words occurring in the corpus).\n\nThe vector space model has the following limitations:\n 1. Long documents are poorly represented because they have poor similarity values (a small scalar product and a large dimensionality)\n 2. Search keywords must precisely match document terms; word substrings might result in a "false positive match"\n 3. Semantic sensitivity; documents with similar context but different term vocabulary won't be associated, resulting in a "false negative match".\n 4. The order in which the terms appear in the document is lost in the vector space representation.\n', 'The definition of term depends on the application. Typically terms are single words, keywords, or longer phrases. If the words are chosen to be the terms, the dimensionality of the vector is the number of words in the vocabulary A document is represented as a vector. Each dimensions corresponds to a separate terms. If a term occurs in the document, its value in the vector is non-zero.\nRelevancy rankings of documents in a keyword search can be calculated, using the assumptions of document similarities theory, by comparing the deviation of angles between each document vector and the original query vector where the query is represented as same kind of vector as the documents.\nLIMITATION:\n\nThere is some limitation of vector space model.\nModels based on and extending the vector space model include:\n•\tGeneralized vector space model.\n•\t(enhanced) Topic-based Vector Space Model [1] (eTVSM) – Extends the vector space model by removing the constraint that the term-vectors be orthogonal. In contrast to the generalized vector space model the (enhanced) Topic-based Vector Space Model does not depend on concurrence-based similarities between terms. The enhancement of the enhanced Topic-based Vector Space Model (compared to the not enhanced one) is a proposal on how to derive term-vectors from an Ontology. \n\n', 'The vector space model (also called, term vector model) is an algebraic model used to represent text documents, as well as any objects in general, as vectors of identifiers. It is used in information retrieval and was first used in the SMART Information Retrieval System.\n\nA document is represented as a vector and each dimension corresponds to a separate term. If a term appears in the document then its value in the vector is non-zero. Many different ways of calculating these values, also known as (term) weights, have been developed. One of the best known methods is called tf-idf weighting.\n\nThe definition of term depends on the application but generally terms are single words, keywords, or longer phrases. If the words are chosen to be the terms, the dimensionality of the vector is the number of words in the vocabulary, which is the number of distinct words occurring in the corpus.\n\nThe vector space model has several disadvantages. Firstly, long documents are represented badly because they have poor similarity values. Secondly, search keywords must accurately match document terms and substrings of words might result in a "false-positive match". Thirdly, documents with similar context but different term vocabulary will not be associated, resulting in a "false-negative match". Finally, the order in which the terms appear in the document is lost in the vector space representation. \n', 'Vector space model is an algebraic model for representing text documents (and in general, any objects) as vectors of identifiers, such as, for example, index terms. Its first use was in the SMART Information Retrieval System. It is used in informa

tion filtering, information retrieval, indexing and relevancy rankings.

A document is represented as a vector, and each dimension corresponds to a separate term. If a term occurs in the document, its value in the vector is non-zero. Several different ways of computing these values, also known as (term) weights, have been developed. The definition of term depends on the application. Typically terms are single words, keywords, or longer phrases. If the words are chosen to be the terms, the dimensionality of the vector is the number of words in the vocabulary (the number of distinct words occurring in the corpus).

One of the best known schemes is tf-idf weighting, proposed by Salton, Wong and Yang. In the classic vector space model, the term specific weights in the document vectors are products of local and global parameters.

Relevancy rankings of documents in a keyword search can be calculated, using the assumptions of document similarities theory, by comparing the deviation of angles between each document vector and the original query vector where the query is represented as same kind of vector as the documents.

The vector space model has the following limitations:

- \* Search keywords must precisely match document terms; word substrings might result in a "false positive match";
- \* Semantic sensitivity; documents with similar context but different term vocabulary won't be associated, resulting in a "false negative match";
- \* The order in which the terms appear in the document is lost in the vector space representation;
- \* Long documents are poorly represented because they have poor similarity values (a small scalar product and a large dimensionality).

The vector space model is where each document is viewed as a bag of words, where there order has little significance. Each document is a vector where each word is a dimension. The vector is then constructed of the frequency of each word (dimension). The drawback to this approach is that the length of the document has an impact on the vector, to compensate for this you can compute the cosine similarity between your two comparison documents. This will find the difference between the two vectors (the dot product), ignoring the size of them.

In order to query the search space, the query can also be represented as a vector, then you find the document whose vector has the greatest cosine similarities to your query. There are a number of weighting schemes which can be incorporated in order to increase the accuracy of the vectors.

There are some drawbacks with this approach, Computing the cosine similarities between each vector can be expensive as the number of dimensions can be in the thousands, To tackle this problem you can use inverted indexes and then a series of heuristics in order to improve on this.

to top

An algebraic model for representing text documents and any objects in general is known by the name Vector space model. It represents these as vectors of identifiers, index terms are one illustration of these. The Vector Space model was first used in the SMART Information Retrieval System, and it is utilised variously in indexing, information filtering, indexing and information retrieval.

A document has representation as a vector. Every dimension is precisely related to a separate term. The way in which term is defined depends entirely on the application: typically 'terms' are either single words, keywords or longer phrases. The dimensionality of the vector is the number of words in the vocabulary, if it is the words that are chosen to be the terms. So the same rule applies with keywords and indeed longer phrases.

If a term occurs in the document, its value in the vector is non-zero. Several different ways of computing these values, additionally known as (term) weights, have been developed. One of the most famous schemes is tf-idf weighting.

The representation of a set of documents as vectors in a common vector space is known as the vector space model and is fundamental to a host of information retrieval (IR) operations including scoring documents on a query, document classification, and document clustering. We first develop the basic ideas underlying vector space scoring; a pivotal step in this development is the view of queries as vectors in the same vector space as the document collection.

The vector space model is an algebraic model used to represent text documents (and any objects, generally) as vectors of identifiers, for

instance index terms. Its applications include information filtering, information retrieval, indexing and relevancy rankings. With reference to this model, documents are represented as vectors. Each dimension corresponds to a separate term. The value of a vector is non-zero if a term occurs in the document. Several different ways have been developed of calculating these values (also known as term weights). One of the best known schemes is tf-idf (term frequency-inverse document frequency) weighting.

The model can be used to determine the relevancy rankings of documents in a keyword search, using the assumptions of document similarities theory, by comparing the original query vector (where the query is represented as same kind of vector as the documents) and the deviation of angles between each document vector.

The classic vector space model was put forward by Salton, Wong and Yang and is known as term frequency-inverse document frequency model. In this classic model the term specific weights in the document vectors are products of local and global parameters. In a simpler Term Count Model the term specific weights are just the counts of term occurrences and therefore do not include the global parameter.

The algebraic model for representing text documents and objects as vectors of identifiers is called the vector space model. It is used in information filtering, indexing, relevancy rankings and information retrieval. It was first used in the SMART Information Retrieval System.

When a document is represented as a vector, each dimension corresponds to a separate term. A term which occurs in the document has a value in the vector of non-zero. Other ways of computing these values, or weights, have been developed. The most popular is tf-idf weighting.

Depending on the application, the definition of term varies. Single words, keywords and occasionally longer phrases are used for terms. The dimensionality of the vector, if words are used as terms, is the total number of words available for use. By using the assumptions of the document similarities theory, the relevancy rankings of documents in a keyword search can be worked out by comparing the deviation of angles between vectors both within the document and the original query where the vectors of both are the same type.

The limitations of the vector space model are thus. Due to poor similarity values long documents are poorly represented. False positive matches may be returned if search keywords do not precisely match document terms. False negative matches could be returned when documents share a context but have different term vocabulary. Vector space representation results in the loss of the order which the terms are in the document.

Within Information Retrieval each document in a set can be represented as a point in high-dimensional vector space, this representation is called the vector space model. Information Retrieval queries are also represented as vectors in the same vector space; these are then used in conjunction with the document vectors to find relevant documents. The two vectors are compared and the documents with a higher document-query similarity are ranked higher in terms of relevance. There are a variety of techniques that can be used to compare the two vectors; the most frequently used method for the vector space model is the Cosine Coefficient, which calculates the angle between the two vectors and produces a value between 0 and 1.

A Vector space model (or term vector model) is an algebraic way of representing text documents (and any objects, in general) as vectors of identifiers, such as index terms. It is used in information filtering, information retrieval, indexing and relevancy rankings. Its first application was in the SMART Information Retrieval System.

A document can be represented as a vector. Every dimension relates to a different term. If a term appears in the document, the terms value in the vector is non-zero. Many different methods of calculating these values, sometimes known as (term) weights, have been developed. tf-idf weighting is one of the most well known schemes. (see below example).

The definition of a term depends on the application. Normally a term is a single word, keyword, or a longer phrase. If the words are chosen to be the terms, the dimensionality of the vector is the number of words in the vocabulary (the number of distinct words occurring in the corpus).

The vector space model has some limitations:

1. Long



er documents are represented poorly because the documents have poor similarity values (namely a small scalar product and a large dimensionality)\n

- 2.\tSearch keywords have to precisely match document terms; word substrings could potentially result in a "false positive match"\n3.\tSemantic sensitivity; documents with a similar context, but different term vocabulary won't be associated, resulting in a "false negative match".\n4.\tThe order in which terms appear in the document is lost in a vector space representation.\n

'A Vector space model is an algebraic model for representing text documents as vectors of identifiers. A possible use for a vector space model is for retrieval and filtering of information. Other possible uses for vector space models are indexing and also to rank the relevancy of differing documents.\nTo explain further vector space models, basically a document is characterized by a vector. With each separate term corresponding to the differing dimensions. There has been multiple ways of trying to compute the different possible values for vector space models with the most recognised being the tf-idf weighting.\nThe differing application has a direct influence on what the definition of the term means. A normal term is usually a single word, keywords or longer phrases. The number of unique words in the vocabulary denotes the dimensionality, if words are used for the terms.\nHowever whilst vector space modelling is useful there are 4 key problems with using it, they are; that the order of the terms are lost, keywords must be precise if searched for, bigger documents have a poor similarity value due to being poorly represented and two documents based on the same topic won't be associated if term vocabulary differs.\n', 'In the vector space model (VSM), documents take the form of "bags of words" - a standard information retrieval approach which represents documents as in a mathematical "bag" structure, recording what terms are present and how often they occur. \n\nThe vector space model is used in information retrieval to determine how similar documents are to one another, and how similar documents are to a search query. \n\nIn a collection of documents, each document can be viewed as a vector of  $n$  values (the terms in the document), where each term is an axis. Queries can also be represented as vectors on this vector space model, and so deciding which document matches the query the closest becomes a matter of selecting the document vector which is nearest to the query vector. \n\nThe query vector is compared to each document vector in turn using a "vector similarity measure", which is the cosine of the angle between the query vector and the document vector. \n\nThis equation is calculated by dividing the dot product of the query vector and the document vector by the modulus of the query vector multiplied by the modulus of the document vector. The denominator takes into account differences in the length of the vector, and has the effect of "normalising" the length. Whichever document returns the highest cosine similarity score is considered to be the closest matching document to the query. \n', 'Information retrieval (IR) is the science of searching for documents, for information within documents and for metadata about documents, as well as that of searching relational databases and the World Wide Web. IR is interdisciplinary, based on computer science, mathematics, library science, information science, information architecture, cognitive psychology, linguistics, statistics and physics. There is overlap in the usage of the terms data retrieval, document retrieval, information retrieval, and text retrieval, but each also has its own body of literature, theory, praxis and technologies. \nAutomated information retrieval systems are used to reduce what has been called "information overload". Many universities and public libraries use IR systems to provide access to books, journals and other documents. \n', 'Vector space model, or term vector model as it is also known, is an algebraic model for representing objects (although it is mainly used for text documents) as vectors of identifiers; for example, index terms. It is used in information retrieval and filtering, indexing and relevancy rankings, and was first used in the SMART Information Retrieval System.\n\nA document is represented as a vector, with each dimension corresponding to a separate term. If a term occurs in the document, the value will be non-zero in the vector. Many

different ways of computing these values (aka (term) weights) have been developed; one of the best known schemes is tf-idf weighting.

The way that a 'term' is defined depends on the application. Typically, terms are single words, keywords, or sometimes even longer phrases. If the words are chosen as the terms, the number of dimensions in the vector is the number of distinct words in the corpus.

Relevancy ranks for documents, in a keyword search, can be calculated; this uses the assumptions of document similarities theory, by comparing the difference of angles between each document vector and the original query vector, where the query is represented as the same format vector as the documents.

Generally, it is easier to calculate the cosine of the angle between the vectors instead of the angle itself. A zero value for the cosine indicates that the query and document vector are orthogonal and so had no match; this means the query term did not exist in the document being considered.

However, the vector space model has limitations. Long documents are poorly represented due to their poor similarity values (a small scalar product and a large dimensionality); search keywords must match precisely the document terms; word substrings might result in a "false positive match"; similar context documents but different term vocabulary won't be associated, leading to a "false negative match"; and the order that the terms appear in the document is not represented in the vector space model.

There are a large number of models used in solving the problem of Information Retrieval and they are all based on one of three mathematical bases: set theory, algebra and probabilistic. The vector space model is one of these methods, and it is an algebraic model.

In the vector space model a document is represented as a vector. Within this vector, each dimension corresponds to a separate term (where a term is typically a single word, keyword or phrase.) If the term doesn't occur within the document, the value in the vector is zero. If a term occurs in the document, its value is non-zero.

To calculate how relevant each document is in a keyword search the cosine value of the angle between the vectors is easier to calculate instead of the actual angle.

The vector space model, however, is not without its limitations: they have small similarity values, long documents are poorly represented; the order of words does not matter; false positive matches may be brought about by terms contained within words themselves; and documents that should match but use different semantics may return false negative matches. There are a number of other models that are based on or extend the vector space model, and these are designed to try to eradicate these problems.

Using the vector space model for Information Retrieval models all pages and queries as high-dimensional sparse vectors. Each item in the vector represents a different keyword.

The similarity between two pages or a query and a page can be computed using the dot product formula to find the cosine between them. This represents the angle between them, but in n-dimensional space. Results will range from -1 to 1, with 1 being a close match. Normally the vector will not have any negative values, so results will always be greater than or equal to 0. The cosine is computed using:  $\cos A = (|a||b|)/(a.b)$ .

The vector space model (or term vector model) is an algebraic model for representing text documents (and any objects, in general) as vectors of identifiers, such as index terms. It is used in information filtering, information retrieval, indexing and relevancy rankings. It was used in the first time in the SMART Information Retrieval System.

A document is represented as a vector. Each and every dimension corresponds to a separate term. If a term exists in a document, its value in the vector is not equal to zero. A couple of different algorithms of computing these values, also known as (term) weights, have been created. One of the most popular schemes is tf-idf weighting.

The definition of term is dependent on the application. Typically terms are keywords, single words or longer phrases. Provided that words are selected to be the terms, the dimensionality of the vector is equal to the number of words in the vocabulary.

It is easiest to calculate the cosine of the angle between the vectors instead of the angle by the formula:

$$\cos(\theta) = \frac{v_1 \cdot v_2}{(|v_1| |v_2|)}$$

A null cosine value

blue says that the query and document vector were orthogonal and had no match which means that no term of the query was ever encountered in the document.

The vector space model are the documents which are represented as "bags of words". The basic idea is to represent each document as a vector of certain weighted word frequencies. In order to do so, the following parsing and extraction steps are needed.

1. Ignoring case, extract all unique words from the entire set of documents.
2. Eliminate non-content-bearing "stopwords" such as "a", "and", "the", etc. For sample lists of stopwords.
3. For each document, count the number of occurrences of each word.
4. Using heuristic or information-theoretic criteria, eliminate non-content-bearing "high-frequency" and "low-frequency" words.
5. After the above elimination, suppose unique words remain. Assign a unique identifier between 1 and  $n$  to each remaining word, and a unique identifier between 1 and  $m$  to each document.

In vector space model, the documents from which the information is to be retrieved are represented as vectors. The term weighting identifies the success or failure of the vector space method. Terms are basically the words or any indexing unit used to identify the contents of a text. Furthermore, a term weighting scheme plays an important role for the similarity measure. The similarity measures largely identify the retrieval efficiency of a particular information retrieval system.

This largely depends on formulas. Where the formulas depend only on the frequencies within the document and they not depend on inter-document frequencies. The main components of the formulas are as follows:

**Binary:** Binary formula gives every word that appears in a document equal relevance. This can be useful when the number of times a word appears is not considered important.

**Term frequency:** This formula counts how many times the term occurs in a document. The more times a term  $t$  occurs in document  $d$  the more likely it is that  $t$  is relevant to the document. Used alone, favors common words and long documents. This formula gives more credit to words that appears more frequently, but often too much credit.

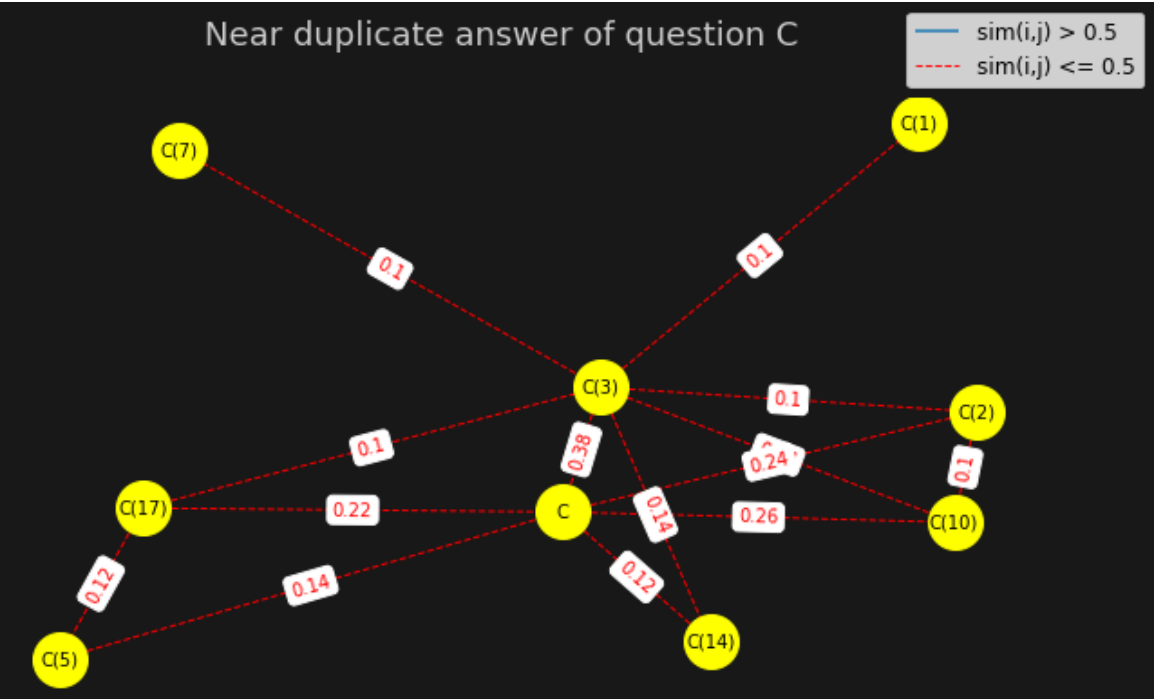
**Augmented normalized term frequency** This formula tries to give credit to any word that appears and then give some additional credit to words that appear frequently.

**Logarithmic term frequency** Logarithms are a way to de-emphasize the effect of frequency. Literature proposes log and alternate log as the most used.

#task: 20

Signatures metric: 20 x 100

#permutations used to create signatures: 100



**Question D:**The result illustrates that there are 3 pairs of answers which can be identified as near duplicate answers as shown below:

- 1. Answer D and D(14) with Jaccard similarity of 0.7
- 2. Answer D and D(18) with Jaccard similarity of 0.68
- 3. Answer D(18) and A(14) with Jaccard similarity of 0.62

In [ ]:

```
LSH_task('D')
```

[In probability theory, Bayes' theorem (often called Bayes' law after Rev Thomas Bayes) relates the conditional and marginal probabilities of two random events. It is often used to compute posterior probabilities given observations. For example, a patient may be observed to have certain symptoms. Bayes' theorem can be used to compute the probability that a proposed diagnosis is correct, given that observation. (See example 2) As a formal theorem, Bayes' theorem is valid in all common interpretations of probability. However, it plays a central role in the debate around the foundations of statistics: frequentist and Bayesian interpretations disagree about the ways in which probabilities should be assigned in applications. Frequentists assign probabilities to random events according to their frequencies of occurrence or to subsets of populations as proportions of the whole, while Bayesians describe probabilities in terms of beliefs and degrees of uncertainty. The articles on Bayesian probability and frequentist probability discuss these debates in greater detail.

Bayes' theorem relates the conditional and marginal probabilities of events A and B, where B has a non-vanishing probability:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Each term in Bayes' theorem has a conventional name:

- $P(A)$  is the prior probability or marginal probability of A. It is "prior" in the sense that it does not take into account any information about B.
- $P(A|B)$  is the conditional probability of A, given B. It is also called the posterior probability because it is derived from or depends upon the specified value of B.
- $P(B|A)$  is the conditional probability of B given A.
- $P(B)$  is the prior or marginal probability of B, and acts as a normalizing constant.

Intuitively, Bayes' theorem in this form describes the way in which one's beliefs about observing 'A' are updated by having observed 'B'.

'Bayes' Theorem' or 'Bayes' Rule', or something called Bayesian reasoning

The Bayesian Conspiracy is a multinational, interdisciplinary, and shadowy group of scientists that controls publication, grants, tenure, and the illicit traffic in grad students. The best way to be accepted into the Bayesian Conspiracy is to join the Campus Crusade for Bayes in high school or college, and gradually work your way up to the inner circles.

Bayes' Theorem

Let  $A$  and  $B$  be sets. Conditional probability requires that

- $A \cap B$  denotes intersection ("and"), and also that
- Therefore,
- Now, let
- so  $E$  is an event in  $\mathcal{E}$  and for  $\omega \in \Omega$ , then
- $P(E|A)$
- But this can be written
- so

This paper proposes a new measure called scaled inverse document frequency (SIDF) which evaluates the conditional specificity of query terms over a subset  $S$  of  $D$  and without making any assumption about term independence.  $S$  can be estimated from search results, OR searches, or computed from inverted index data. We have evaluated SIDF values from commercial search engines by submitting queries relevant to the financial investment domain. Results compare favorably across search engines and queries. Our approach has practical applications for 'real-world' scenarios like in Web Mining, Homeland Security, and keyword-driven marketing research scenarios.

'Bayes' theorem was named after Rev Thomas Bayes and is a method used in probability theory. This theorem aims to relate the conditional and marginal probabilities of two random events occurring, and given various observations is frequently used to compute subsequent probabilities. Bayes' theorem is also often known as Bayes' law.

An example of where Bayes' theorem may be used is in the following extract: "Suppose there exists a school with forty percent females and sixty percent males as students. The female students can only wear skirts or trousers in equal numbers whereas all the male students can only wear trousers. An observer randomly sees a student from a distance and all he can see is that this student is wearing trousers. What is the probability this student is female?"

There is a debate amongst frequentists and Bayesians about how Bayes' theorem plays a major role around the beginnings of statistical mathematics. Frequentist and Bayesian explanations do not agree about the ways in which probabilities should be assigned. This is primarily because Bayesians assign probabilities in terms of beliefs whereas frequentists assign probabilities

bilities to random events according to the frequencies of them occurring.

'Bayes' theorem relates the conditional and marginal probabilities of two random events. For example, a person may be seen to have certain medical symptoms; Bayes' theorem can then be used to compute the probability that, given that observation, the proposed diagnosis is the right one.

Bayes' theorem forms a relationship between the probabilities of events A and B. Intuitively, Bayes' theorem in this form describes the way in which one's recognition of 'A' are updated by having observed 'B'.

$$P(A | B) = P(B | A) P(A) / P(B)$$

$P(A|B)$  is the conditional probability of A given B. It is derived from or depends upon the specified value of B, therefore it is also known as the posterior probability.

$P(B|A)$  is the conditional probability of B given A.

$P(A)$  is the prior probability A. It doesn't take into account any information about B, so it is "prior".

$P(B)$  is the prior or marginal probability of B, and acts to normalise the probability.

To derive the theorem, we begin with the definition of conditional probability. By combining and re-arranging these two equations for A and B, we get a lemma called product rule for probabilities. Provided that  $P(B)$  is not a zero, dividing both sides by  $P(B)$  renders us with Bayes' theorem.

In probability theory; Bayes theorem (often called Bayes law after Rev Thomas Bayes) relates the conditional and marginal probabilities of two random events. It is used to compute posterior probabilities given observations. For example; a person may be observed to have certain symptoms. Bayes theorem can be used to compute the probability that a proposed diagnosis is correct.

As a formal theorem Bayes theorem is valid in all common interpretations of probability. However, it plays a central role in the debate around the foundations of statistics: frequentist and Bayesian interpretations disagree about the ways in which probabilities should be assigned to each other. Bayesians describe probabilities in terms of beliefs and degrees of uncertainty, While frequentists assign probabilities to random events according to their frequencies of occurrence or to subsets of populations as proportions of the whole. The articles on Bayesian probability and frequentist probability discuss these debates in detail.

'Bayes' theorem in connection with conditional probabilities is of fundamental importance, since it permits a calculation of  $PROB(AB)$  from  $PROB(BA)$ . Statistical information that is often gathered in great volume can therefore be avoided.

'Bayes Theorem is an important theorem relating conditional probabilities, it allows us to calculate  $PROB(A|B)$  from  $PROB(B|A)$ . Bayes Theorem is important because it can save us from gathering vast amounts of statistical evidence. The main theory is  $PROB(A|B) = PROB(B|A) * PROB(A) / PROB(B)$ , it means Using  $PROB(WIN|RAIN)$  from earlier, we can find the probability that it rained on a day that Harry won a race.

"Bayes' theorem relates the conditional and marginal probabilities of two random events and is named after the Reverend Thomas Bayes (1702–1761), who studied how to compute a distribution for the parameter of a binomial distribution. It is valid in all common interpretations of probability. It plays a central role in the debate around the foundations of statistics: frequentist and Bayesian interpretations disagree about the ways in which probabilities should be assigned in applications. Frequentists assign probabilities to random events according to their frequencies of occurrence or to subsets of populations as proportions of the whole, while Bayesians describe probabilities in terms of beliefs and degrees of uncertainty. Applications of Bayes' theorem often assume the philosophy underlying Bayesian probability that uncertainty and degrees of belief can be measured as probabilities. One of Bayes' results (Proposition 5) gives a simple description of conditional probability, and shows that it can be expressed independently of the order in which things occur:

If there be two subsequent events, the probability of the second given the probability of both together  $P(B|A)$ , and it being first discovered that the second event has also happened, from hence I guess that the first event has also happened, the probability I am right [i.e., the conditional probability of the first event being true given that the second has also happened] is  $P(A|B)$ .

b. \nNote that the expression says nothing about the order in which the events occurred; it measures correlation, not causation.\n", 'Bayes' theorem relates the conditional and marginal probabilities of two random events. It is mainly used to calculate the probability of one event's outcome given that a previous event happened. For example, the probability that a doctor's diagnosis is correct given that the doctor had previously observed symptoms in the patient. Bayes' theorem can be used for all forms of probability, however it is currently at the centre of a debate concerning the ways in which probabilities should be assigned in applications. \nThe theorem states that the probability of Event A happening given Event B is the probability of B given A multiplied by the probability of A regardless of B all divided by the probability of B regardless of A which acts as a normalising constant. Bayes' theorem formed in this way basically details how one's beliefs about Event A are renewed or updated knowing that Event B happened. When calculating conditional probabilities such as these, it is often useful to create a table containing the number of occurrences, or relative frequencies, of each outcome for each of the variables independently.\n\n", 'Bayes Theorem is a mathematical formula used to calculate conditional probabilities. Given the probability of event A given event B, Bayes Theorem can be used to calculate the probability of B given A. This is achieved using the conditional probability of B given A and the prior probabilities of both events A and B. For example: suppose there is a bag of coloured balls with 25 red ones and 75 black ones. Lucky Joe likes to predict the colour of the ball he selects and he is 80% accurate. Joe records all of his results and about 0.5% of the time he accidentally records the wrong results. Using all of this information more probabilities can be inferred, including using Bayes Theorem to calculate various probabilities like Joe recording correctly if he guesses correctly or Joe recording incorrectly when his guess was correct (and other like combinations). \n", " In probability theory, Bayes' theorem (often called Bayes' law after Rev Thomas Bayes) relates the conditional and marginal probabilities of two random events. It is often used to compute posterior probabilities given observations. For example, a patient may be observed to have certain symptoms. Bayes' theorem can be used to compute the probability that a proposed diagnosis is correct, given that observation.\nAs a formal theorem, Bayes' theorem is valid in all common interpretations of probability. However, it plays a central role in the debate around the foundations of statistics: frequentist and Bayesian interpretations disagree about the ways in which probabilities should be assigned in applications. \nSuppose there is a co-ed school having 60% boys and 40% girls as students. The girl students wear trousers or skirts in equal numbers; the boys all wear trousers. An observer sees a (random) student from a distance; all they can see is that this student is wearing trousers. What is the probability this student is a girl? The correct answer can be computed using Bayes' theorem.\nThe event A is that the student observed is a girl, and the event B is that the student observed is wearing trousers. To compute  $P(A|B)$ , we first need to know:  $P(B|A')$ , or the probability of the student wearing trousers given that the student is a boy. This is given as  $1 \cdot P(A)$ , or the probability that the student is a girl regardless of any other information. Since the observer sees a random student, meaning that all students have the same probability of being observed, and the fraction of girls among the students is 40%, this probability equals  $0.4 \cdot P(A')$ , or the probability that the student is a boy regardless of any other information ( $A'$  is the complementary event to A). This is 60%, or  $0.6 \cdot P(B|A)$ , or the probability of the student wearing trousers given that the student is a girl. As they are as likely to wear skirts as trousers, this is 0.5", '\nIn probability theory, Bayes\' theorem also called Bayes\' law after Rev Thomas Bayes compares the conditional and marginal probabilities of two random events. It is often used to calculate posterior probabilities given observations. For example, a patient may be observed to have certain symptoms. Bayes\' theorem can be used to calculate the likelihood that a proposed analysis is accurate, given that observatio



n. \nAs an official theorem, Bayes\' theorem is valid in all universal interpretations of probability. However, it plays a fundamental role in the debate around the foundations of statistics: frequentist and Bayesian interpretations disagree about the ways in which probabilities should be assigned in applications.\nFrequentists assign probabilities to random events according to their frequencies of happening or to subsets of populations as proportions of the whole. Whilst Bayesians describe probabilities in terms of beliefs and degrees of uncertainty. The articles on Bayesian probability and frequentist probability discuss these debates in greater detail.\nBayes\' theorem compares the conditional and marginal probabilities of events A and B, where B has a non-vanishing probability.\nEach term in Bayes\' theorem has a conventional name:\nP(A) is the previous probability of A. It is "previous" in the sense that it does not take into account any information about B.\nP(A|B) is the conditional probability of A, given B. It is also called the subsequent probability because it is derived from or depends upon the specified value of B.\nP(B|A) is the conditional probability of B given A.\nP(B) is the previous.\n', "In probability theory, Bayes' theorem relates the conditional and marginal probabilities of two random events. It is often used to compute posterior probabilities given observations.\n\nBayes' theorem is expressed mathematically as:\n
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$
\n\nwhere P(A|B) is the conditional probability of A given B, P(A) is the prior probability of A, P(B) is the prior probability of B, and P(B|A) is the conditional probability of B given A. \n\nBayes' theorem relates the conditional and marginal probabilities of two random events P(A) and P(B), and is valid in all common interpretations of probability. For example, in a school made up of 3/5 boys and 2/5 girls, the girls wear skirts or trousers in equal numbers and the boys all wear trousers. If a student is observed from a distance wearing trousers, Bayes theorem can be used to determine the probability of this student being a girl. \n\nP(A) is the probability of the student being a girl (which is 2/5). P(B|A) is the probability of the student wearing trousers given that the student is a girl, which is 0.5 P(B) is the probability of a random student wearing trousers, which can be calculated as  $P(B) = P(B|A)P(A) + P(B|A')P(A')$  where ' denotes a complementary event, which is 0.8. \nTherefore the probability using the formula is 0.25. \n\nBayes theorem is often used to compute posterior probabilities given observations, for instance the probability that a proposed medical diagnosis is correct, given certain observed symptoms. \n", 'The Probability of an event happening mean considering the likelihood of or the number of the instance occurring, and dividing this value by the total number of events. The equation for this calculation would look as follows:\n
$$Probability(P) = \frac{\text{number of instance}}{\text{total number of events}}$$
\n\nthe other hand Probability Theory (P) usually involves assigning values to events. For example:\nP=1: event is certain to occur\nP=0: event is certain NOT to occur\nP=0.5: event occurs half of the time. \n\nThere is also Conditional Probability which is usually interested in the way variables relate to each other. Bayes Theorem is the name given to an important theorem relating\nConditional probabilities and it can be seen as a way of understanding how the probability that a theory is true, is affected by a new piece of evidence. It has been used in a wide variety of contexts, ranging from marine biology to the development of "Bayesian" spam blockers for email systems. \n', 'In probability theory, Bayes\' theorem (often called Bayes\' law after Rev Thomas Bayes) relates the conditional and marginal probabilities of two random events. It is often used to compute posterior probabilities given observations (for example, a patient may be observed to have certain symptoms). Bayes\' theorem can be used to compute the probability that a proposed diagnosis is correct, given that observation.\n\nAs a formal theorem, Bayes\' theorem is valid in all common interpretations of probability. However, it plays a central role in the debate around the foundations of statistics; frequentist and Bayesian interpretations disagree about the ways in which probabilities should be assigned in applications. Frequentists assign probabilities to random events according to their f

frequencies of occurrence or to subsets of populations as proportions of the whole, while Bayesians describe probabilities in terms of beliefs and degrees of uncertainty. The articles on Bayesian probability and frequentist probability discuss these debates in greater detail.

Bayes' theorem relates the conditional and marginal probabilities of events A and B, where B has a non-vanishing probability:
 
$$P(A|B) = (P(B|A) \times P(A)) / P(B)$$
 Each term in Bayes' theorem has a conventional name:

- $P(A)$  is the prior probability or marginal probability of A. It is "prior" in the sense that it does not take into account any information about B.
- $P(A|B)$  is the conditional probability of A, given B. It is also called the posterior probability because it is derived from or depends upon the specified value of B.
- $P(B|A)$  is the conditional probability of B given A.
- $P(B)$  is the prior or marginal probability of B, and acts as a normalizing constant.

Intuitively, Bayes' theorem in this form describes the way in which one's beliefs about observing 'A' are updated by having observed 'B'. Bayes' theorem (often called Bayes' law) connects the conditional and marginal probabilities of two arbitrary events. One of its uses is calculating posterior probabilities given observations.

Bayes' theorem plays a key role in the debate around the principles of statistics: frequentist and Bayesian interpretations disagree about the ways in which probabilities should be assigned in applications.

Bayes' theorem is useful in evaluating the result of drug tests. If a test can identify a drug user 99% of the time, and can identify a non-user as testing negative 99% of the time, it may seem to be a relatively accurate test. However, Bayes' theorem will reveal the flaw that despite the apparently high accuracy of the test, the probability that an employee who tested positive actually did use drugs is only about 33%.

In probability theory, the prior and conditional probabilities of two random events are related by Bayes' theorem. The theorem is often used when we have observations and wish to compute posterior probabilities.

For example, given an observation that a patient is seen to have certain symptoms, we can use Bayes' theorem to compute the probability that a suggested diagnosis is correct.

$P(A)$  is the prior probability of A.  $P(A|B)$  is the conditional probability of A given B.  $P(B|A)$  is the conditional probability of B given A.  $P(B)$  is the prior probability of B, and must be non-zero. Bayes' theorem is given by
 
$$P(A|B) = (P(B|A)P(A)) / (P(B))$$
 In probability theory, Bayes' theorem (or Bayes' law after Rev Thomas Bayes) provides relation between the conditional and marginal probabilities of two random events. It is usually used to calculate posterior probabilities given observations. For example: a patient might be observed to show certain symptoms. Bayes' theorem could be used to compute the probability that a certain diagnosis is right, given that observation.

Since it is a formal theorem, Bayes' theorem holds in all popular interpretations of probability.

Bayes' theorem relates the conditional and marginal probabilities of events a and b, where b has a non-vanishing probability:
 
$$P(a|b) = P(a)bP(a)/P(b)$$
 Terms in Bayes' theorem are named by a convention:

- $P(A)$  is the prior probability or marginal probability of A. It does not take into account any information about B and therefore is considered "prior".
- $P(A|B)$  is the conditional probability of A, given B. It is derived from or depends upon the specified value of B. Usually it is called the posterior probability.
- $P(B|A)$  is the conditional probability of B given A.
- $P(B)$  (a.k.a. the normalizing constant) is the prior or marginal probability of B.

Obviously, Bayes' theorem describes the way in which one's assumptions about observing the event 'a' are changed by having observed the event 'b'.

In probability theory, Bayes' theorem relates the conditional and marginal probabilities of two random events. It is usually used to compute posterior probabilities given observations. For instance, a patient may be observed to have certain symptoms. Bayes' theorem can be used to compute the probability that a proposed diagnosis is correct, given that observation.

As a formal theorem, Bayes' theorem is valid in all common interpretations of probability. However, it plays a central role in the debate around the foundations of statistics: frequentist and Bayesian

sian interpretations disagree about the ways in which probabilities should be assigned in applications. The articles on Bayesian probability and frequentist probability discuss these debates in greater detail. Frequentists assign probabilities to random events according to their frequencies of occurrence or to subsets of populations as proportions of the whole. At the same time, Bayesians describe probabilities in terms of beliefs and degrees of uncertainty.

Bayes' theorem relates the conditional and marginal probabilities of events A and B, where B has a non-vanishing probability:

Each term in Bayes' theorem has a conventional name:

- $P(A)$  is the prior probability or marginal probability of A. It is "prior" in the sense that it does not take into account any information about B.
- $P(A|B)$  is the conditional probability of A, given B. It is also called the posterior probability because it is derived from or depends upon the specified value of B.
- $P(B|A)$  is the conditional probability of B given A.
- $P(B)$  is the prior or marginal probability of B, and acts as a normalizing constant.

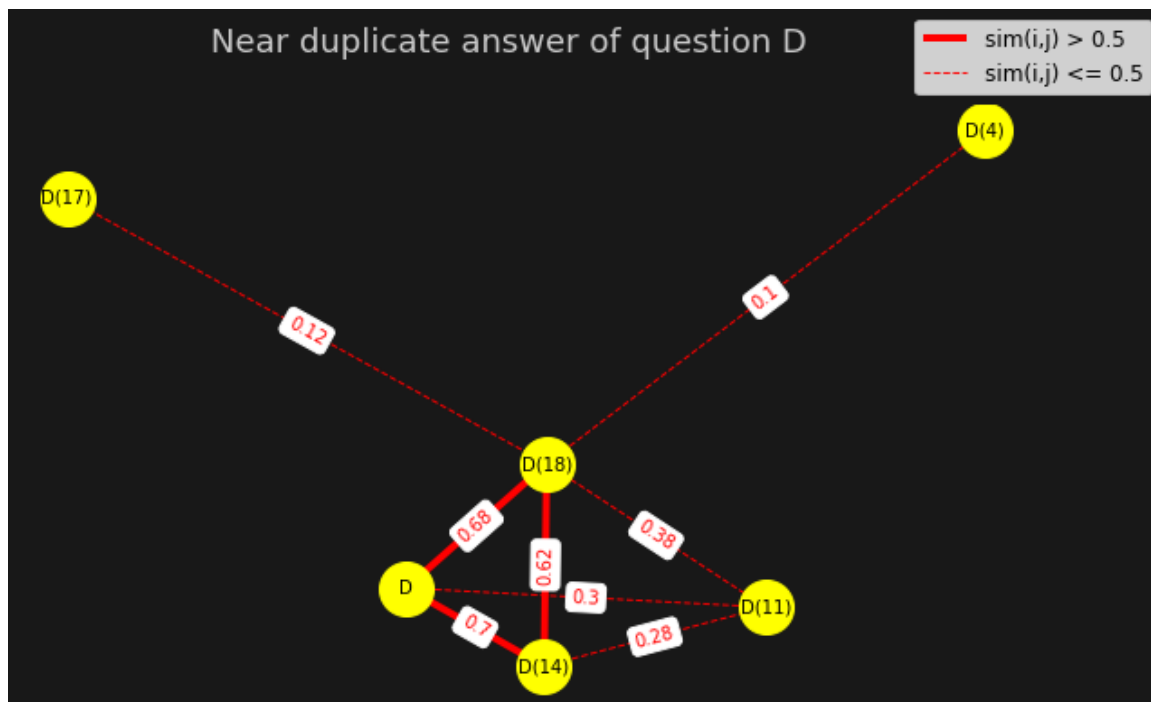
Intuitively, Bayes' theorem in this form describes the way in which one's beliefs about observing 'A' are updated by having observed 'B'.

Bayes' Theorem is a simple mathematical formula used for calculating conditional probabilities. Bayes' Theorem is a theorem of probability theory originally stated by the Reverend Thomas Bayes. It figures prominently in subjectivist or Bayesian approaches to epistemology, statistics, and inductive logic. It can be seen as a way of understanding how the probability that a theory is true is affected by a new piece of evidence. It has been used in a wide variety of contexts, ranging from marine biology to the development of "Bayesian" spam blockers for email systems. In the philosophy of science, it has been used to try to clarify the relationship between theory and evidence. Many insights in the philosophy of science involving confirmation, falsification, the relation between science and pseudoscience, and other topics can be made more precise, and sometimes extended or corrected, by using Bayes' Theorem. Subjectivists, who maintain that rational belief is governed by the laws of probability, lean heavily on conditional probabilities in their theories of evidence and their models of empirical learning. Bayes' Theorem is central to these enterprises both because it simplifies the calculation of conditional probabilities and because it clarifies significant features of subjectivist position. Indeed, the Theorem's central insight – that a hypothesis is confirmed by any body of data that its truth renders probable – is the cornerstone of all subjectivist methodology. ']

#task: 20

Signatures metric: 20 x 100

#permutations used to create signatures: 100



**Question E:** The result illustrates that there are no pairs of answers which can be identified as near duplicate answers

In [ ]:

```
LSH_task('E')
```

['In mathematics and computer science, dynamic programming is a method of solving problems that exhibit the properties of overlapping subproblems and optimal substructure (described below). The method takes much less time than naive methods.\n

The term was originally used in the 1940s by Richard Bellman to describe the process of solving problems where one needs to find the best decisions one after another. By 1953, he had refined this to the modern meaning. The field was founded as a systems analysis and engineering topic that is recognized by the IEEE. Bellman's contribution is remembered in the name of the Bellman equation, a central result of dynamic programming which restates an optimization problem in recursive form.\n

The word "programming" in "dynamic programming" has no particular connection to computer programming at all, and instead comes from the term "mathematical programming", a synonym for optimization. Thus, the "program" is the optimal plan for action that is produced. For instance, a finalized schedule of events at an exhibition is sometimes called a program. Programming, in this sense, means finding an acceptable plan of action, an algorithm.\n

Optimal substructure means that optimal solutions of subproblems can be used to find the optimal solutions of the overall problem. For example, the shortest path to a goal from a vertex in a graph can be found by first computing the shortest path to the goal from all adjacent vertices, and then using this to pick the best overall path, as shown in Figure 1. In general, we can solve a problem with optimal substructure using a three-step process:\n

1. Break the problem into smaller subproblems.\n
2. Solve these problems optimally using this three-step process recursively.\n
3. Use these optimal solutions to construct an optimal solution for the original problem.\n

The subproblems are, themselves, solved by dividing them into sub-subproblems, and so on, until we reach some simple case that is solvable in constant time.\n

Figure 2. The subproblem graph for the Fibonacci sequence. That it is not a tree but a DAG indicates overlapping subproblems.\n

To say that a problem has overlapping subproblems is to say that the same subproblems are used to solve many different larger problems. For example, in the Fibonacci sequence,  $F_3 = F_1 + F_2$  and  $F_4 = F_2 + F_3$  – computing each number involves computing  $F_2$ . Because both  $F_3$  and  $F_4$  are needed to compute  $F_5$ , a naive approach to computing  $F_5$  may end up computing  $F_2$  twice or more. This applies whenever overlapping subproblems are present: a naive approach may waste time recomputing optimal solutions to subproblems it has already solved.\n

In order to avoid this, we instead save the solutions to problems we have already solved. Then, if we need to solve the same problem later, we can retrieve and reuse our already-computed solution. This approach is called memoization (not memorization, although this term also fits). If we are sure we won't need a particular solution anymore, we can throw it away to save space. In some cases, we can even compute the solutions to subproblems we know that we'll need in advance.\n

', 'Dynamic programming is a method of providing solutions to potential problems exhibiting the properties of overlapping subproblems and optimal structure. This is highly used in dynamic programming. The advantage being the less time consumption in comparison to other amateur methods.\n

It has to be kept in mind that the term programming in the field has got nothing to do with computer programming at all. On the other hand it is derived from the term mathematical programming which is a similar word used for optimization. Here by meaning that a program can be an optimal plan for the produced action. The typical example could be of a finalized schedule of events at an exhibition. This leads to the concept of programming being a helper in finding an acceptable plan of action, which can also be termed as an algorithm\n

The subproblems are, themselves, solved by dividing them into sub-subproblems, and so on, until we reach some simple case that is solvable in constant time.\n

Overlapping subproblems means that the same subproblems are used to solve many different larger problems. Example could be of Fibonacci sequence;  $F_3 = F_1 + F_2$  and  $F_4 = F_2 + F_3$  – computing each number involves computing  $F_2$ . Because both  $F_3$  and  $F_4$  are needed to compute  $F_5$ , a naive approach to computing  $F_5$  may end up computing  $F_2$  twice or more. It means that whenever we encounter with

overlapping subproblems, a naive approach may waste time recomputing optimal solutions to the already solved subproblems.

'Dynamic programming is a method for efficiently solving a broad range of search and optimization problems which exhibit the characteristics of overlapping. Dynamic programming. Design technique, like divide-and-conquer method.'

The leading and most up-to-date textbook on the far-ranging algorithmic methodology of Dynamic Programming, which can be used for optimal control, ...

The word Programming in the name has nothing to do with writing computer programs. Mathematicians use the word to describe a set of rules which anyone can follow to solve a problem. They do not have to be written in a computer language.

Dynamic programming was the brainchild of an American Mathematician, Richard Bellman, who described the way of solving problems where you need to find the best decisions one after another. In the forty-odd years since this development, the number of uses and applications of dynamic programming has increased enormously.

For example, in 1982 David Kohler used dynamic programming to analyse the best way to play the game of darts.

In recent years, dynamic programming languages develop very fast, especially PHP and Ruby. There is no doubt that they have already become the first choice for many programmers when developing web applications. When you learn a new natural language and you start to use it you naturally, you find yourself using new concepts and paradigms that enrich the use of the language you already know; expect the same result with computer languages.

'Dynamic Programming is an algorithm design technique used for optimisation problems, such as minimising or maximising. Like divide and conquer, Dynamic Programming solves problems by combining solutions to sub-problems. However, unlike divide and conquer, sub-problems are not always independent as sub-problems may share sub-sub-problems but solution to one sub-problem may not affect the solutions to other sub-problems of the same problem.'

There are four steps in Dynamic Programming:

1. Characterise structure of an optimal solution.
2. Define value of optimal solution recursively.
3. Compute optimal solution values either top-down with caching or bottom-up in a table.
4. Construct an optimal solution from computed values.

An example of the type of problem for which Dynamic Programming may be used is: given two sequences,  $X=(x_1, \dots, x_m)$  and  $Y=(y_1, \dots, y_n)$  find a common subsequence whose length is maximum.

Dynamic Programming reduces computation by solving sub-problems in a bottom-up fashion and by storing solution to a sub-problem the first time it is solved. Also, looking up the solution when a sub-problem is encountered again helps reduce computation. However, the key in Dynamic Programming is to determine the structure of optimal solutions.

'Dynamic programming is a method for solving mathematical programming problems that exhibit the properties of overlapping subproblems and optimal substructure. This is a much quicker method than other more naive methods. The word "programming" in "dynamic programming" relates optimization, which is commonly referred to as mathematical programming. Richard Bellman originally coined the term in the 1940s to describe a method for solving problems where one needs to find the best decisions one after another, and by 1953, he refined his method to the current modern meaning.

Optimal substructure means that by splitting the programming into optimal solutions of subproblems, these can then be used to find the optimal solutions of the overall problem. One example is the computing of the shortest path to a goal from a vertex in a graph. First, compute the shortest path to the goal from all adjacent vertices. Then, using this, the best overall path can be found, thereby demonstrating the dynamic programming principle. This general three-step process can be used to solve a problem:

1. Break up the problem into different smaller subproblems.
2. Recursively use this three-step process to compute the optimal path in the subproblem.
3. Construct an optimal solution, using the computed optimal subproblems, for the original problem.

This process continues recursively, working over the subproblems by dividing them into sub-subproblems and so forth, until a simple case is reached (one that is easily solvable).

'In computer science; dynamic programming is a way of solving

ing problems consist of overlapping subproblems and optimal substructure. The method is more efficient than naive methods. The term was first coined in the 1940s by Richard Bellman to describe the process of solving problems where you need to find the best decisions consecutively. In 1953 he had refined this to the modern meaning. The field was founded as a systems analysis and engineering topic that is recognized by the IEEE. Bellman equation is a central result of dynamic programming which restates an optimization problem in recursive form. Dynamic programming has little connection to computer programming at all, and instead comes from the term mathematical programming, a synonym for optimization. Thus, the program is the best plan for action that is produced. For instance, an events schedule at an exhibition is sometimes called a program. Programming means finding a plan of action. Dynamic programming (DP) is an extremely powerful, general tool for solving optimization difficulties on left-right-ordered item, for example character strings. It is similar to divide and conquer, however is differentiated as its subproblems are not independent. It is easily applicable, in relative terms, once understood. However until one has witnessed enough examples, it looks like magic. DP minimizes computation by solving subproblems from the base upwards, storing solution to a subproblem when it is initially conquered, and looking up the solution when the subproblem is experienced for a second time. Dynamic programming is a method of solving problems that exhibit the properties of overlapping subproblems and optimal substructure (described below). The method takes much less time than naive methods. The word "programming" in "dynamic programming" has no particular connection to computer programming at all, and instead comes from the term "mathematical programming", a synonym for optimization. Thus, the "program" is the optimal plan for action that is produced. For instance, a finalized schedule of events at an exhibition is sometimes called a program. Programming, in this sense, means finding an acceptable plan of action, an algorithm. Dynamic programming is an algorithmic technique used to solve certain optimization problems where the object is to find the best solution from a number of possibilities. It uses a so called 'bottom-up' approach, meaning that the problem is solved as a set of sub-problems which in turn are made up of sub-sub-problems. Sub-problems are then selected and used to solve the overall problem. These sub-problems are only solved once and the solutions are saved so that they will not need to be recalculated again. Whilst calculated individually, they may also overlap. When any sub-problem is met again, it can be found and re-used to solve another problem. Since it searches all possibilities, it is also very accurate. This method is far more efficient than recalculating and therefore considerably reduces computation. It is widely used in computer science and can be applied for example, to compress data in high density bar codes. Dynamic programming is most effective and therefore most often used on objects that are ordered from left to right and whose order cannot be rearranged. This means it works well on character chains for example. In mathematics and computer science, dynamic programming is a method of solving problems that exhibit the properties of overlapping subproblems and optimal substructure. The term was originally used in the 1940s by Richard Bellman to describe the process of solving problems where one needs to find the best decisions one after another. By 1953, he had refined this to the modern meaning. Bellman's contribution is remembered in the name of the Bellman equation, a central result of dynamic programming which restates an optimization problem in recursive form. The word "programming" in "dynamic programming" has no particular connection to computer programming at all, and instead comes from the term "mathematical programming", a synonym for optimization. Thus, the "program" is the optimal plan for action that is produced. For instance, a finalized schedule of events at an exhibition is sometimes called a program. Programming, in this sense, means finding an acceptable plan of action, an algorithm. Dynamic programming usually takes one of two approaches, the top-down approach, the problem is broken into sub problems, and these sub problems are solved and the solution



ons remembered, in case they need to be solved again. This is recursion and memorization combined together and the bottom-up approach, all sub problems that might be needed are solved in advance and then used to build up solutions to larger problems. This approach is slightly better in stack space and number of function calls, but it is sometimes not intuitive to figure out all the sub problems needed for solving the given problem.

Some programming languages can automatically memorize the result of a function call with a particular set of arguments, in order to speed up call-by-name. Some languages make it possible portably (e.g. Scheme, Common Lisp or Perl), some need special extensions. This is only possible for a referentially transparent function.

'Dynamic programming is a faster method of solving problems that make use of optimal substructure, overlapping sub-problems and memoization. It has no relationship to computer programming; instead it is a process of finding a satisfactory algorithm.

Optimal substructure is the process of using the optional solutions to sub problems to find the optimal solution to the overall problem. When the same sub problem solutions can be used to solve various bigger problems it is said to have overlapping-sub problems. Memoization is used in order to save time the solutions are stored rather than be recomputed. A solution can be disposed of once we are positive that it will no longer be required, in some cases a solution to a future problem can be computed in advance.

There are two main approaches for dynamic programming. The first is the bottom up approach. Although it is not always simple to find all of them, any required sub problems are solved in advance and then used to create solutions to larger problems. The other method is the top down approach which is a method that combines memorization and recursion. The main problem is divided into sub problems which are solved and stored for future use.

'Dynamic Programming is a very powerful mathematical technique, often utilised in programming, for solving optimization problems. Normally, minimizing or maximizing.

Greedy algorithms focus on making the best local choice at each decision making stage. Without a proof of correctness, such an algorithm is likely to fail. With Dynamic Programming, we can design our own algorithm which searches for all possibilities (which ensures correctness) whilst storing the results to avoid having to recomputed (leading to computational efficiency).

Dynamic Programming solves problems by combining the solutions of subproblems. These subproblems are not, however, independent. Subproblems can share subsubproblems, but the solution to one subproblem doesn't necessarily affect the solutions to other subproblems stemming from the same problem.

Dynamic Programming reduces computation time by solving subproblems in a bottom-up way. It stores the solution to a subproblem the first time it is solved, meaning that it can look up the solution when that subproblem is encountered subsequently.

The key to Dynamic Programming is to find the structure of optimal solutions. The steps required are as follows:

1. Generalise the structure of an optimal solution
2. Recursively define the value of an optimal solution
3. Compute the optimal solution values either top-down (with caching), or bottom-up using a table
4. Generate the optimal solution of these computed values

'In mathematics and computer science, dynamic programming is a method of solving problems, that exhibit the properties of overlapping subproblems and optimal substructure. The method takes much less time than naive methods.

The term was originally used in the 1940s to describe the process of solving problems where one needs to find the best decisions one after another.

The field was founded as a systems analysis and engineering topic that is recognized by the IEEE

The word "programming" in "dynamic programming" has no particular connection to computer programming at all, and instead comes from the term "mathematical programming", a synonym for optimization. Thus, the "program" is the optimal plan for action that is produced. For instance, a finalized schedule of events at an exhibition is sometimes called a program. Programming, in this sense, means finding an acceptable plan of action, an algorithm.

Optimal substructure means that optimal solutions of subproblems can be used to find the optimal solutions of the overall

problem. For example, the shortest path to a goal from a vertex in a graph can be found by first computing the shortest path to the goal from all adjacent vertices, and then using this to pick the best overall path.

In general, we can solve a problem with optimal substructure using a three-step process:

1. Break the problem into smaller subproblems.
2. solve these problems optimally using this three-step process recursively.
3. Use these optimal solutions to construct an optimal solution for the original problem.

The subproblems are, themselves, solved by dividing them into sub-subproblems, and so on, until we reach some simple case that is solvable in constant time.

'Dynamic programming is a problem-solving method which solves recursive problems. The term is derived from mathematical programming which is commonly referred to as optimisation, hence dynamic programming is an optimal method of solving the problems and takes much less time than naïve methods.

Dynamic programming uses the properties of optimal substructure, overlapping subproblems and memoization to create an algorithm to solve such problems. Optimal substructure means that the structure of the problem is made up of sub-problems which can be used to find the solution to the problem overall. A problem with overlapping subproblems means that the same subproblems may be used to solve many different larger problems. Each sub-problem is solved by being divided into sub-subproblems, until a case is reached which is solvable in constant time. Memoization stores solutions which have already been computed in order to reduce unnecessary re-computation.

Dynamic programming can be divided into two main approaches: top-down and bottom-up. The top-down approach breaks the problem into subproblems, which are solved and remembered, using a combination of memoization and recursion. The bottom-up approach solves all subproblems that might be need in advance, and then uses these solutions to build up the solutions to the bigger problem.

'Dynamic Programming (DP) is in basic terms an algorithm design technique that is used for optimization problems and often involves minimizing or maximizing.

Furthermore, by combining solutions to subproblems, DP solves problems. Subproblems may include and contain many other subsubproblems and even in such cases, the solution to one subproblem may not affect the solutions to other subproblems involved in the same problem.

By solving subproblems in a bottom-up fashion, which is basically when storing solution to a subproblem the first time it is solved and looking up to find the solution when a subproblem is come across once more, this would cause DP to reduce computations.

The following is a generalization path to be taken in Dynamic Programming:

- Firstly it is needed to Characterize the structure of an optimal solution.
- Secondly to define the value of the optimal solution recursively.
- Furthermore, to compute the optimal solution values either by following a top-down method with caching, or a bottom-up method in a table.
- The last point would be to construct an optimal solution from the computed values.

"In the field of computer science, term 'dynamic programming' relates to the style of programming that breaks a large problem down into smaller subproblems, and generally allows for the finding of the optimal solution. When the problem is split into subproblems, these themselves may be split into smaller problems, and so on, until they cannot be reduced any more.

It is also common for dynamic programming to make use of recursion, and the saving of previous results for faster computation later; this also leads to higher efficiency, as calculations are not being redone. For example, when a problem is reduced into sub problems, and those are then reduced further, it may be that there are common subsubproblems, and so only one calculation needs to be done and the result saved to help solve more than one subproblem.

An example of this gain in efficiency is a path-finding problem. If there are two distinct routes in a network of 10 nodes, tagged A to J, then if the two routes share a common section (say, between nodes B and D), the cost of that section should be calculated for the first route and saved. Then, when the second route is being processed, the cost of B to D does not need to be calculated again.

In general, dynamic programming is used on optimisation problems, where the most efficient solution is needed. Areas where this so

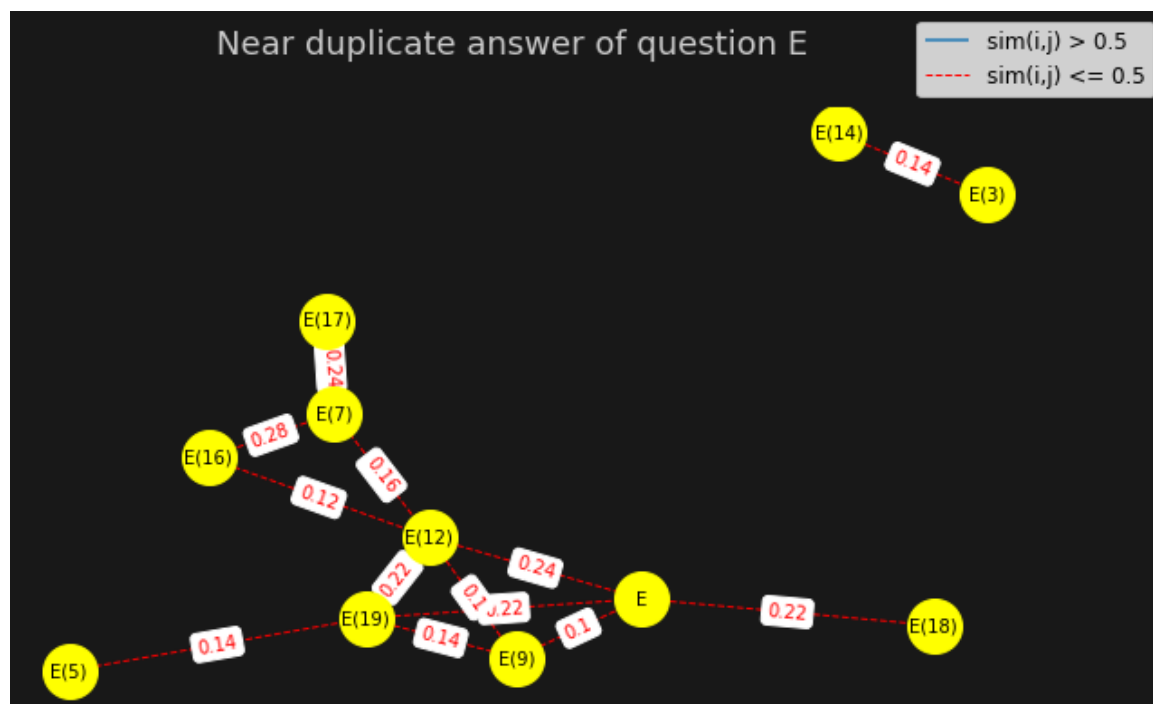
rt of programming is useful is in AI, computer graphics, compression routines, and biomedical applications.", 'Dynamic Programming is a method of solving problems that exhibit the properties of overlapping subproblems and optimal substructure. The term was originally used in the 1940s by Richard Bellman.\n\nThe word "programming" in "dynamic programming" has no particular connection to computer programming at all, and instead comes from the term "mathematical programming", a synonym for optimization. The "program" is the optimal plan for action that is produced.\n\nFor instance, a finalized schedule of events at an exhibition is sometimes called a program. Programming, in this sense, means finding an acceptable plan of action. \n\nTo say that a problem has overlapping subproblems is to say that the same subproblems are used to solve many different larger problems. Optimal substructure means that optimal solutions of subproblems can be used to find the optimal solutions of the overall problem.', 'In computer science and mathematics, dynamic programming\nis a method of problem solving that utilizes the properties\nof overlapping subproblems and optimal substructure. And thus\nthe method takes much less time than more naive methods.\n\nIn "dynamic programming", the word "programming" has no\nreal connection to computer programming at all, it actually\ncomes from the term "mathematical programming",\na synonym for optimisation. Thus, the "program" is the optimal\nplan of action that is being produced. For example, a\nschedule of events at an exhibition is sometimes called a\nprogramme. Programming, in this sense, means finding an\nacceptable plan, an algorithm.', 'In mathematics and computer science, dynamic programming is a method of solving problems that exhibit the properties of overlapping subproblems and optimal substructure.\n\nThe word "programming" in "dynamic programming" has no particular connection to computer programming at all, and instead comes from the term "mathematical programming", a synonym for optimization. Programming, in this sense, means finding an acceptable plan of action, an algorithm.\n\nOptimal substructure means that optimal solutions of subproblems can be used to find the optimal solutions of the overall problem. In general, we can solve a problem with optimal substructure using a three-step process:\n\n1. Break the problem into smaller subproblems.\n2. Solve these problems optimally using this three-step process recursively.\n3. Use these optimal solutions to construct an optimal solution for the original problem.\n\nThe subproblems are, themselves, solved by dividing them into sub-subproblems, and so on, until we reach some simple case that is solvable in constant time.\n\nTo say that a problem has overlapping subproblems is to say that the same subproblems are used to solve many different larger problems. For example, in the Fibonacci sequence,  $F_3 = F_1 + F_2$  and  $F_4 = F_2 + F_3$  – computing each number involves computing  $F_2$ . Because both  $F_3$  and  $F_4$  are needed to compute  $F_5$ , a naive approach to computing  $F_5$  may end up computing  $F_2$  twice or more. This applies whenever overlapping subproblems are present: a naive approach may waste time recomputing optimal solutions to subproblems it has already solved.\n\nIn order to avoid this, we instead save the solutions to problems we have already solved. Then, if we need to solve the same problem later, we can retrieve and reuse our already-computed solution. If we are sure we won't need a particular solution anymore, we can throw it away to save space. In some cases, we can even compute the solutions to subproblems we know that we'll need in advance.\n\nDynamic programming makes use of:\n\nOverlapping subproblems\nOptimal substructure\nMemoization\n\nDynamic programming usually takes one of two approaches:\n\nTop-down approach\nBottom-up approach\n\n', 'In mathematics and computer science, dynamic programming is a methodology of the solution of the problems that exhibit the properties of overlapping subproblems and optimal substructure (described below). The methodology takes much less time rather than naive methods.\n\nThe term was originally used during the 1940s by Richard Bellman to describe the process of solving problems where one needs to find the best decisions one after another. By 1953, he had refined this to the modern meaning. The field was founded as a systems analysis and engineering topic that is recognized by the IEEE. Bellman

\ 's contribution is remembered in the name of the Bellman equation, a central result of dynamic programming, which restates an optimization problem in a recursive form. The word "programming" in "dynamic programming" has no particular connection to computer programming in general, and instead of this it comes from the term "mathematical programming", a synonym for optimization. Therefore, the "program" is the optimal plan for action that is produced. For example, a finalized schedule of events at an exhibition is sometimes called a program. Optimal substructure means that optimal solutions of subproblems can be used to find the optimal solutions of the overall problem. For instance, the shortest path to a goal from a vertex in a graph can be found by first computing the shortest path to the goal from all adjacent vertices. After this, it is using this to pick the best overall path. In a word, we can solve a problem with optimal substructure using a three-step process.

#task: 20

Signatures metric: 20 x 100

#permutations used to create signatures: 100



## 2. Package: datasketch // Dataset: News headlines

We will implement MinHash LSH Forest, which takes a MinHash data sketch of the query set, and returns the top-k matching sets that have the approximately highest Jaccard similarities with the query set, to create news headlines recommendation

**Step 1: Install and Import packages**

In [ ]:

```
pip install datasketch
```

Requirement already satisfied: datasketch in /usr/local/lib/python3.6/dist-packages (1.5.3)

Requirement already satisfied: numpy>=1.11 in /usr/local/lib/python3.6/dist-packages (from datasketch) (1.19.5)

In [ ]:

```
import numpy as np
import pandas as pd
import re
import time
from datasketch import MinHash, MinHashLSHForest
```

**Step 2: Acquire data using web scraping package (BeautifulSoup)**

In [ ]:

```
from bs4 import BeautifulSoup
import requests

url1 = 'https://www.reuters.com/news/archive/technologynews?view=page&page=1&pageSize=10'
url2 = 'https://www.reuters.com/news/archive/technologynews?view=page&page=2&pageSize=10'
url3 = 'https://www.reuters.com/news/archive/technologynews?view=page&page=3&pageSize=10'
url4 = 'https://www.reuters.com/news/archive/technologynews?view=page&page=4&pageSize=10'
url5 = 'https://www.reuters.com/news/archive/technologynews?view=page&page=5&pageSize=10'

def web_scraping_news(url):
    r = requests.get(url)
    html = r.text
    soup = BeautifulSoup(html, 'lxml')
    div_tag = soup.find_all('h3', attrs={'class': "story-title"})
    news = [i for i in div_tag]

    return news
```

In [ ]:

```
news1 = web_scraping_news(url1)
news2 = web_scraping_news(url2)
news3 = web_scraping_news(url3)
news4 = web_scraping_news(url4)
news5 = web_scraping_news(url5)
news = news1 + news2 + news3 + news4 + news5
titles = [str(i) for i in news]

#Extract News Headline
headlines = [a.strip(''''<h3 class="story-title">
                                     '').rstrip('</') for a
in titles]

for headline in headlines:
    print(headline)
```

U.S. thanks Taiwan for support for auto chips in key trade meeting  
Robinhood lifts trading restrictions on all stocks, including GameStop  
PayPal says to shut domestic payments business in India  
Shares of Tencent-backed Kuaishou triple in HK debut as retail frenzy continues

Australian drone firm reshapes strategy over Google pull-out threat  
Facebook faces a reckoning in Myanmar after blocked by military  
GameStop, 'Reddit rally' stocks slide more, Yellen vows scrutiny  
Yellen seeks to 'understand deeply' GameStop frenzy as market regulators meet

In GameStop saga, U.S. regulator examining all aspects and parties: sources

Samsung considers Austin for \$17 billion chip plant, eyes tax breaks of at least \$806 million: documents

Biden calls for expanded efforts to protect LGBTQ rights globally

Biden set to accept more refugees after years of Trump restrictions

U.S. House punishes Republican congresswoman over incendiary remarks

Amazon bucks UK's grim labour market with 1,000 apprenticeships

Looming Apple privacy changes weigh on Snap despite revenue growth

SoftBank third-quarter earnings recovery seen driven by IPO boom

Activision Blizzard's annual sales forecast tops estimates on 'Call of Duty' boost

Pinterest beats estimates on ad spending recovery, strong user growth

T-Mobile beats quarterly postpaid phone additions estimates

Reddit rally' stocks bounce on day after selloff, then dip after hours

U.S. Treasury's Yellen to meet financial regulators Thursday to discuss volatility

Koss family rakes it in from Reddit-fueled rally

Instagram removes hundreds of accounts tied to username hacking

U.S. Senate Democrats push ahead on road to new COVID-19 relief

'America is back' - Biden touts muscular foreign policy in first diplomatic speech

Trump rejects call to testify at his impeachment trial

BNY Mellon, Google Cloud technology to predict Treasury settlement failures

Texas court considers hearing on changing venue of Google antitrust case

Exclusive-White House pulled into unionization effort at Amazon facility

Google phone cameras will read heart, breathing rates with AI help

GameStop rises, AMC dips in early U.S. premarket trading

Xilinx to supply chips to Fujitsu for U.S. 5G network gear

Hit by shortage, Volkswagen demands boost to Europe chip sector

Elon Musk, back on Twitter, turns his attention to Dogecoin

EU electric and plug-in hybrid car sales jump to over 1 million in 2020

Exclusive: China's Ant to hive off credit data in revamp; sees IPO in 2 years - sources

Biden says war in Yemen 'has to end,' U.S. will continue to support Saudi Arabia

Biden to pursue arms control, seeks to engage China: U.S. envoy

In uneasy truce, House Republicans fail to punish Greene or Cheney

Nokia warns of "challenging" year as it plays catch-up

Mazda expects chip shortage to affect about 7,000 vehicles in February

Uber's Mideast business Careem sees recovery slowing as infections rise

NatWest latest UK bank to switch to Mastercard debit cards from Visa

Ethereum scales record peak before futures launch

Five things to watch in Reddit stocks trading mania

Taiwan says auto chip shortage not a main topic for coming U.S. meeting

Parler CEO John Matze says he was fired by board

Robinhood to allow buying fractional shares of GameStop, AMC

Alibaba sets initial price guidance on \$5 billion bond offering: term sheet

U.S. House Republican leader does not plan to oust Cheney from leadership

post: CNN

Pentagon, stumped by extremism in ranks, orders stand-down in next 60 days

Biden decides to stick with Space Force as branch of U.S. military

Two Google engineers resign over firing of AI ethics researcher Timnit Geb  
ru

Amazon plans AI-powered cameras in delivery vans to improve driver safety

Qualcomm shares drop as chip supply constraints hold back sales

Ebay earnings beat on pandemic-driven surge in online shopping; shares soa  
r

PayPal profit tops estimates as pandemic drives online spending to record  
levels

SEC studies social media posts for signs of fraud in GameStop frenzy: Bloo  
mberg

Number of GameStop shares shorted edges higher: S3 Partners

Spotify outlook weakens as pandemic uncertainty persists

Ant Group reaches deal with China regulators on restructuring - source

Daimler to spin off truck unit, sharpen investor focus on Mercedes-Benz

Yellen calls for 'acting now - and acting big' on pandemic relief

Schumer, after Biden meeting, says Democrats united on a 'bold' COVID-19 b  
ill

Biden tells congressional Democrats would consider limits on who gets COVI  
D-19 checks

### Step 3: Create functions to perform news headlines recommendations



In [ ]:

```
def preprocess(text, char_ngram=5):  
    return set(text[head:head + char_ngram] for head in range(0, len(text) - char_ngram  
)  
)  
  
def get_forest(data, perms):  
    start_time = time.time()  
  
    minhash = []  
  
    for text in data['title']:  
        tokens = preprocess(text)  
        m = MinHash(num_perm=perms)  
        for s in tokens:  
            m.update(s.encode('utf8'))  
        minhash.append(m)  
  
    forest = MinHashLSHForest(num_perm=perms)  
  
    for i,m in enumerate(minhash):  
        forest.add(i,m)  
  
    forest.index()  
  
    print('It took %s seconds to build forest.' %(time.time()-start_time))  
  
    return forest  
  
def predict(text, data, perms, num_results, forest):  
    start_time = time.time()  
  
    tokens = preprocess(text)  
    m = MinHash(num_perm=perms)  
    for s in tokens:  
        m.update(s.encode('utf8'))  
  
    idx_array = np.array(forest.query(m, num_results))  
    if len(idx_array) == 0:  
        return None # if your query is empty, return none  
  
    result = data.iloc[idx_array]['title']  
  
    print('It took %s seconds to query forest.' %(time.time()-start_time))  
  
    return result
```

In [ ]:

```
data = pd.DataFrame(headlines,columns= ['title'])
data.head()
```

Out[ ]:

	title
0	U.S. thanks Taiwan for support for auto chips ...
1	Robinhood lifts trading restrictions on all st...
2	PayPal says to shut domestic payments business...
3	Shares of Tencent-backed Kuaishou triple in HK...
4	Australian drone firm reshapes strategy over G...

In [ ]:

```
forest = get_forest(data, 100)
```

It took 0.1442406177520752 seconds to build forest.

#### Step 4: Predict news headlines recommendations for the given title.

In [ ]:

```
title = "Stocks explained: What's going on with GameStop?"
result = predict(title, data, 100, 10, forest)
print('\n Top Recommendation(s) is(are) \n', result)
```

It took 0.007291078567504883 seconds to query forest.

```
Top Recommendation(s) is(are)
1    Robinhood lifts trading restrictions on all st...
11   Wisconsin governor clashes with lawmakers over...
Name: title, dtype: object
```

In [ ]:

```
title = "What is this?"
result = predict(title, data, 100, 10, forest)
print('\n Top Recommendation(s) is(are) \n', result)
```

```
Top Recommendation(s) is(are)
None
```