

# CS101 - Unit 6: How to Have Infinite Power

## Contents

- 1 [Infinite Power](#)
- 2 [Long Words](#)
  - 2.1 [Quiz: Long Words](#)
- 3 [Counter](#)
  - 3.1 [Quiz: Counter](#)
  - 3.2 [Quiz: Expanding our Grammar](#)
- 4 [Recursive Definitions](#)
- 5 [Ancestors](#)
  - 5.1 [Quiz: Ancestors](#)
- 6 [Recursive Procedures](#)
  - 6.1 [Quiz: Recursive Factorial](#)
- 7 [Palindromes](#)
  - 7.1 [Quiz: Palindromes](#)
- 8 [Recursive v. Iterative](#)
- 9 [Bunnies](#)
  - 9.1 [Quiz: Bunnies](#)
- 10 [Divide and Be Conquered](#)
  - 10.1 [Quiz: Counting Calls](#)
  - 10.2 [Quiz: Faster Fibonacci](#)
- 11 [Ranking Web Pages](#)
- 12 [Popularity](#)
  - 12.1 [Quiz: Good Definitions](#)
- 13 [Circular Definitions](#)
  - 13.1 [Quiz: Circular Definitions](#)
- 14 [Relaxation](#)
  - 14.1 [Quiz: Relaxation](#)
- 15 [Page Rank](#)
- 16 [Altavista](#)
  - 16.1 [Quiz: Altavista](#)
- 17 [Urank](#)
  - 17.1 [Quiz: Implementing Urank](#)
- 18 [Computing Page Rank](#)
- 19 [Formal Calculations](#)
- 20 [Computer Rank](#)
  - 20.1 [Quiz: Finishing Urank](#)
- 21 [Search Engine](#)

## Infinite Power

Welcome to unit 6! After this unit you will have learned all of the technical aspects that you will be tested on in the final exam. Unit 7 will consist of field trips and interviews, which will put what you have learned in context.

The big idea that will be introduced in unit 6 is recursive definitions, which you will learn how to use as a method for increasing your page ranking - being able to find the best page to respond to the query. The real goal of this unit is to give you infinite power!

Recall, that in unit 2 when you learned about procedures, you were told that the **if** statement gave you enough to write every possible computer program, which is infinitely powerful. Then, you learned how to use the **while** loop to go on. If you were infinitely powerful just knowing the **if** statement then you should not have needed to learn the **while** loop. You should have been able to build it from the things you have already seen - and it turns out that you can!

In this unit you will learn how to build up your own powerful control structures without using anything other than procedures. You will see that you can build up these control structures, as powerful as the **while** loop, from nothing but the procedures, if, and arithmetic and comparison operations that you learned in unit 2.

The point of learning this is not to be able to replace procedures, but to learn a new way of thinking called recursive definitions, which is a very powerful tool for solving problems.

## Long Words

### Quiz: Long Words

This is kind of a trick quiz. Don't worry if you're not a native English speaker. This quiz is just as hard for them as it is for you!

What's the longest word in the English language?

1. honorificabilitudinitatibus
2. antidisestablishmentarianism
3. hippopotomonstrosesquippedaliophobia
4. pneumonoultramicroscopicsilicovolcanoconiosis
5. None of the above

[Answer](#)

## Counter

A word is something that has meaning that is understood by the speakers of that

words language. A word could be defined as what is in a dictionary, but there are a lot of things that are words, but that are not in the dictionary.

There is a rule that says that for a word, you can make a new word by adding counter in front of the old one. The notation used is the BNF (Backus Naur Form) replacement grammar, which was introduced in unit 1. If you need a refresher, please see unit 1, sections 9-11. Recall that the basic property of a BNF grammar is to replace what is on the left by what is on the right.

- Word --> counter-Word

The meaning of the new word is something that goes against, or counter to the original word.

If you start with the word intelligence, (intelligence as in spycraft not smart), you can use the rule to replace intelligence with counter-intelligence, which means trying to thwart the intelligence of the enemy.

You can continue, replacing counter-intelligence with counter-counter-intelligence, which means trying to thwart the enemy's counter-intelligence. Repeating this again, you get a word that isn't used but that still has a sensible meaning - counter-counter-counter-intelligence.

- Word --> counter-Word
- intelligence
- counter-intelligence
- counter-counter-intelligence
- counter-counter-counter-intelligence

One of the long words in the first quiz, hippopotomonstrosesquippedaliophobia, means fear of long words. If you know what a word means, even if you've never seen counter in front of it, you can guess the meaning. Counter-hippopotomonstrosesquippedaliophobia means something that goes against the fear of long words. It could be a medication that can cure someone from the fear of long words.

If you add another counter in front of counter-hippopotomonstrosesquippedaliophobia, then you get counter-counter-hippopotomonstrosesquippedaliophobia, which is something that goes against counter-hippopotomonstrosesquippedaliophobia. Maybe coffee stops the medication from working and so coffee is a counter-counter-hippopotomonstrosesquippedaliophobia!

## Quiz: Counter

If the only rule we have for making words is this one:

- Word --> counter-Word

how many words can we make, starting from Word.

1. None
2. 1
3. 2
4. Infinitely Many

[Answer](#)

## Quiz: Expanding our Grammar

For this quiz, an extra rule is added.

How many different words can we make starting from Word using only these two rules:

- Word --> counter-Word
- Word --> hippopotomonstrosesquippedaliophobia
- None
- 1
- 2
- Infinitely Many

[Answer](#)

## Recursive Definitions

Recursive definitions work for things other than words, but you're probably most familiar with them from language. You will learn how to use them in procedures and in later courses you will see how to use them to define data structures. A lot of things in computing are defined in terms of recursive definitions.

A recursive definition has two parts: the base case and the recursive case.

In the previous example, with respect to Word, the second rule:

- Word --> hippopotomonstrosesquippedaliophobia

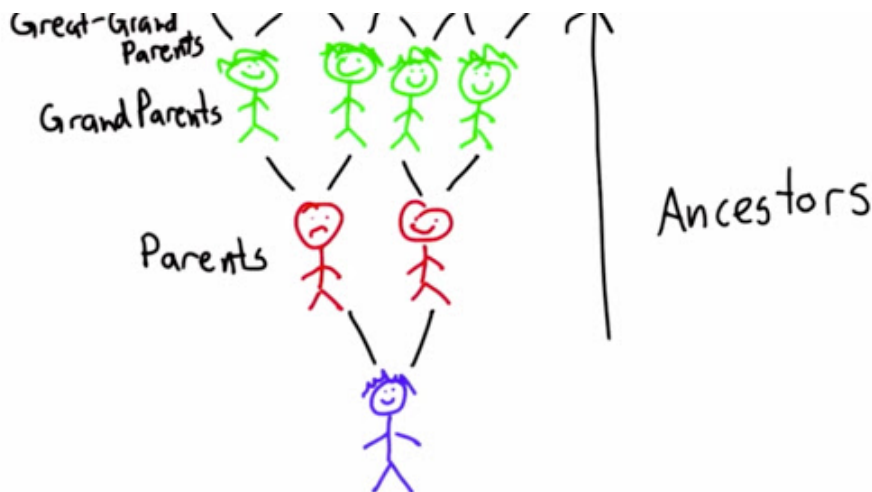
was the base case. It is a starting point and it is important that it is not defined in terms of itself. For programs it is usually going to be the smallest input, or the simplest input. The base case must be something that you already know how to define. You must already know the answer and not need to do anything to work it out.

The recursive case is defined in terms of itself, but not itself exactly. It is defined in terms of a smaller version of itself, as progress must be made towards the base case. You will see what this means in programs soon, but first another example -

not in terms of a program - to get a better idea of how things can be defined recursively.

## Ancestors

How can you define who your ancestors are?



Your parents are your ancestors, but they are not your only ancestors. Your parents have parents - your grandparents, who are also your ancestors. Your grandparents also have parents who are your ancestors too, and so on.

### Quiz: Ancestors

Which of these is the *best* definition of ancestors?

1. Ancestor --> Parent of Ancestor
2. Ancestor --> Parent
  - Ancestor --> Parent of Ancestor
3. Ancestor --> Parent
  - Ancestor --> Parent of Parent Ancestor --> Parent of Parent of Ancestor

[Answer](#)

## Recursive Procedures

You have seen how to use recursive definitions to make Words and define concepts like Ancestors. Now you'll see how to use recursive definitions to define a procedure. In unit 2, the factorial  $n$  was defined as the number of ways to arrange  $n$

items, which means that the input is  $n$ . This can be calculated as:

- $factorial(n) = n * (n-1) * (n-2) * \dots * 1$

This equation is not a very precise mathematical definition because of the dot, dot, dot. Humans understand it correctly, but it's not precise mathematically. Using a recursive definition allows a factorial to be defined precisely. For this, a base case is needed. This should be the simplest input. It is something for which the answer is already known. For factorial and for many procedures involving numbers, the simplest input is 0. The factorial of 0 is defined as 1, that is:

- $factorial(0) = 1$ ,

which is the base case.

In the imprecise definition:

- $factorial(n) = n * (n-1) * (n-2) * \dots * 1$ ,

the product:

- $(n-1) * (n-2) * \dots * 1$  is just the factorial of  $n-1$ ,

that is:  $factorial(n-1)$ .

This means you can write the recursive case as:

- for any integer  $n > 0$ ,  $factorial(n) = n * factorial(n-1)$ .

It makes sense when you think of the way to arrange  $n$  items. There are  $n$  ways to pick the first item, and then there are  $n-1$  items remaining. There are  $factorial(n-1)$  ways to arrange these  $n-1$  items.

Base case:  $factorial(0) = 1$  Recursive case:  $factorial(n) = n * factorial(n-1)$  for  $n > 0$

## Quiz: Recursive Factorial

Define a procedure, **factorial**, that takes a natural number as its input, and outputs the number of ways to arrange the input number of items.

You've already seen how to do this using a while loop. Your goal here is to define that procedure without using a while loop, that is, to define it using a recursive definition.

Note that a **natural number** is defined to be a positive whole number.

[Answer](#)

## Palindromes

Here's another example of defining a recursive procedure. A **palindrome** is a string that reads the same way forwards and backwards.

For example: String 'level' is a palindrome (if you read 'level' forwards, you get 'level' and if you read it backwards, you get exactly the same string). Some other typical examples of palindromes include:

- Any single letter is a palindrome. For e.g. 'a' (If you read 'a' forwards, you get 'a' and if you read 'a' backwards, you again get 'a')
- An empty string is also a palindrome (i.e. ""). If you read empty string " forwards, you get the empty string and if you read empty string backwards, you get an empty string.



Additionally:

- <http://norvig.com/palindrome.html>: This is an encyclopediac text on this topic by Prof. Peter Norvig.
- [http://en.wikipedia.org/wiki/Palindrome#Long\\_palindromes](http://en.wikipedia.org/wiki/Palindrome#Long_palindromes): You might also want to include some interesting palindromes/facts from this Wikipedia page.

## Quiz: Palindromes

Define a procedure, **is\_palindrome**, that takes as input a string, and outputs a Boolean to indicate if the input string is a palindrome.

First, try to think on your own and define the procedure that tests whether an input string is a palindrome or not. This is a pretty tough question. There are easy ways to reverse the string and check to see if it is same as the original string in Python. But since you have not studied them yet in this course, write this procedure using what you already know about Python. Here are some hints for defining this procedure:

- You might have already noticed that there is one simple case (i.e. empty string

") where you know that the string is a palindrome. So, this will be your base case. If the input to the procedure, **is\_palindrome** is an empty string, the result of **is\_palindrome** is *True*.

Note: When you write recursive procedures on numbers, the base case is often some small number (like 0 or 1). When you write recursive procedures on strings, the base case is more likely to be the simplest string (which is the empty string ").

- What will you do if the string is not an empty string? Well, one way you can solve this is by looking at the first and last letter of string. If these two are equal, then it might be a palindrome. It will be a palindrome only if all the letters left over in the middle are also a palindrome. In order to check whether remaining middle of string is a palindrome, you can recursively call **is\_palindrome**.

Base case: "" → True  
 Recursive case:  
 if first and last characters don't match → False  
 if they do match, is middle a palindrome?

## Recursive v. Iterative

Any procedure that you write recursively, you can also write without using a recursive definition. Here is another way to define **is\_palindrome**:

```
def iter_palindrome(s):
    for i in range(0, len(s) / 2):
        if s[i] != s[-(i + 1)]:
            return False
    return True
```

**iter\_palindrome** is written using a **for** loop. You loop using variable **i** in range 0 to length of input string divided by 2 (i.e. this loop is going through halfway of string **s**). Inside the loop, there is an **if** test that checks if the character at position **i** is different from the character in position **-(i + 1)** (i.e. counting from back of string **i**<sup>th</sup> positions away). If those characters are different, you have bumped into a mismatch and we return *False*. If they are not different, continue going through the loop. If you reach the end of the **for** loop without finding any differences, you know that it is a palindrome and you return *True*.

This alternative way to define **is\_palindrome** is more complicated to understand. If you want to test a very long palindrome, **iter\_palindrome** (i.e. iterative version) will be much more efficient than **is\_palindrome** (recursive version). There are a few reasons why:



1. Inside **is\_palindrome**, when we call **is\_palindrome** recursively:

```
return is_palindrome(s[1:-1])
```

This recursive version keeps making a new string every time you make a recursive call. This creates a new string and this step is pretty expensive.

1. Recursive calls themselves are fairly expensive. There are languages which make recursive calls really cheap, but Python is not one of them. For most procedures, the recursive way is often the most elegant and easiest way to return a correct result. But if you are worried about performance and want your procedure to work on really large inputs, you are better off finding a non-recursive way to define that procedure.

## Bunnies

**Fibonacci Numbers** are one of the most interesting things in mathematics. Once you know about them, you will start to see them all over the place, both in nature and design.

The name comes from Leonardo da Pisa, who is also known as Fibonacci. In 1202 he published a book called, *Liber Abaci*. The root, abaci, is the same for the word abacus, the calculating machine. *Liber Abaci* is loosely translated as the "book of calculation." The book introduced Indian mathematics to the West, particularly, Arabic numerals. Arabic numerals soon replaced the Roman numeral system, which had been widely used. In his book, Fibonacci showed how much easier it is to do calculations using numbers in the decimal system where the position of the number indicates its value. He showed this by introducing problems and using calculation to solve them.

The problem that became known as the Fibonacci Numbers, was one of the problems in his book. He posed the problem like this:

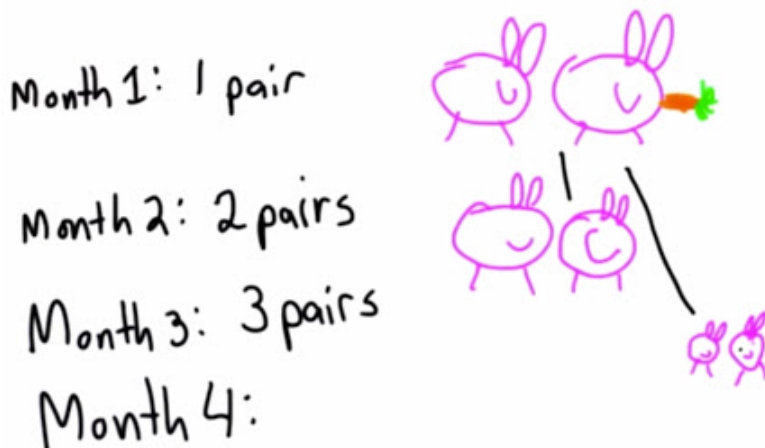
In the beginning there is one pair of rabbits. It takes one month for a rabbit to mature, and one month for a rabbit to produce offspring.

# Fibonacci Numbers



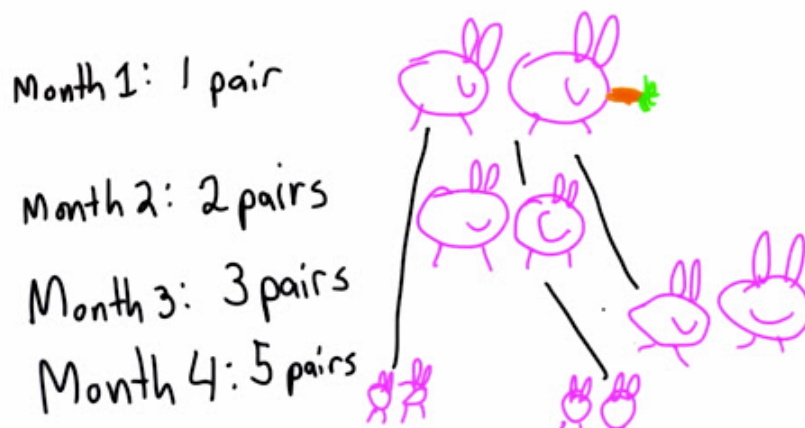
Every month a mature rabbit will produce a new pair of rabbits. Notice how in month three, since it takes a month for the rabbits to reach maturity (and only in maturity can the rabbits reproduce), there is only one set of offspring from the original pair of rabbits.

## Fibonacci Numbers

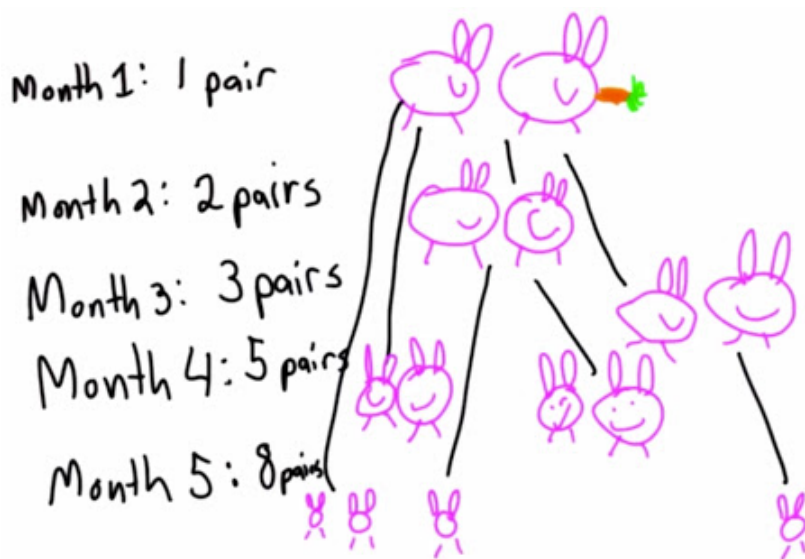


Assume that for the sake of this model, rabbits never die. In month four, one pair of baby rabbits will mature, while the two mature pairs of rabbits will reproduce. This makes a total of five pairs of rabbits.

## Fibonacci Numbers



This keeps going, since the model assumes rabbits never die, that every month a pair of mature rabbits produces a pair of rabbit babies and that it takes one month for rabbit babies to mature. Therefore, in month five the three mature pairs of rabbits will all produce a pair of offspring, making eight pairs of rabbits.



While this model is not realistic, it is an interesting mathematical model. This model can be written in a more formal way.

Each month the number of rabbits is the number of rabbit from the previous month, plus all of the rabbits that were mature, which is the number of rabbits from two months ago. After this observation you can predict the number of rabbit pairs there will be in month six using the same formula: the number of rabbit pairs in month five plus the number of rabbit pairs in month four. This makes 13 pairs of rabbits.

This is the model that Fibonacci developed. So, here is the question: can you figure out how many pairs of rabbits there will be at the end of month  $n$ , given any number  $n$ ?

Here is how you can define this mathematically:

- **fibonacci (0) = 0**
- **fibonacci (1) = 1**

This is different from the other recursive definitions you've seen in that there are two base cases. From here you can define every other fibonacci number recursively, starting from these base cases. So, the fibonacci number  $n$ , where  $n$  is some whole number greater than one, is equivalent to the sum of all the babies from the previous month, **fibonacci (n-1)**, plus all of the new babies. The number of new babies is the number of rabbit pairs two months ago (all of the mature rabbit pairs), **fibonacci (n-2)**. You can write this:

- $n > 1$
- $\text{fibonacci}(n) = \text{fibonacci}(n-1) + \text{fibonacci}(n-2)$

This equation defines every fibonacci number in terms of the two base cases and

the one recursive case.

$$\begin{array}{l}
 \text{fibonacci}(0) = 0 \\
 \text{fibonacci}(1) = 1
 \end{array}
 \left. \vphantom{\begin{array}{l} \text{fibonacci}(0) = 0 \\ \text{fibonacci}(1) = 1 \end{array}} \right\} \text{two base cases}$$

$$\text{recursive case } \left( \begin{array}{l} \text{fibonacci}(n) = \text{fibonacci}(n-1) \\ n > 1 \end{array} \right. + \text{fibonacci}(n-2)$$

## Quiz: Bunnies

Define a procedure, **fibonacci**, that takes a natural number (any whole number zero or higher) as its input, and outputs the value of that fibonacci number.

[Answer](#)

## Divide and Be Conquered

Here is the procedure you just defined as an answer to the previous quiz question:

```
def fibonacci(n):
    if n == 0:
        return 0
    if n == 1:
        return 1
    return fibonacci(n-1) + fibonacci(n-2)
print fibonacci(0) # print first base case
0
```

or

```
print fibonacci(1) # print second base case
1
```

and

```
print fibonacci(2) # plug 2 into final return statement
1
print fibonacci(3)
2
print fibonacci(4)
3
print fibonacci(5)
5
print fibonacci(10)
```

```

55
print fibonacci(24) # number of rabbits in two years
46368
print fibonacci(36) # number of rabbits in three years

```

If you try running the code above in your Python interpreter, you'll notice that the time to compute "fibonacci" this way gets longer as input numbers get larger. The reason for this is that you are making a lot of redundant computations. If you look at the procedure, you'll notice that "fibonacci(n)" recursively calls "fibonacci(n - 1)" and "fibonacci(n - 2)" every time the base case conditions are not satisfied.

At the start call "fibonacci(36)", which is broken down into recursive calls to "fibonacci(35)" and "fibonacci(34)". Now "fibonacci(35)" recursively calls "fibonacci(34)" and "fibonacci(33)", while "fibonacci(34)" recursively calls "fibonacci(33)" and "fibonacci(32)", and so on. See the tree below:

The procedure needs to do a lot of computations and it will take a long time to get calls to the base cases ("fibonacci(0)" and "fibonacci(1)"), which are the only places where the procedure stops making more recursive calls. If you look at the figure above, you'll notice that: . We need to evaluate fibonacci(32) 5 times. . We need to evaluate fibonacci(33) 3 times. . We need to evaluate fibonacci(34) 2 times. . We need to evaluate fibonacci(35) 1 time. . We need to evaluate fibonacci(36) 1 time.

Do you see a pattern in the list above? Try to see if you can solve the next question related to this.

## Quiz: Counting Calls

How many times will you need to evaluate "fibonacci(30)" in evaluating "fibonacci(36)"?

Figure this out without drawing the whole tree. Think about what you read in last section. Do you notice a pattern in the figure above that could help you to answer this question?

[Answer](#)

## Quiz: Faster Fibonacci

Define a faster "fibonacci" procedure that will enable you to compute "fibonacci(36)".

Your procedure should estimate the number of rabbits after 36 months, according to Fibonacci's model.

Hint: You'll need a "while" loop where you use variables to keep track of the previous two numbers, "fibonacci(n-1)" and "fibonacci(n-2)", so you can compute

the next one by adding those. You'll also need to figure out how to keep the variables up-to-date to maintain the previous two numbers each time you go through the loop.

Test your code on small numbers before trying `fibonacci(36)`. If you do it this way, you should be able to compute values of the Fibonacci Numbers for much higher input than you could with a recursive definition.

[Answer](#)

## Ranking Web Pages

Having survived the bunny uprising, you're ready to move on to the main goal of the class, which is to return the best page that matches the search query rather than returning all the pages. It's important to do this well. This is something that really distinguished Google from earlier search engines. They had a much smarter way of ranking pages. Often the first or second item in the returned search was what the user was looking for.

To recap from earlier units; first, you learned to build a crawler (units 1-3). The crawler followed all the links in the web pages and built an index. After units 4-5, you had an index, which was a hash table, where you could look up a keyword. You could find the entry where the keyword appears, and find the list of all the urls of all the pages that contain that keyword.

The order in which the pages appear in the list of urls associated with a keyword is the order the pages were crawled. This process says nothing about which pages are best. In the early days of the web, when there weren't many pages, this maybe wasn't too much of a problem since only a few pages might match a given keyword. Those days are long gone and now there could be thousands, if not millions of pages containing a given keyword. A good search engine ranks the pages so that the one at the front of the list is the one the user most likely wants.

The problem of deciding how to rank the pages leads to the question of how to decide popularity, which is the topic of the next section.

## Popularity

Consider a typical group of friends in middle school. One way to decide popularity is to look at friendship links. Friendship links go in one direction. Just because Bob is friends with Alice does not mean Alice is friends with Bob.

Is having a lot of friends enough to make you popular? No, it's not. You have to have the right sort of friends. It's no good to have lots of friends with no friends, you have to have friends who are popular.

Popularity is about having lots of friends who have lots of friends.

Initially, you can define popularity as the number of friends a person has. From the diagram above, it's the number of arrows pointing towards a person.

```
. 'popularity(p)' = '# of people who are friends with p'
```

This isn't quite right though because it doesn't take into consideration the number of friends of those friends. To do this, you could sum all the popularities of all the friends of a person. In mathematical notation:

```
. 'popularity(p) = sigma over f E friends of p    popularity(f)'
```

The notation  $\sum_f$ , is the summation sign, which tells you to sum up the "popularity(f)". The text under the symbol tells you what values of "f" to include in the sum. It tells you to sum the popularity of each friend, "f" of "p". If you're unfamiliar with the mathematical notation, here's the same thing in Python pseudocode.

Pseudocode is an outline of the code which is written for human readability rather than for a computer. [Read more on pseudocode.](#)

In the code below, you know that friends p means the friends of p, but it's not actually defined anywhere so the computer will not be able to run it. 

```
python def popularity(p): score = 0 for f in friends(p): score = score + popularity(f) return score
```

## Quiz: Good Definitions

Is this a good recursive definition? For something to be a good definition it has to provide a meaningful answer for all possible inputs.

1. Yes
2. No

[Answer](#)

## Circular Definitions

How can this problem be fixed? With all the other recursive definitions you had a base case - a way to stop.

For the recursive factorial definition, you predefined the value at 0 to be 1, that is, **factorial(0) --> 1.**

For the palindromes, you defined an empty string to be a palindrome, that is, **palindrome('') --> True.**

For the Fibonacci sequence, you had two base cases.

For all of these definitions you had a starting point that was not defined in terms of the thing you're defining. That is why they were good recursive definitions. You had a base case. Maybe inventing a base case will solve the popularity problem.

Assume Alice has popularity 1, and try that as a base case. For the mathematical definition, this is:

- $\text{popularity}(\text{'Alice'}) = 1$
- $\text{popularity}(p) = \sum \text{popularity}(f)$  over friends  $f$  of  $p$

For the code, you need to add the base case, which is an if statement, to see if the person you're checking the popularity of is Alice.

```
def popularity(p):
    if p == 'Alice':
        return 1
    score = 0
    for f in friends(p):
        score = score + popularity(f)
    return score
```

## Quiz: Circular Definitions

Would this definition work?

1. Only if everyone is friends with 'Alice'.
2. Only if no one is friends with 'Alice'.
3. Only if there is a friendship path from everyone to 'Alice'. (There is some way to follow links from every person in the graph to get to 'Alice')
4. Only if there are no cycles in the graph. (There is no way to start from one person and end up back at the same person by following friendship links.)
5. No.

[Answer](#)

## Relaxation

There was no sensible base case that provides a good recursive definition. Instead, an algorithm called the **relaxation algorithm** can be used. The basic idea is simple. Start with a guess and then loop where you do something to improve the guess. There isn't a good stopping place yet, or a clear starting place, like setting the popularity of Alice. Each time you go through the loop, the guess will be refined, and at some point you'll stop and take that to be the result you want. In summary:

```
# start with a guess
while not done:
```



make the guess better

The procedure will have an extra parameter, which is a time step:

`popularity(<time step> ,<person>) --> score`

The base case will be to set the popularity for everyone at time 0 to 1. For the recursive step, the popularity of each of their friends at the previous time step,  $t-1$  is summed. In mathematical terms this is:

Base case:  $popularity(0, p) = 1$  Recursive step:  $popularity(t, p) = \sum_{f \in \text{friends}(p)} popularity(t-1, f)$  for  $t > 0$ .

In Python code:

```
def popularity(t,p):
    if t == 0: # base case, at time step 0
        return 1 # the score is always 1
    else:
        score = 0
        for f in friends(p): # summing over the friends
            score = score + popularity(t-1,f) # adding the popularity at the
        return score
```

So, now you have a new definition written in both mathematical notation and in Python code.

## Quiz: Relaxation

Is this a good recursive definition? (By this, it is not meant whether this is a good definition of popularity. For the way in which popularity is defined, for all possible inputs, that is, all possible values for  $t$  and  $p$ , does it give a result?)

1. yes
2. only if people can't befriend themselves
3. only if everyone has at least one friend
4. only if everyone is more popular than 'Alice'

[Answer](#)

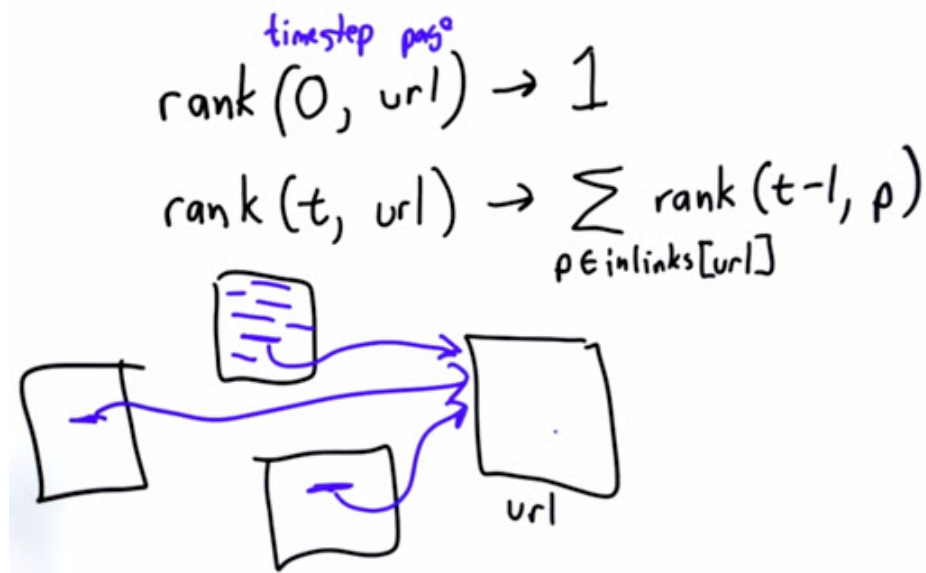
## Page Rank

Ranking web pages is the same as measuring popularity for people. Links on the web are analogous to friendships in that model. And links from some pages count for more than links from others.

This model is a random web surfer who starts at a random page and then follows the links at random. The popularity of a page is the probability that the random

surfer reaches a particular page.

The rank function is defined recursively over time. At time 0, the rank of a url is 1. At time  $t$ , the rank of a url is the sum of the ranks for all pages that link to that url.



In addition, the rank contributed by each page is inversely weighted by the number of outlinks from that page. So, divide each rank in the sum by the number of outlinks from that page.

$$\sum_{p \in \text{inlinks}[\text{url}]} \text{rank}(t-1, p) / \text{outlinks}[p]$$

## Altavista

On this model, pages with no links have a rank of 0, which makes it very hard to start a new page. So, instead each page will have some starting rank greater than 0.

Since the model represents the probability that a random surfer reached a given page, the ranks should be a probability distribution. This means that the ranks for all pages will sum up to 1. And so at time 0, instead of 1, the rank is  $1/N$  for each page, where  $N$  is the number of pages.

A **damping constant** is used to diminish the raw values of the ranking algorithm. In the model, this represents the probability that a page was reached via following

a link. Set the damping constant,  $d$ , to 0.8 for now.

At time  $t$ , add  $(1-d)/N$  to each rank that will represent the probability that the page was not reached via following a link. Multiply the first term, the sum, by  $d$ . This will make the ranks a probability distribution at any given time step.

## Quiz: Altavista

What is [AltaVista](#)?

1. The view from the Udacity headquarters.
2. The most popular web search engine in 1998.
3. Spanish for "You're Terminated, Baby!"
4. A small town in Virginia.

[Answer](#)

## Urank

Since [PageRank](#) is a registered trademark of Google, the algorithm will be called URank instead.

URank needs to keep track of which pages link to which pages, so you'll need a data structure to keep track of which pages link to which other pages. You'll use a directed graph. A **directed graph** is a data structure where nodes are linked to other nodes, and the links only go one way. (see [http://en.wikipedia.org/wiki/Directed\\_graph](http://en.wikipedia.org/wiki/Directed_graph))

So **crawl\_web** will produce a graph in addition to an index, where the graph gives a mapping from each page to all the pages it links to. The graph will be a dictionary, since it is a mapping from individual URLs to lists of URLs.

Also, add the variable **outlinks** which stores the return value of **get\_all\_links(content)** so it can be used for both **tocrawl** and **graph**. Adding the line to update the graph will be a quiz.

## Quiz: Implementing Urank

For this quiz, add one line of code which will update the graph for each page crawled.

```
def crawl_web(seed): # returns index, graph of inlinks
    tocrawl = [seed]
    crawled = []
    graph = {} # <url>, [list of pages it links to]
    index = {}
    while tocrawl:
```

```

page = tocrawl.pop()
if page not in crawled:
    content = get_page(page)
    add_page_to_index(index, page, content)
    outlinks = get_all_links(content)

    #Insert Code Here
    union(tocrawl, outlinks)
    crawled.append(page)
return index, graph

```

[Answer](#)

## Computing Page Rank

Much like the Fibonacci solution before, you can use an iterative solution to implement the recursive definition, instead of a recursive solution. Use a loop to update the page ranks for each time step.

The output of **compute\_ranks** is a dictionary mapping each URL to its rank, which is a number.

For the homework, the function **lookup\_best** will find the highest-ranking page for a given keyword, using both the index and the page ranks.

## Formal Calculations

Remember how the ranking function was defined; **npages** refers to the number of pages:

$$\begin{aligned}
 \text{rank}(0, \text{url}) &= 1/\text{npages} \\
 \text{rank}(t, \text{url}) &= (1-d)/\text{npages} + \\
 &\quad \text{sum}([\text{for each page 'p' that links to URL,} \\
 &\quad d * \text{rank}(t-1, p) / (\text{number of outlinks from p})])
 \end{aligned}$$

Since the ranks should not depend on the order that the pages were examined by the algorithm, you need to keep track of the ranks at the last time step. Keep two separate dictionaries, **ranks** and **newranks**, where **newranks** is the working space for each time step. This is similar to the trick used earlier for the iterative Fibonacci solution.

## Computer Rank

The code for this section is shown below under the quiz. **d** is the damping factor. **numloops** is the number of times to do our "relaxation". Changing the number of loops can give different results. **npages** is the number of pages in the graph, which

is given by **len(graph)**.

Initially, all ranks are set to **1.0/npages** (remember: the decimal point to use floating point arithmetic), matching the base case of the recursive definition.

Then, the algorithm loops through **numloops** times, updating the rank for each page in the graph. **newrank** is initialized to **(1-d)/npages**, and then the quiz will be to update **newrank** with the sum of the inlink ranks. Then **newrank** is stored in the **newranks** dictionary.

When the **for** loop is finished, assign **newranks** to **ranks**, since the calculations are finished for that time step. At the end of the function, return **ranks**.

## Quiz: Finishing Urank

Update **newrank** based on the values of the previous iteration, **ranks**, and the incoming links from **graph**.

```
def compute_ranks(graph):
    d = 0.8 # damping factor
    numloops = 10

    ranks = {}
    npages = len(graph)
    for page in graph:
        ranks[page] = 1.0 / npages

    for i in range(0, numloops):
        newranks = {}
        for page in graph:
            newrank = (1 - d) / npages

            #Insert Code Here

            newranks[page] = newrank
        ranks = newranks
    return ranks
```

[Answer](#)

## Search Engine

Congratulations! You have now built a search engine! You learned how to collect a corpus using a web crawler, how to build an index and how to make it faster. Most recently you learned how to rank the results. Your search engine does page ranking better than any search engines before 1998.

Now, there is one thing left to learn - how to use the ranks.

For your homework you are asked to use the ranks to get the best result. To find the best result, just using the dictionary of ranks is not enough. Instead, you need to use the dictionary to find the result that matches the query with the best rank.

There are still a few problems left to solve before you can build a search engine to compete with Google. Finding a name for your search engine is probably the hardest problem. Yoogoo? [DuckDuckFind](#)?

Another problem to solve is actually getting your search engine on the web so that other people can send queries to it. You can learn how to do this in the upcoming Web Applications Course.

In unit 7, you will get prepared for the final exam and you will see some interesting examples of using computing in context.