```ruby
 1  require 'nokogiri'
 2  require 'open-uri'
 3
 4  =begin
 5  Ejercicio 8. Construir un Web Crawler (rastreador, araña, o robot) en Ruby capaz de analizar
 6  una página web concreta y que cree una base de datos con la información obtenida
 7  =end
 8
 9  =begin
10      Usage: Gets source code and urls from a page
11      Name of method: getSoruceCode
12      Date of creation: 29/04/2021
13      Members: Roberto Jiménez y Alberto Pérez
14      Last modification: 06/05/2021
15      Parameters:
16          Entry:
17              - URL: Parameter of the program
18          Out:
19              - list with urls without parsing
20  =end
21
22  def getSourceCode url
23      # Gets source code from page
24      parsed_data = Nokogiri::HTML.parse(URI.open(url))
25      list = []
26      puts parsed_data.class
27
28      #Gets all urls  from source code
29      tags = parsed_data.xpath("//a")
30      tags.each do |tag|
31          list.append("#{tag[:href]}\t#{tag.text}")
32      end
33      return list
34  end
35
36  =begin
37      Usage: Insert data received into a DB
38      Name of method: insertDB
39      Date of creation: 29/04/2021
40      Members: Roberto Jiménez y Alberto Pérez
41      Last modification: 06/05/2021
42      Parameters:
43          Entry:
44              - list: list with urls
45              - url: URL from page to analyze
46          Out:
47              - None, output is saved in DB
48  =end
49
50  def insertDB list, url
51      #Insert data into DB
52      `echo -n #{url} > aux.csv`
53      `echo " ~" >> aux.csv`
54      `echo #{list} >> aux.csv`
55      `sudo mysql --local-infile=1 -h "localhost" -u rober "-prober"  < insertarBD.sql`
56  end
57
58  #Parse information with Grep command
59  list = getSourceCode ARGV[0]
```

```ruby
60  parsedList = `echo #{list} | grep -E 'http[s]*://[a-zA-Z0-9./?@=%&:_#-]*'`
61  finalList =  parsedList.split(',')
62  insertDB finalList, ARGV[0]
63  for url in finalList
64      puts url
65  end
66
67
```