

Web Scrapping

with



ABOUT ME

Hengki Sihombing

Building Karejo.com - Organizer **JakartaJS**

twitter @hengkiardo

github @aredo

hengki@karejo.com

Schedule

- **What is Web Scrapping**
- **Why we do that**
- **How we do that in Node.js**
- **Our Target**
- **Demo Code**

What is Web Scrapping

Web scraping is a computer software technique of extracting information from websites.

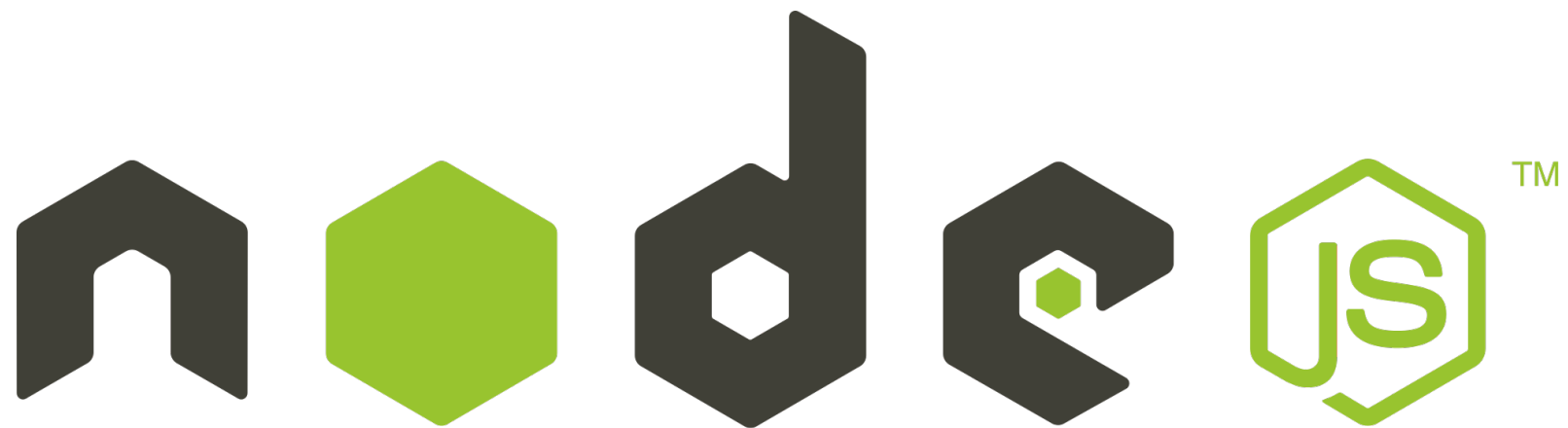
Why we do that

- **To get more accessible data**
- **Improve our hacking skill**
- **Have fun and build some nice project**

Some product doing Web scrapping

- **Wego, Skyscanner**
- **Flipboard, Instapaper, Nuzzel**
- **Telunjuk.com**
- **Karejo.com**
- **etc..**

How we do with



Our Dependencies

Request

Cheerio

Async

Mongoose

Express

Agenda


```
1 {
2   "name": "jakartajs-des-2015",
3   "version": "1.0.0",
4   "main": "app.js",
5   "scripts": {
6     "start": "nodemon app.js"
7   },
8   "repository": {
9     "type": "git",
10    "url": "git@github.com:aredo/jakartajs-des-2015.git"
11  },
12  "author": "Hengki Sihombing <hengki.sihombing@gmail.com>",
13  "license": "(ISC OR GPL-3.0)",
14  "dependencies": {
15    "agenda": "^0.6.28",
16    "agenda-ui": "0.0.7",
17    "async": "^1.4.2",
18    "cheerio": "^0.19.0",
19    "express": "^4.13.3",
20    "moment": "2.10.6",
21    "mongoose": "^4.1.11",
22    "mongoose-timestamp": "^0.4.0",
23    "request": "^2.65.0",
24    "lodash": "^3.10.1"
25  }
26 }
27
```

Step by step Scapping

```
[  
    "Requesting Pages",  
    "Parsing HTML",  
    "Traversing Dom",  
    "Collecting Data",  
    "Save to MongoDB"  
]
```

Request

```
1 var Request = require('request')
2
3 var options = {
4   url: url,
5   headers: {
6     'User-Agent': 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_10_1
7       AppleWebKit/537.36 (KHTML, like Gecko) Chrome/44.0.2403.125
8       Safari/537.36',
9   }
10 };
11
12 Request(options, function (err, response, body) {
13   if (err) {
14     return throw new Error()
15   }
16   callback(err, body)
17 })
```

Cheerio

```
1 var cheerio = require('cheerio')
2
3 var html = `<!DOCTYPE html>
4 <html lang="en" class="no-js">
5   <head>
6   </head>
7   <body>
8     <h2 class="title">Hello world</h2>
9   </body>
10 </html>`
11
12 $ = cheerio.load(html)
13
14 var title = $('h2').text()
15
16 console.log(title)
```

Our Target

SEMUA KATEGORI

Cari produk, kategori, atau merk

Cari



VOUCHER RP 50.000
DAPATKAN DISINI

SISTEM OPERASI

- ☐ Android (6130)
- ☐ Android 2.2 Froyo (5)
- ☐ Android 2.3 (4)
- ☐ Android 2.3 Gingerbread (3)
- ☐ Android 4.0 Ic (2)
- ☐ Android 4.0 Ice Cream (3)
- ☐ Android 4.1 Jelly Bean (25)
- ☐ Android 4.2 Jelly Bean (13)

more...

HARGA

6000 RP - 4646464

JUMLAH PORT USB

- ☐ 1 (11)
- ☐ 2 (10)
- ☐ 3 (2)

TIPE KABEL

- ☐ Aux (117)
- ☐ DVI (24)
- ☐ HDMI (273)
- ☐ Lainnya (113)
- ☐ Lighting (301)
- ☐ Micro USB (3035)
- ☐ USB (2932)

Handphone & Tablet

840427 Results

ATUR BERDASARKAN:

Terpopuler

Grid List



Mi 4i - 16GB - Putih

RP 2.999.000

- 13%

RP 2.599.000

★★★★★ (122 reviews)

Seller lainnya dari RP 2.622.000



Lenovo A7000 Special Edition - 16 GB - ...

RP 2.399.000

- 4%

RP 2.299.000

DECLEN100 RP 2.199.000

★★★★★ (196 reviews)



Alcatel Flash 2 - 16 GB - Volcanic Gre ...

RP 2.499.000

- 20%

RP 1.999.000

★★★★★ (167 reviews)



Infinix Note 2 X600 4G LTE - 16 GB - A ...

RP 1.999.000

- 5%

RP 1.899.000

★★★★★ (21 reviews)



Samsung Galaxy J5 Dual SIM - 8 GB - Putih

RP 2.899.000

- 9%

RP 2.649.990

★★★★★ (7 reviews)



Lenovo A6000 - 4G LTE - 8 GB - Hitam

RP 1.699.000

- 19%

RP 1.375.000

★★★★★ (196 reviews)



Infinix Android One X510 Hot 2 - 16GB ...

RP 1.499.000

- 13%

RP 1.299.000

★★★★★ (72 reviews)



Delcell Power Bank Note Polymer Batter ...

RP 229.000

- 31%

RP 159.000

★★★★★ (22 reviews)

DemoCode

Agenda

```
1  var config = {
2    mongodb: 'mongodb://127.0.0.1/jakartajs-des-2015',
3  }
4
5  var agendaOptions = {
6    address: config.mongodb,
7    collection: "cron"
8  }
9
10 var agenda = new Agenda({ db: agendaOptions })
11
12 agenda.define('lazada.co.id - handphone tablet', function (job, done) {
13
14   require('./handphone-tablet.js')(done, job)
15
16 })
17
18 agenda.every('10 minutes', 'lazada.co.id - handphone tablet');
19
20 agenda.start()
```


Please!!



Let's Join Us



Software Engineer
Front-End Developer
UI/UX Designer

let me know by email :

hengki@karejo.com

Thank you

<https://jakartajs-join.herokuapp.com>
meetup.com/JakartaJS