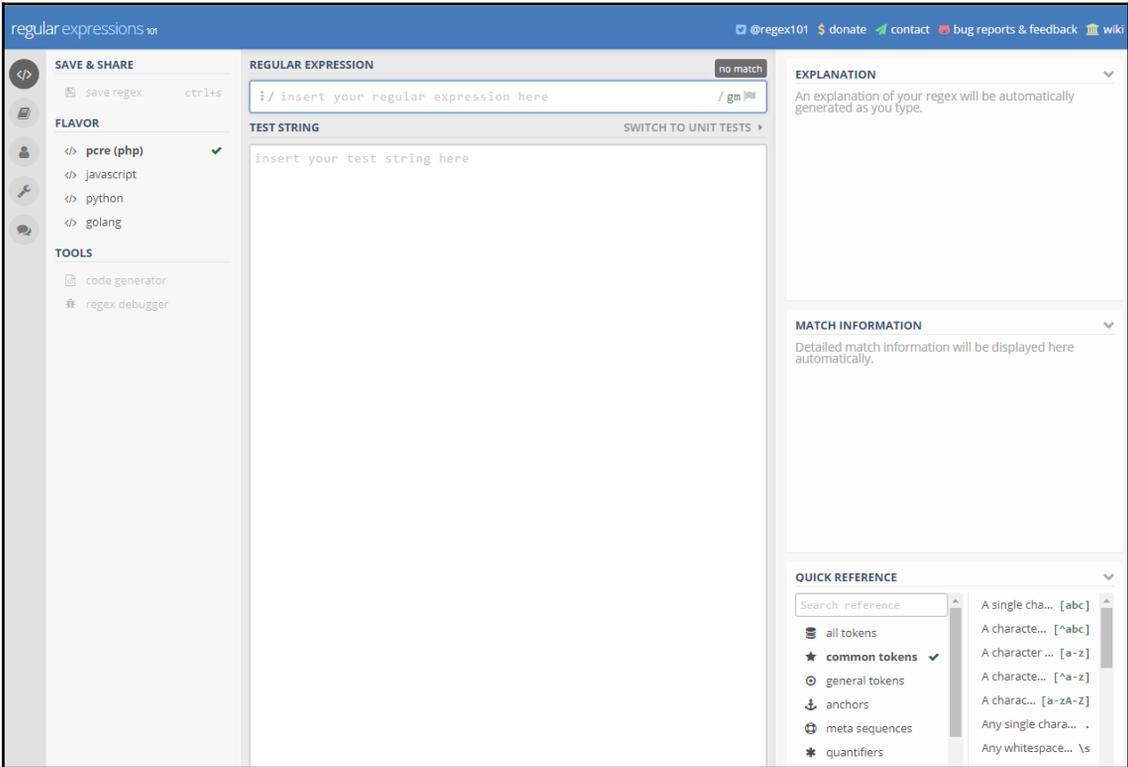
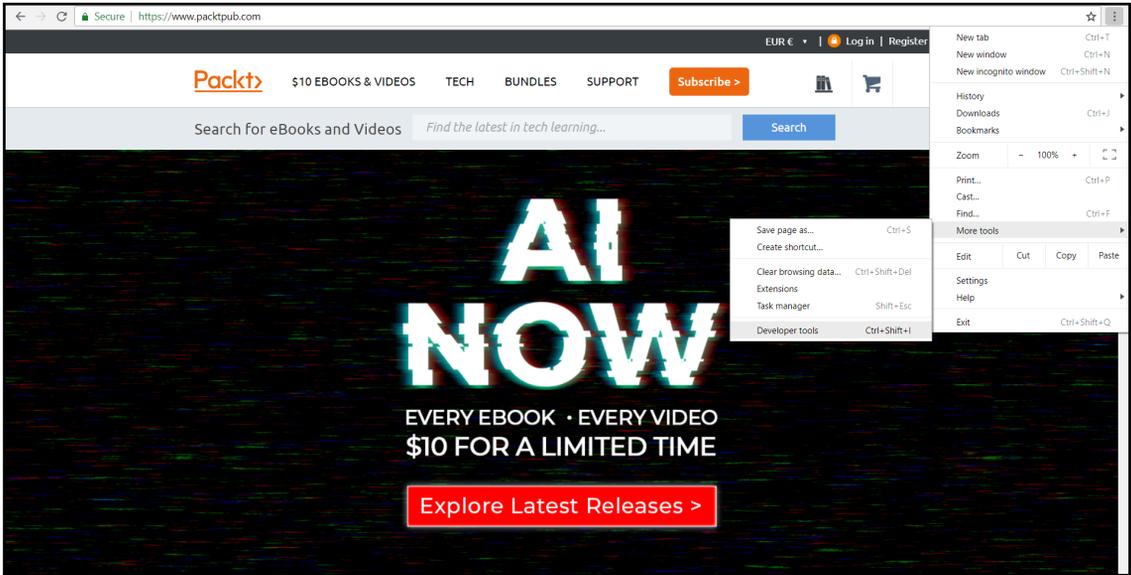


Chapter 1: Introduction to Web Scraping

```
≥ $x('//*[@id="menu-packt"]/span')
< ▼ [span.menu-text] ⓘ
  ▼ 0: span.menu-text
    accessKey: ""
    assignedSlot: null
    ▶ attributeStyleMap: StylePropertyMap {size: 0}
    ▶ attributes: NamedNodeMap {0: class, class: class, length: 1}
    autocapitalize: ""
    baseURI: "https://www.packtpub.com/"
    childElementCount: 0
    ▶ childNodes: NodeList [text]
    ▶ children: HTMLCollection []
    ▶ classList: DOMTokenList ["menu-text", value: "menu-text"]
    className: "menu-text"
    clientHeight: 0
```

Chapter 2: XML Path Language and Regular Expression Language





```
Elements Console Sources Network Performance Memory Application Security Audits
▼<div id="main-container" class="front not-logged-in page-index no-sidebars">
  ▼<div class="section-inner cf">
    ::before
    ▼<div class="menu-bar-popup-search-inner">
      ▼<form accept-charset="UTF-8" method="post" id="packt-libraries-main-search-form">
        ▼<div>
          ▼<div id="menu-bar-search" class="cf">
            ::before
            ▼<div id="menu-bar-search-title">
              ...
              <div class="form-item">Search for eBooks and Videos</div> == #0
            </div>
            ▶<div id="menu-bar-search-box">...</div>
            ▶<div id="menu-bar-search-button">...</div>
            <span class="magnifier" onclick="Packt.util.submitClosestForm(this)"></span>
            <input type="hidden" name="sort" id="edit-sort-1" value="0">
            <input type="hidden" name="types" id="edit-types-1" value="0">
            <input type="hidden" name="forthcoming" id="edit-forthcoming-1" value="1">
            <input type="hidden" name="available" id="edit-available-1" value="1">
            <input type="hidden" name="count" id="edit-count-1" value="20">
          html #ppv4 div #page #main-container div div.menu-bar-popup-search-inner form#packt-libraries-main-search-form div div#menu-bar-search.cf div#menu-bar-search-title div.form-item
```

regular expressions @regex101 [\\$ donate](#) [contact](#) [bug reports & feedback](#) [wiki](#)

SAVE & SHARE
 save regex ctrl+s

FLAVOR
 pcre (php) ✓
 javascript
 python
 golang

TOOLS
 code generator
 regex debugger

REGULAR EXPRESSION 3 matches, 576 steps (~394ms)
`/([a-z]+)([a-z]+\.[a-z]+)|([a-z]+\.[a-z]+)/gm`

TEST STRING SWITCH TO UNIT TESTS
 Olgun Aydin, info@olgunaydin.com, Gdansk, Poland
 Olgun Aydin, olgunaydinn@gmail.com, Gdansk, Poland
 Olgun Aydin, olgun.aydin@olgunaydin.com, Gdansk, Poland

EXPLANATION
 /([a-z]+)([a-z]+\.[a-z]+)|([a-z]+\.[a-z]+)/gm
 1st Capturing Group ([a-z]+)([a-z]+\.[a-z]+)
 1st Alternative ([a-z]+)
 2nd Capturing Group ([a-z]+\.[a-z]+)
 Match a single character present in the list below [a-z]
 Quantifier — Matches between one and unlimited times, as many times as possible, giving back as needed (greedy)
 a-z a single character in the range between a and z (index)

MATCH INFORMATION
 Match 1
 Full match 13-32 "info@olgunaydin.com"
 Group 1 13-17 "info"
 Group 2 13-17 "info"
 Match 2
 Full match 62-83 "olgunaydinn@gmail.com"
 Group 1 62-73 "olgunaydinn"

QUICK REFERENCE
 Search reference
 all tokens A single character of: a, b or c [abc]
 common tokens ✓ A character except: a, b or c [^abc]
 general tokens A character in the range: a-z [a-z]
 anchors A character not in the range... [^a-z]
 meta sequences A character in the range: ... [a-zA-Z]
 Any single character

SUBSTITUTION

regular expressions @regex101 [\\$ donate](#) [contact](#) [bug reports & feedback](#) [wiki](#)

SAVE & SHARE
 save regex ctrl+s

FLAVOR
 pcre (php) ✓
 javascript
 python
 golang

TOOLS
 code generator
 regex debugger

REGULAR EXPRESSION 3 matches, 491 steps (~2ms)
`/([a-z]+|[0-9+])([a-z]+\.[a-z]+|[0-9+])|([a-z]+\.[a-z]+|[0-9+])/gm`

TEST STRING SWITCH TO UNIT TESTS
 Olgun Aydin, info@olgunaydin.com, Gdansk, Poland
 Olgun Aydin, olgunaydinn88@gmail.com, Gdansk, Poland
 Olgun Aydin, olgun.aydin35@olgunaydin.com, Gdansk, Poland

EXPLANATION
 /([a-z]+|[0-9+])([a-z]+\.[a-z]+|[0-9+])|([a-z]+\.[a-z]+|[0-9+])/gm
 1st Capturing Group ([a-z]+|[0-9+])([a-z]+\.[a-z]+|[0-9+])
 1st Alternative ([a-z]+|[0-9+])
 2nd Capturing Group ([a-z]+\.[a-z]+|[0-9+])
 Match a single character present in the list below [a-z]
 Quantifier — Matches between one and unlimited times, as many times as possible, giving

MATCH INFORMATION
 Match 1
 Full match 13-33 "info@olgunaydin.com"
 Group 1 13-18 "info1"
 Group 2 13-18 "info1"
 Match 2
 Full match 63-86 "olgunaydinn88@gmail.com"
 Group 1 63-76 "olgunaydinn88"
 Group 2 63-76 "olgunaydinn88"

QUICK REFERENCE
 Search reference
 all tokens A single character of: a, b or c [abc]
 common tokens ✓ A character except: a, b or c [^abc]
 general tokens A character in the range: a-z [a-z]
 anchors A character not in the range... [^a-z]
 meta sequences A character in the range: ... [a-zA-Z]
 Any single character

SUBSTITUTION


```

> $x('//*[@id="menu-books"]')
< [div#menu-books.menu-item-text] ⓘ
  ▼ 0: div#menu-books.menu-item-text
    accessKey: ""
    align: ""
    assignedSlot: null
    ▶ attributeStyleMap: StylePropertyMap {size: 0}
    ▶ attributes: NamedNodeMap {0: class, 1: id, class: class, id: id, length: 2}
    autocapitalize: ""
    baseURI: "https://www.packtpub.com/"
    childElementCount: 1
    ▶ childNodes: NodeList(3) [text, span.menu-text, text]
    ▶ children: HTMLCollection [span.menu-text]
    ▶ classList: DOMTokenList ["menu-item-text", value: "menu-item-text"]
    className: "menu-item-text"
    clientHeight: 65
    clientLeft: 0
    clientTop: 0
    clientWidth: 156
    contentEditable: "inherit"
    ▶ dataset: DOMStringMap {}
    dir: ""
    draggable: false
    ▶ firstChild: text
    ▶ firstElementChild: span.menu-text
    hidden: false

```

The screenshot shows the browser's developer tools with the 'Elements' panel open. The DOM tree is expanded to show the following structure:

```

<div id="page-header">
  <div id="account-bar" class="cf"></div>
  <div id="menu-bar">
    <div class="section-inner cf">
      ::before
      <div id="menuIcon"></div>
      <div id="menu-links" class="">
        <a href="/"></a>
        <a href="/all"></a>
        <div class="menu-item-tab">
          <div class="menu-item-text" id="menu-books">
            <span class="menu-text">Books & Videos</span>
          </div>
        </div>
      </div>
      <a href="/tech"></a>
      <a href="/?></a>
      <a href="/books/content/support"></a>
      <a href="https://hub.packtpub.com" target="blank"></a>
      <a id="menu-packtlib-href" href="mailto:" target="blank"></a>
    </div>
  <div id="menu-icons"></div>
  ::after

```

The breadcrumb at the bottom of the developer tools reads: `html.js > body > ppv4.with-logo > div#respoPage > div#page > div#page-header > div#menu-bar > div#section-inner.cf > div#menu-links > a > div#menu-item-tab > div#menu-books.menu-item-text > span.menu-text`

```
> $x('//div[@class="menu-item-text"]/span/text()')
< ▼ (18) [text, text, text] ⓘ
  ▶ 0: text
  ▶ 1: text
  ▶ 2: text
  ▶ 3: text
  ▶ 4: text
  ▶ 5: text
  ▶ 6: text
  ▶ 7: text
  ▶ 8: text
  ▶ 9: text
  ▶ 10: text
  ▶ 11: text
  ▶ 12: text
  ▶ 13: text
  ▶ 14: text
  ▶ 15: text
  ▶ 16: text
  ▶ 17: text
  length: 18
  ▶ __proto__: Array(0)
```

```
> $('//div[@id="menu-books"]/span/text()')
< ▼ [text] ⓘ
  ▼ 0: text
    assignedSlot: null
    baseURI: "https://www.packtpub.com/"
    ▶ childNodes: NodeList []
    data: "Books & Videos"
    firstChild: null
    isConnected: true
    lastChild: null
    length: 14
    nextElementSibling: null
    nextSibling: null
    nodeName: "#text"
    nodeType: 3
    nodeValue: "Books & Videos"
    ▶ ownerDocument: document
    ▶ parentElement: span.menu-text
    ▶ parentNode: span.menu-text
    previousElementSibling: null
    previousSibling: null
    textContent: "Books & Videos"
    wholeText: "Books & Videos"
    ▶ __proto__: Text
    length: 1
    ▶ __proto__: Array(0)
```

Chapter 3: Web Scraping with rvest

DEVVERİ.COM Boğulacaksan büyük veride boğul!

Big Data Hadoop NoSQL Yazarlar

Amazon EMR ile Spark

18 Ocak 2018 Hakan İtler Cloud 0



Bu yazıda Amazon EMR üzerinde bir Spark uygulamasının nasıl çalıştırabileceğinden bahsedeceğim. Eğer EMR ile ilgili bir önceki yazıyı okumadıysanız bu yazıyı, AWS Big Data teknolojileri ile ilgili genel bilgi için de bu yazıyı okuyabilirsiniz. EMR üzerinde çalıştıracağımız örnek uygulamada daha önce defalarca kullandığım NYSE verisini kullanacağım. Tab karakterleri ile ayrılmış bu dosya içerisinde günlük borsa [...]

AWS, Big Data, Cloud, EMR, Hadoop, S3, Spark

Amazon EMR

13 Ocak 2018 Hakan İtler Cloud 0

Amazon EMR Nedir? Amazon Elastic MapReduce (EMR), büyük veri işlemeyi kolaylaştırmak amacıyla Amazon tarafından yönetilen, içerisinde Hadoop, Spark gibi açık kaynaklı büyük veri teknolojilerini içeren bir servistir. Aslında temelde AWS üzerinde Hadoop kümesi kurmak için tek tek sunucuları açmak, gerekli yazılımları yüklemek gibi işlemleri otomatik olarak yapmaktadır. Bu sayede tek tık ile bir kümeyi kurabileceğiniz [...]

Amazon, AWS, Big Data, EMR, Hadoop, Spark

AWS ile Big Data

11 Ocak 2018 Hakan İtler Cloud 0



Bulut teknolojilerinin öncüsü olan Amazon Web Servisleri bize birçok büyük veri teknolojisini esnek ve uygun maliyetli olarak test etme ve kullanma şansı sağlıyor. Amazon'un bize sunduğu veri toplama, işleme, saklama, analiz etme ve arşivleme amacıyla tasarlanmış büyük veri servislerini şöyle listeleyebiliriz: Amazon Kinesis Amazon Elastic MapReduce (EMR) Amazon Athena Amazon Machine Learning Amazon DynamoDB Amazon [...]

Amazon, Athena, AWS, Big Data, DynamoDB, ElasticSearch, EMR, Kinesis, Machine Learning, Quicksight, Redshift

Apache Hadoop 3.0

10 Ocak 2018 Hakan İtler Hadoop 0



Uzun zamandır 2.x sürümüyle devam eden Hadoop projesinde 13 Aralık 2017 tarihinde yeni sürüm olan Hadoop 3.0 yayınlandı. Bu yeni sürümde ne gibi özellikler olduğunu bu yazıda özetlemeye çalışacağız; Java Update: Bütün proje minimum Java sürümü Java 8 olacak şekilde derlendi. Dolayısı ile Hadoop 3.0 kullanmak isteyenlerin Java sürümünü de yükseltmesi gerekecek. HDFS Erasure Coding: Yeni [...]

Kategoriler

- Big Data (11)
- Cloud (3)
- docker (1)
- Doğal Dil İşleme (2)
- ElasticSearch (4)
- Graph (1)
- Haberler (7)
- Hadoop (24)
- HBase (1)
- Kitap (1)
- Lucene / Solr (3)
- Nosql (12)
- Ölçeklenebilirlik (2)
- Polyglot (1)
- Sunum (1)
- Veri Bilimi (2)
- Veri Madenciliği (4)
- Yapay Öğrenme (3)

Son Yazılar

- Amazon EMR ile Spark
- Amazon EMR
- AWS ile Big Data
- Apache Hadoop 3.0
- Big Data Teknolojilerine Hızlı Giriş
- Günlük Hayatta Yapay Zekâ Teknikleri – Yazı Dizisi (1)

Kategoriler

Big Data (11)

Cloud (3)

docker (1)

Dođal Dil İşleme (2)

ElasticSearch (4)

Graph (1)

Haberler (7)

Hadoop (24)

HBase (1)

Kitap (1)

Lucene / Solr (3)

Nosql (12)

Ölçeklenebilirlik (2)

Polyglot (1)

Sunum (1)

Veri Bilimi (2)

Veri Madenciliđi (4)

Yapay Öğrenme (3)

```
> urLs
```

```
[1] "http://devveri.com/"
```


Amazon EMR ile Spark

18 Ocak 2018 Hakan İtler Cloud 0



Bu yazıda Amazon EMR üzerinde bir Spark uygulamasının nasıl çalıştırabileceğinden bahsedeceğim. Eğer EMR ile ilgili bir önceki yazıyı okumadıysanız bu yazıyı, AWS Big Data teknolojileri ile ilgili genel bilgi için de bu yazıyı okuyabilirsiniz. EMR üzerinde çalıştıracağımız örnek uygulamada daha önce defalarca kullandığım NYSE verisini kullanacağım. Tab karakterleri ile ayrılmış bu dosya içerisinde günlük borsa [...]

AWS, Big Data, Cloud, EMR, Hadoop, S3, Spark

Amazon EMR

13 Ocak 2018 Hakan İtler Cloud 0

Amazon EMR Nedir? Amazon Elastic MapReduce (EMR), büyük veri işlemeyi kolaylaştırmak amacıyla Amazon tarafından yönetilen, içerisinde Hadoop, Spark gibi açık kaynaklı büyük veri teknolojilerini içeren bir servistir. Aslında temelinde AWS (servisler: Hadoop, Kinesis) bu makale için tek kaynak olan aramak, nesnelik yazılımları, yüklemek gibi işlemleri

Kategoriler

Big Data (11)
Cloud (3)
docker (1)
Doğal Dil İşleme (2)
ElasticSearch (4)
Graph (1)
Haberler (7)
Hadoop (24)
HBase (1)
Kıtao (1)

```
Elements Console Sources Network Performance Memory Application Security Audits
<!-- Start: Post -->
<div class="post-1970 post type-post status-publish format-standard has-post-thumbnail hentry category-cloud tag-aws tag-big-data tag-cloud tag-emr tag-hadoop tag-s3 tag-spark">
  <h2>
    <a href="http://devveri.com/cloud/amazon-emr-ile-spark" rel="bookmark" title="Permanent Link to Amazon EMR ile Spark">Amazon EMR ile Spark/</a> == $0
  </h2>
  <p class="post-meta"></p>
  
  <p></p>
  <p class="more"></p>
  <p class="tags"></p>
</div>
<!-- End: Post -->
<!-- Start: Post -->
<div class="post-1909 post type-post status-publish format-standard has-post-thumbnail hentry category-cloud tag-amazon tag-aws tag-big-data tag-emr tag-hadoop tag-spark"></div>
<!-- End: Post -->
<!-- Start: Post -->
<div class="post-1934 post type-post status-publish format-standard has-post-thumbnail hentry category-cloud tag-amazon tag-athena tag-aws tag-big-data tag-dynamodb tag-elasticsearch tag-emr tag-kinesis tag-machine-learning tag-quicksight tag-redshift"></div>
<!-- End: Post -->
<!-- Start: Post -->
<div class="post-1928 post type-post status-publish format-standard has-post-thumbnail hentry category-hadoop tag-erasure-coding tag-hadoop tag-hadoop-3-0 tag-yarn"></div>
<!-- End: Post -->
<!-- Start: Post -->
<div class="post-1899 post type-post status-publish format-standard has-post-thumbnail hentry category-big-data tag-aws tag-big-data tag-buyuk-veri tag-hadoop tag-spark"></div>
<!-- End: Post -->
<!-- Start: Post -->
<div class="post-1881 post type-post status-publish format-standard hentry category-yapay-ogrenme tag-classification tag-regression tag-siniflandirma tag-support-vector-machines tag-svm tag->
html body.home.blog div.body div.content div.main div.post-1970.post.type-post.status-publish.format-standard.has-post-thumbnail.hentry.category-cloud.tag-aws.tag-big-data.tag-cloud.tag-emr.tag-hadoop.tag-s3.tag-spark h2 a
```

Console



top

Filter

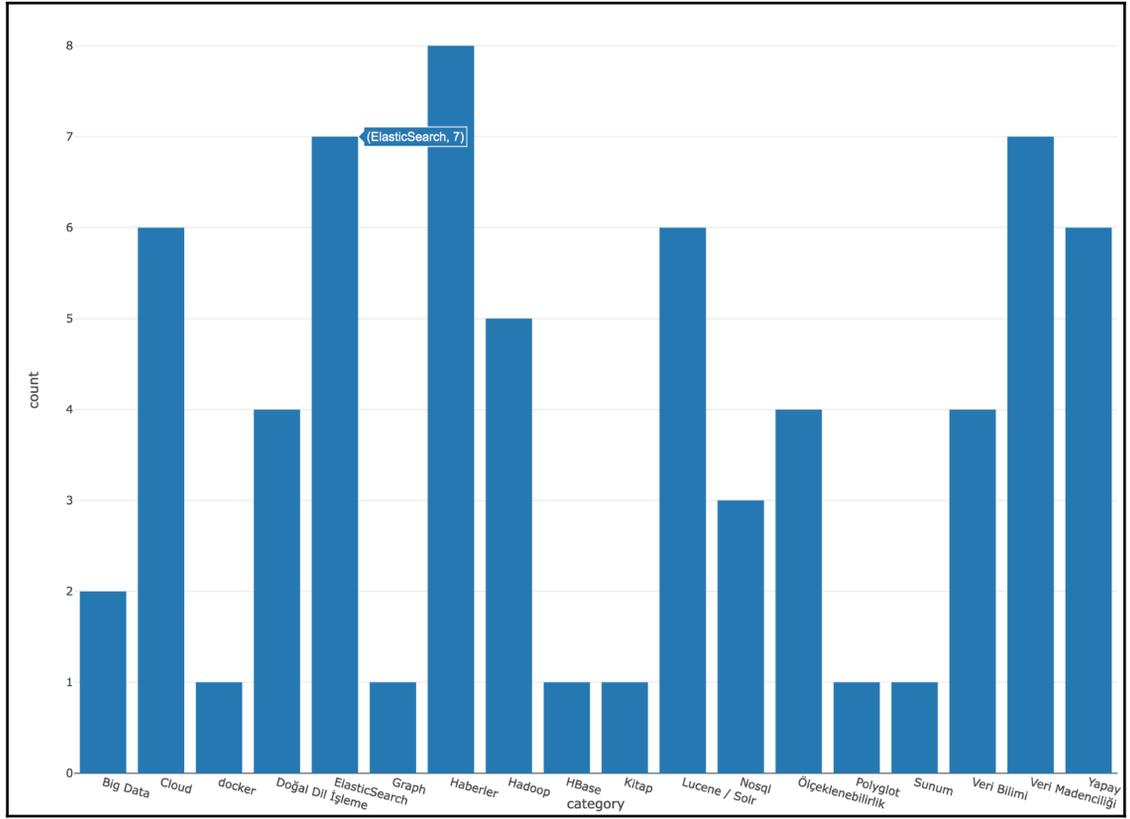
Default levels

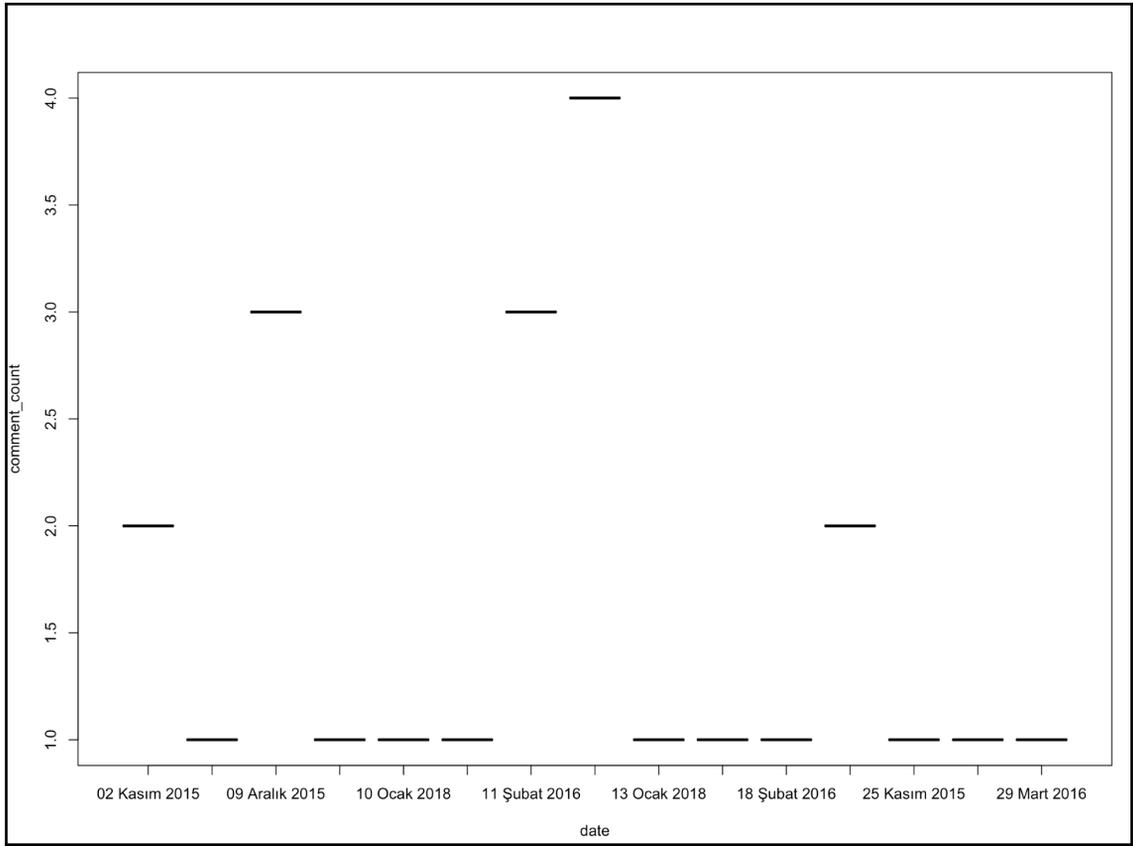
Group similar

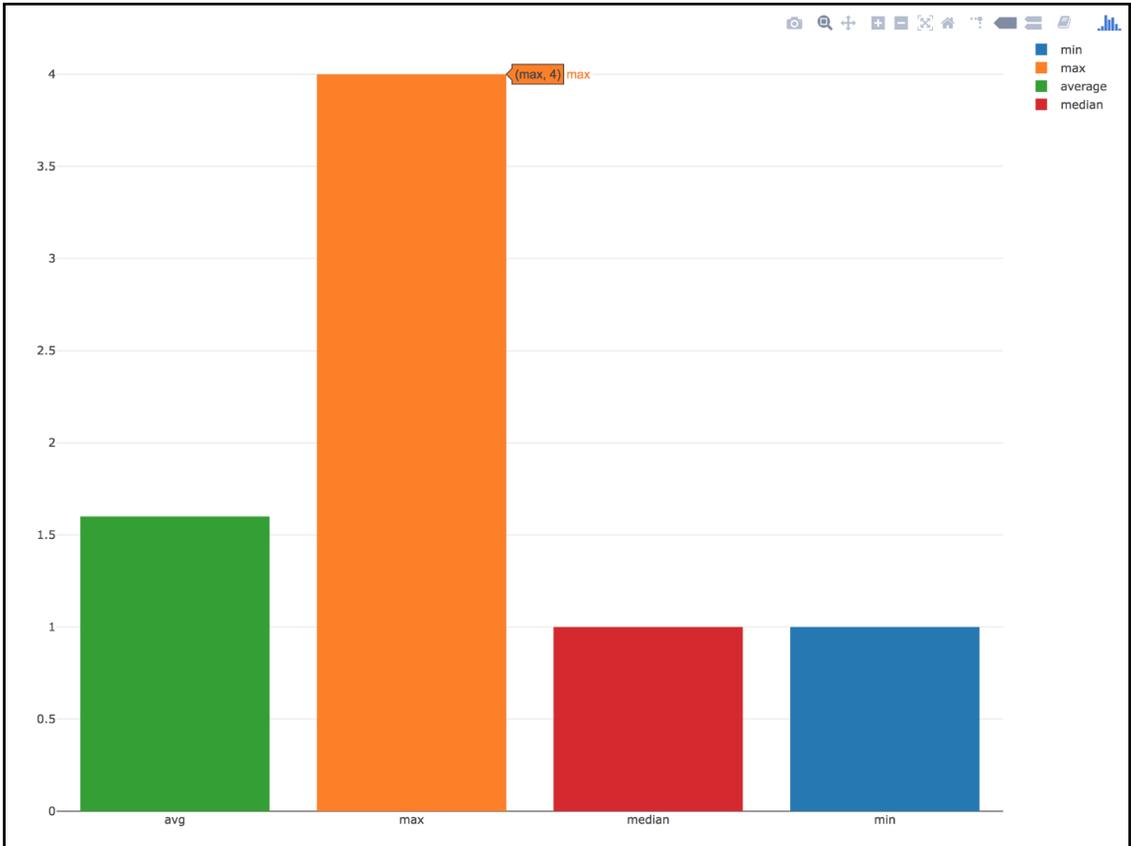
> \$x({'/html/body/div[3]/div/div[1]/div/h2/a/text()')}

< ▶ (15) [text, text, text, text, text, text, text, text, text, text, text, text, text, text]


```
Console
top Filter Default levels Group similar
> $('html/body/div[3]/div/div[1]/div/p[1]/span[4]/a/text()')
< ▼ (15) [text, text, text, text, text, text, text, text, text, text, text, text] ⓘ
  ▼ 0: text
    assignedSlot: null
    baseURI: "http://devveri.com/"
    ▶ childNodes: NodeList []
    data: "0"
    firstChild: null
    isConnected: true
    lastChild: null
    length: 1
    nextElementSibling: null
    nextSibling: null
    nodeName: "#text"
    nodeType: 3
    nodeValue: "0"
    ▶ ownerDocument: document
    ▶ parentElement: a
    ▶ parentNode: a
    previousElementSibling: null
    previousSibling: null
    textContent: "0"
    wholeText: "0"
    ▶ __proto__: Text
```







Chapter 4: Web Scraping with Relenium



```
Console Terminal x
~/...
* DONE (R Selenium)
Warning message:
In strptime(x, fmt, tz = "GMT") :
  unknown timezone 'zone/tz/2018e.1.0/zoneinfo/Europe/Warsaw'
> library("R Selenium")
> rD <- rsDriver()
checking Selenium Server versions:
BEGIN: PREDOWNLOAD
BEGIN: DOWNLOAD
Creating directory: /Users/olgunaydin/Library/Application Support/b...
Downloading binary: https://www.googleapis.com/download/storage/v1/...

Creating directory: /Users/olgunaydin/Library/Application Support/b...
Downloading binary: https://www.googleapis.com/download/storage/v1/...

Creating directory: /Users/olgunaydin/Library/Application Support/b...
Downloading binary: https://www.googleapis.com/download/storage/v1/...
```

Chrome is being controlled by automated test software.

Himanshu Sharma on Instagram x

Secure https://www.instagram.com/p/BiFW2XFD8CM/?hl=en&tagged=packtpub

Instagram

Search

Log In Sign Up

Packt Publishing Limited ★★★★★ 5

Kali Linux - An Ethical Hacker's Cookbook

#1 Best Seller in Networks & System Administration

80% off



Format **Paperback** >

M.R.P.: ₹-4,441
Price: ₹-949

Oxhimanshu • Follow

Oxhimanshu To everyone who said I can't do it. Started from the bottom now we're here. #nevergiveup #believeinyourself #neversettle #keepgrinding #bestseller #packtpub #author #hackers

Buy it from here -> <https://goo.gl/f9u2ky>

sri_shyam Wow!! ❤️👍👎

hmalviya9 Great stuff Himanshu . You deserve the success

Oxhimanshu @hmalviya9 Bhai thanks :)

zaainlive Himaanshu!! Love you bro !! Bless you

rea2der You are a champ

anushka_anu01 #proudbitch 🙌👍 #keepgoing

Oxhimanshu @zaainlive thank you bro

46 likes

APRIL 27

Log in to like or comment.

Chrome is being controlled by automated test software.



Instagram

Search

Log In

Sign Up

Packt Publishing Limited

Kali Linux - An Ethical Hacker's Cookbook

#1 Best Seller In Networks & System Administration

80% off



Format Paperback

M.R.P.: ₹4,444
Price: ₹949

★★★★★ 5



Ozhimanshu • Follow

do it.
Started from the bottom now we're here.
#nevergiveup #believeinyourself
#neversettle #keepgrinding #bestseller
#packtpub #author #hackers

Buy it from here -> <https://goo.gl/f9u2ky>
sri_shyam Wow! 🤔🤔

hmalviya9 Great stuff Himanshu . You
deserve the success

Ozhimanshu @hmalviya9 Bhai thanks :)
zaainlive Himaanshu!! Love you bro !! Bless
you

rea2der You are a champ

anushka_anu01 #proudbitch 🤔🤔
#keepgoing

Ozhimanshu @zaainlive thank you bro
dhruv_2204 Buying it! 🤔



46 likes

APRIL 27

Log in to like or comment.

← → ↻ 🔒 Not Secure | devveri.com

Chrome is being controlled by automated test software.

DEVVERİ.COM

Boğulacaksan büyük veride boğul!

Big Data

Hadoop

NoSQL

Yazarlar

Amazon EMR ile Spark

📅 18 Ocak 2018

👤 Hakan İtler

📁 Cloud

💬 0



Bu yazıda Amazon EMR üzerinde bir Spark uygulamasının nasıl çalıştırabileceğinden bahsedeceğim. Eğer EMR ile ilgili bir önceki yazıyı okumadıysanız bu yazıyı, AWS Big Data teknolojileri ile ilgili genel bilgi için de bu yazıyı okuyabilirsiniz. EMR üzerinde çalıştıracağımız örnek uygulamada daha önce defalarca kullandığım NYSE verisini kullanacağım. Tab karakterleri ile ayrılmış bu dosya içerisinde günlük borsa [...]

🔗 AWS, Big Data, Cloud, EMR, Hadoop, S3, Spark



Amazon EMR ile Spark

18 Ocak 2018 Hakan İtler Cloud 0

Bu yazıda Amazon EMR üzerinde bir Spark uygulamasının nasıl çalıştırabileceğinden bahsedeceğim. Eğer EMR ile ilgili bir önceki yazıyı okumadıysanız [bu yazıyı](#), AWS Big Data teknolojileri ile ilgili genel bilgi için de [bu yazıyı](#) okuyabilirsiniz.



EMR üzerinde çalıştıracığımız örnek uygulamada daha önce defalarca kullandığım NYSE verisini kullanacağım. Tab karakterleri ile ayrılmış bu dosya içerisinde günlük borsa işlemleri ile ilgili veriler bulunuyor. Spark 2.x ile birlikte gelen CSV desteği ile kolayca DataFrame haline getireceğimiz veriyi, S3 üzerinden okuyup üzerinde aggregation yapan bir SQL komutu çalıştırıp sonucu Parquet formatında tekrar S3 üzerine yazacağız.

Başlamadan önce işleyeceğimiz dosyayı indirip S3 üzerine atmamız gerekiyor. Bu işlemi elle AWS S3 arayüzünden yapabileceğiniz gibi komut satırından da yapabilirsiniz. Bunun için **awscli** uygulamasının yüklenmiş, ayarlarını yapmış olmanız ve S3 üzerindeki bucket'ı yaratmış olmanız gerekiyor.

```
1 wget https://s3.amazonaws.com/hw-sandbox/tutorial1/NYSE-2000-2001.tsv.gz
2 aws s3 cp NYSE-2000-2001.tsv.gz s3://datapyro-main/test/
```

EMR içerisinde kaynak yönetimi için **YARN** kullandığından Spark uygulamasında master parametresini "yarn" olarak belirtmek gerekiyor. Uygulamamız input parametresi ile S3 üzerinde okuyacağı veriyi, output parametresi ile de S3 üzerinde sonucun yazılacağı yeri alıyor.

Kategoriler

Big Data (11)

Cloud (3)

docker (1)

Doğal Dil İşleme (2)

ElasticSearch (4)

Graph (1)

Haberler (7)

Hadoop (24)

HBase (1)

Kitap (1)

Lucene / Solr (3)

Nosql (12)

Ölçeklenebilirlik (2)

AWS, Big Data, Cloud, EMR, Hadoop, S3, Spark

Amazon EMR

Amazon EMR

13 Ocak 2018 Hakan litter Cloud 0

Amazon EMR Nedir?

Amazon Elastic MapReduce (EMR), büyük veri işlemeyi kolaylaştırmak amacıyla Amazon tarafından yönetilen, içerisinde Hadoop, Spark gibi açık kaynaklı büyük veri teknolojilerini içeren bir servistir. Aslında temelde AWS üzerinde Hadoop kümesi kurmak için tek tek sunucuları açmak, gerekli yazılımları yüklemek gibi işlemleri otomatik olarak yapmaktadır. Bu sayede tek tık ile bir kümeyi kurabileceğiniz gibi, işiniz bittiğinde de yine tek tık ile kümeyi silebilirsiniz. Yani, bu esneklik sayesinde çok daha büyük bir kümeyi veriyi işleyeceğiniz zaman yaratıp, veriyi işledikten sonra da kümeyi kapatarak kaynakları boşuna kullanmadan hesaplı ve hızlı bir şekilde sonuca ulaşabilirsiniz.

Neden Amazon EMR Kullanırım?

Büyük veri teknolojilerine yeni başlayanların en çok sıkıntı çektiği konulardan birisi uygun çalışma ortamına sahip olmamak. Sanal makineler yardımıyla kendi bilgisayarımızda bu teknolojileri öğrenmek mümkün, ancak öğrendiklerinizi gerçek hayatta uygulamak için gerçek bir kümeye sahip olmak maalesef çok mümkün olmuyor. Amazon EMR bu problemi ortadan kaldırıyor. İsteddiğiniz kümeyi kurup, kurcalayıp, kullandığınız kadar ödeyebilirsiniz.

POC aşamasında sunucu satın alıp kaynak ayırmak ile uğraşmak yerine, EMR ile POC çalışmanızı ya da demonuzu yapabilirsiniz.

Firmalar için ise, kullandığınız kadar ödemek gerçekten anlamlı. Örneğin günde bir kere çalışan ve TB'larca veriyi işlemeyi gerektiren batch bir iş için sürekli ayakta duracak bir Hadoop kümesine ihtiyacınız yok. Ücretlendirme kümeyi oluşturan EC2 sunucularının büyüklük ve sayısına göre yapılıyor. Ayrıca **spot instance** kullanma şansınız da var.

Spot instance kabaca, başkaları tarafından kiralanmış ancak boşta duran EC2 sunucularını normalden çok daha ucuza kullanmamızı sağlayan bir özellik. Bu avantaja rağmen ihtiyaç durumunda bu sunucuların sizden geri alınması dezavantajı var. Fakat büyük veri teknolojileri hata toleranslı olarak tasarlandıkları için, uygun tasarlanmış bir küme içerisindeki bazı sunucuların kapanması çalışan sistemi etkilemeyeceği için EMR ile spot instance kullanmak gerçekten iyi bir seçenek.

Kategoriler

Big Data (11)

Cloud (3)

docker (1)

Doğal Dil İşleme (2)

ElasticSearch (4)

Graph (1)

Haberler (7)

Hadoop (24)

HBase (1)

Kitap (1)

Lucene / Solr (3)

Nosql (12)

Ölçeklenebilirlik (2)

Polyglot (1)

Sunum (1)

Veri Bilimi (2)

Veri Madenciliği (4)

Amazon EMR ile Spark

18 Ocak 2018 Hakan İler Cloud 0



Bu yazıda Amazon EMR üzerinde bir Spark uygulamasının nasıl çalıştırabileceğinden bahsedeceğim. Eğer EMR ile ilgili bir önceki yazıyı okumadıysanız bu yazıyı, AWS Big Data teknolojileri ile ilgili genel bilgi için de bu yazıyı okuyabilirsiniz. EMR üzerinde geliştireceğimiz örnek uygulamada daha önce defalarca kullandığım NYSE verisini kullanacağım. Tabii karakterleri ile ayrılmış bu dosya içerisinde günlük borsa [...]

AWS, Big Data, Cloud, EMR, Hadoop, S3, Spark

Amazon EMR

13 Ocak 2018 Hakan İler Cloud 0

Amazon EMR Nedir? Amazon Elastic MapReduce (EMR), büyük veri işlemeyi kolaylaştırmak amacıyla Amazon tarafından yönetilen, içerisinde Hadoop, Spark gibi açık kaynaklı büyük veri teknolojilerini içeren bir servistir. Aslında temelde AWS üzerinde Hadoop kümesi kurmak için tek tek sunucuları açmak, gerekli yazımları yüklemek gibi işlemleri otomatik olarak yapmaktadır. Bu sayede tek tık ile bir kümeyi kurabileceğinizin [...]

Amazon, AWS, Big Data, EMR, Hadoop, Spark

AWS ile Big Data

11 Ocak 2018 Hakan İler Cloud 0



Bu yazı teknolojinin öncüsü olan Amazon Web Services'te bize birçok büyük veri teknolojisini sunarak ve uygun maliyetli olarak test etme ve kullanma fırsatı sağlıyor. Amazon'un bize sunduğu veri toplama, işleme, saklama, analiz etme ve arşivleme amacıyla tasarlanmış büyük veri servislerini şöyle listeleyebiliriz: Amazon Kinesis Amazon Elastic MapReduce (EMR) Amazon Athena Amazon Machine Learning Amazon DynamoDB Amazon [...]

Amazon, Athena, AWS, Big Data, DynamoDB, ElasticSearch, EMR, Kinesis, Machine Learning, QuickSight, Redshift

Kategoriler

Big Data (11)

Cloud (3)

docker (1)

Doğal Dil İşleme (2)

ElasticSearch (4)

Graph (1)

Haberler (7)

Hadoop (24)

HBase (1)

Kıtap (1)

Lucene / Solr (3)

NoSQL (12)

Ölçeklenebilirlik (2)

Polyglot (1)

Sunum (1)

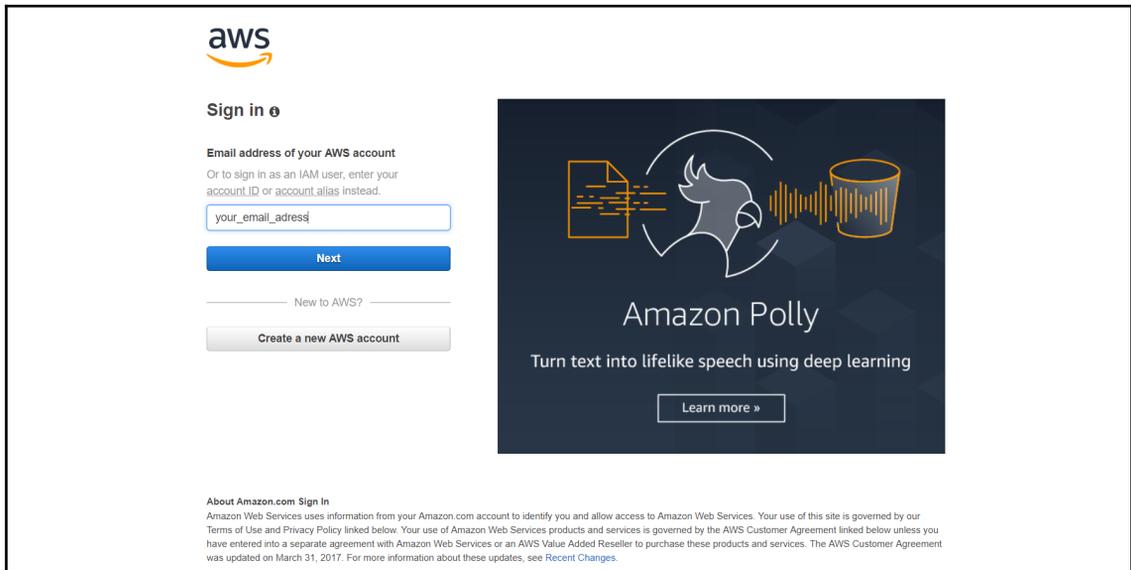
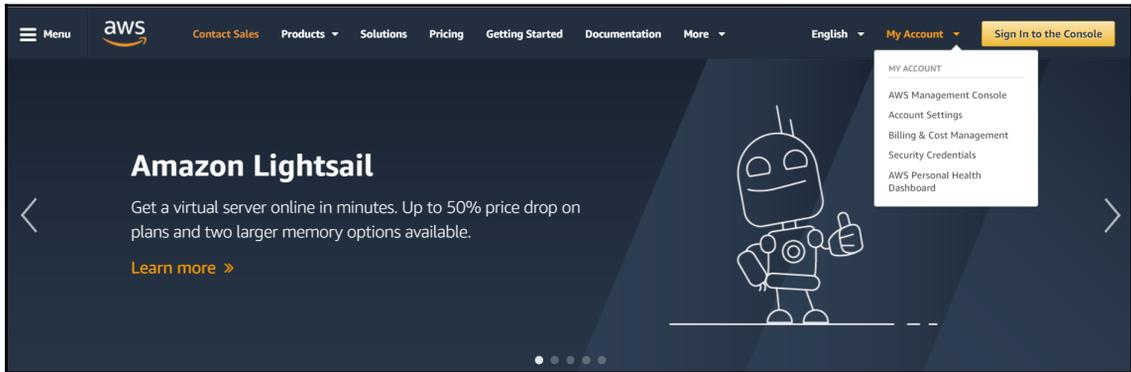
Veri Bilimi (2)

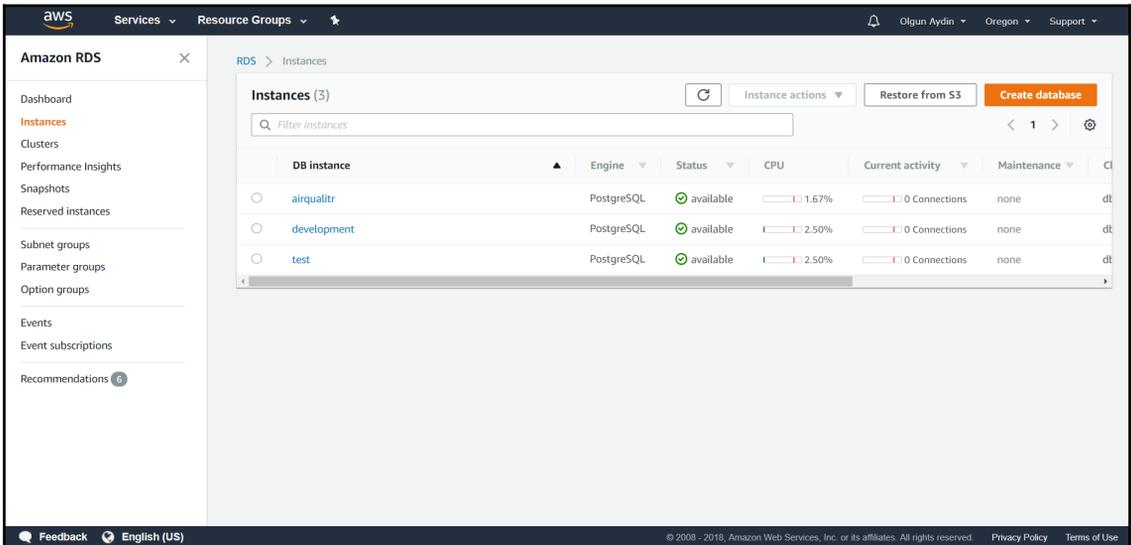
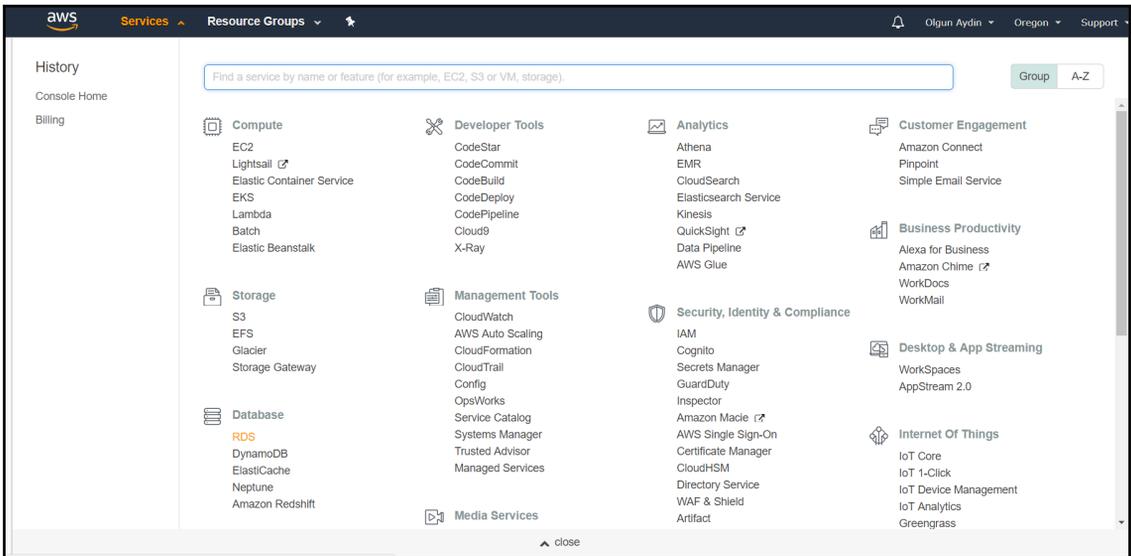
Veri Madenciliği (4)

Yapay Öğrenme (3)

Son Yazılar

Chapter 5: Storing Data and Creating Cronjob





aws Services Resource Groups Olgun Aydin Oregon Support

RDS > Database > Create database

Step 1 **Select engine**

Step 2 Specify DB details

Step 3 Configure advanced settings

Select engine

Engine options

- Amazon Aurora
- MySQL
- MariaDB
- PostgreSQL
- Oracle
- Microsoft SQL Server

PostgreSQL

PostgreSQL is a powerful, open-source object-relational database system with a strong reputation of reliability, stability, and correctness.

- High reliability and stability in a variety of workloads.
- Advanced features to perform in high-volume environments.
- Vibrant open-source community that releases new features multiple times per year.

Feedback English (US) © 2020 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

aws Services Resource Groups Olgun Aydin Oregon Support

RDS > Database > Create database

Step 1 [Select engine](#)

Step 2 **Specify DB details**

Step 3 Configure advanced settings

Specify DB details

Instance specifications

Estimate your monthly costs for the DB Instance using the [AWS Simple Monthly Calculator](#)

DB engine
PostgreSQL

License model [Info](#)
postgresql-license

DB engine version [Info](#)
PostgreSQL 10.4-R1

Free tier

The Amazon RDS Free Tier provides a single db.t2.micro instance as well as up to 20 GiB of storage, allowing new AWS customers to gain hands-on experience with Amazon RDS. Learn more about the RDS Free Tier and the instance restrictions [here](#).

- Only enable options eligible for RDS Free Usage Tier [Info](#)

DB instance class [Info](#)

db.t2.micro — 1 vCPU, 1 GiB RAM

Multi-AZ deployment [Info](#)

- Create replica in different zone
Creates a replica in a different Availability Zone (AZ) to provide data redundancy, eliminate I/O freezes, and minimize latency spikes during system backups.
- No

Storage type [Info](#)

General Purpose (SSD)

Allocated storage

20

GiB

(Minimum: 20 GiB, Maximum: 20 GiB) Higher allocated storage [may improve](#) IOPS performance.

Settings

DB instance identifier [Info](#)

Specify a name that is unique for all DB instances owned by your AWS account in the current region.

test2

DB instance identifier is case insensitive, but stored as all lower-case, as in "mydbinstance". Must contain from 1 to 63 alphanumeric characters or hyphens (1 to 15 for SQL Server). First character must be a letter. Cannot end with a hyphen or contain two consecutive hyphens.

Master username [Info](#)

Specify an alphanumeric string that defines the login ID for the master user.

test

Master Username must start with a letter. Must contain 1 to 63 alphanumeric characters.

Master password [Info](#)

.....

Confirm password [Info](#)

.....

Master Password must be at least eight characters long, as in "mypassword". Can be any printable ASCII character except "/", "", or "@".

Cancel

Previous

Next

Database options

Database name [Info](#)

If you do not specify a database name, Amazon RDS does not create a database.

Port [Info](#)

TCP/IP port the DB instance will use for application connections.

DB parameter group [Info](#)

Option group [Info](#)

IAM DB authentication [Info](#)

Enable IAM DB authentication

Manage your database user credentials through AWS IAM users and roles.

Disable

The screenshot shows the AWS Management Console interface for an Amazon RDS database instance. The top navigation bar includes the AWS logo, 'Services', 'Resource Groups', and user information for 'Olga Rydin' in 'Oregon'. The main content area is titled 'development' and includes a 'Summary' section with the following details:

- Engine: PostgreSQL 9.6.6
- DB instance class: db.t2.micro
- DB instance status: available
- Pending maintenance: none

Below the summary is a 'CloudWatch (20)' section with a legend for 'development' and a search bar. The bottom section contains four performance metrics charts for the time period 08/28 14:00 to 14:30:

- CPU Utilization (Percent):** A line chart showing utilization fluctuating between approximately 1% and 3%.
- DB Connections (Count):** A line chart showing a constant count of 1 connection.
- Free Storage Space (MB):** A line chart showing a constant free space of approximately 4,500 MB.
- Freeable Memory (MB):** A line chart showing a constant freeable memory of approximately 500 MB.

The left sidebar contains navigation options: Dashboard, Instances, Clusters, Performance Insights, Snapshots, Reserved instances, Subnet groups, Parameter groups, Option groups, Events, Event subscriptions, and Recommendations.

The screenshot displays the AWS Management Console for an Amazon RDS instance. The interface is organized into several sections:

- Configurations:** Includes ARN (arn:aws:rds:us-west-2:639049327697:db:development), Engine (PostgreSQL 9.6.6), License Model (Postgresql License), Created Time (Thu Jul 27 14:15:58 GMT+200 2017), DB Name (development), Username (olgun), Option Group (default:postgres9-6), Parameter group (default:postgres9.6 (in-sync)), and Resource ID (db-RR4AH6MEVEZFUJWNZLQFO3JE).
- Security and network:** Shows Availability zone (us-west-2a), VPC (vpc-9c8c6cf8), Subnet group (default), Subnets (subnet-badff9e3, subnet-52bc9325, subnet-ed24ce89), Security groups (rds-launch-wizard-1 (sg-e65b6a80) (active)), and Publicly accessible (Yes).
- Instance and IOPS:** Lists Instance Class (db.t2.micro), Storage Type (General Purpose (SSD)), Storage (5 GiB), Availability and durability (Multi AZ: No, DB instance status: available), Backup and Restore (Automated backups: Enabled (7 Days), Backup window: 13:14-13:44 UTC (GMT)), and Latest restore time (August 28, 2018 at 2:38:44 PM UTC+2).
- Maintenance details:** Shows Auto minor version upgrade (Yes), Maintenance window (tue:06:04-tue:06:34 UTC (GMT)), Pending Modifications (None), Pending maintenance (none), and Encryption details (Encryption enabled: No).

The screenshot shows the RStudio IDE with an R script open. The script contains the following code:

```
18 - {
19   #reading main url
20   h <- read_html(urls[i])
21
22   #getting rating
23   r<- html_nodes(h, xpath = '//div[@itemprop="ratingvalue"]/text()')
24
25
26   #getting date
27   d<- html_nodes(h, xpath = '//div[@itemprop="published"]/text()')
28
29   #saving results
30   date<- rbind(date,as.matrix(as.character(d)))
31   rating<- rbind(rating,as.matrix(as.character(r)))
32 }
33 final_data<- data.frame(date,rating)
34
35
36
37
```

The 'Tools' menu is open, showing options such as 'Install Packages...', 'Check for Package Updates...', 'Version Control', 'Shell...', 'Terminal', 'Addins', 'Keyboard Shortcuts Help', 'Project Options...', and 'Global Options...'. The 'Addins' option is highlighted, and a sub-menu is visible with 'Browse Addins...' selected.

Addins

[? Using RStudio Addins](#)

Package	Name	Description
reprex	Render reprex	Run <code>`reprex::reprex()`</code> to prepare a reproducible example for sharing.
taskscheduleR	Schedule R scripts on Windows	Use Windows task scheduler to schedule your R scripts (Windows)
tfruns	Training Run	Execute a training run with the current source document
tfruns	View Latest Run	View the most recent training run
tfruns	View Run History	View all training runs

[Keyboard Shortcuts...](#)[Execute](#)[Cancel](#)

The screenshot shows the RStudio interface with an R script open. A context menu is displayed over the xpaths in the script. The menu options are:

- REPREX
- Render reprex
- TASKSCHEDULER
- Schedule R scripts on Windows
- TFRUNS
- Use Windows task scheduler to schedule your R scripts (Windows)
- Training Run
- View Latest Run
- View Run History

```
17 for(i in 1:n)
18 {
19   #reading main url
20   h <- read_html(urls[i])
21
22   #getting rating
23   r<- html_nodes(h, xpath = '//div[@itemprop="reviewRating"]//span[@itemprop="ratingValue"]/text()')
24
25   #getting date
26   d<- html_nodes(h, xpath = '//time[@itemprop="datePublished"]/text()')
27
28   #saving results
29   date<- rbind(date,as.matrix(as.character(d)))
30   rating<- rbind(rating,as.matrix(as.character(r)))
31 }
32
33 final_data<- data.frame(date,rating)
34
35
36
```

Task Scheduler

Cancel Schedule your Rscript fast and easy Done

Choose your Rscript

Browse... devveri.R

Upload complete

Rscript repo: location where Rscripts will be copied to schedule + location of logs

C:/Users/olgun_aydin/Documents/R/win-libr;

Schedule:

ONCE

MONTHLY

WEEKLY

DAILY

HOURLY

MINUTE

ONLOGON

ONIDLE

Task checking: does the task already exist?

Schedulingtask for Rscript devveri.R does not exist yet

Start date:

2018-10-22

Hour start

13:47

Additional arguments to Rscript

Date format of your locale

%d/%m/%Y

Create task

Upload and create Stop or Delete