
통계로 보는 역사학

Seongmin Mun

| 목차

- 웹 사이트 분석하기
- 웹 크롤링 수행하기
- 지정된 사이트에서 정보 가져오기

웹 사이트 분석하기

웹사이트 분석하기

웹사이트란?

웹사이트(영어: website, 문화어: 웹브싸이트)는 인터넷 프로토콜 기반의 네트워크에서 도메인 이름이나 IP 주소, 루트 경로만으로 이루어진 일반 URL을 통하여 보이는 웹 페이지 (Web Page)들의 의미 있는 묶음이다.

The screenshot shows the Naver homepage at https://www.naver.com. The page features the Naver logo, a search bar, and various news and service links. The Network tab in the DevTools is open, displaying a list of network requests. The table includes columns for Name, Status, Type, Initiator, Size, Time, and Waterfall. Key requests include images for banners and news cards, and several JavaScript files like nsd1484475.png, jndo_v180212_ie1.js, and various versions of core.min.js. The total transfer time is 8.7 ms and the DOMContentLoaded time is 496 ms.

Name	Status	Type	Initiator	Size	Time	Waterfall
nsd1484475.png	200	png	jndo_v180212_ie1	1.9 KB	33 ms	
nsd14534617.png	200	png	jndo_v180212_ie1	16.4 KB	16 ms	
nsd14463367.png	200	png	jndo_v180212_ie1	16.9 KB	11 ms	
nsd192546763.png	200	png	jndo_v180212_ie1	5.3 KB	12 ms	
nsd1370343431.png	200	png	jndo_v180212_ie1	16.7 KB	13 ms	
nsd113515807.png	200	png	jndo_v180212_ie1	1.9 KB	11 ms	
nsd14405515.png	200	png	jndo_v180212_ie1	16.5 KB	12 ms	
nsd154151664.png	200	png	jndo_v180212_ie1	16.4 KB	12 ms	
nsd162046351.png	200	png	jndo_v180212_ie1	5.3 KB	22 ms	
778ca17a88bef0e3b2b8_201811301101...	206	media	Other	1.5 MB	8.72 s	
fxview?eu=EU10202350&calp=1&oj=bb...	200	text/plain	pc.veta.core.min.js...	241 B	38 ms	
fxview?eu=EU10202350&calp=1&oj=bb...	200	text/plain	pc.veta.core.min.js...	241 B	16 ms	
fxview?eu=EU10202350&calp=1&oj=bb...	200	text/plain	pc.veta.core.min.js...	241 B	14 ms	
fxview?eu=EU10202350&calp=1&oj=bb...	200	text/plain	pc.veta.core.min.js...	241 B	14 ms	
fxview?eu=EU10202350&calp=1&oj=bb...	200	text/plain	pc.veta.core.min.js...	241 B	14 ms	

Console tab highlights the new features in Chrome 70, such as Live Expressions and Autocomplete Conditional Breakpoints.

Bottom right corner: new 70 logo.

서버란?

서버(영어: server)는 클라이언트에게 네트워크를 통해 정보나 서비스를 제공하는 컴퓨터(server computer) 또는 프로그램(server program)을 말한다. 특히, 서버에서 동작하는 소프트웨어를 서버 소프트웨어(server software)라 한다.

외부에서 요청하면 규칙대로 정보를 제공하는 것이다.



브라우저란?

웹 브라우저(Web Browser, 문화어: 열람기)는 웹 서버에서 쌍방향 통신하는 HTML 문서나 파일과 연동하고 출력하는 응용 소프트웨어이다.

서버가 주는 것들을 사용자에게 보여주는 것이다.



웹사이트 분석하기

웹 사이트의 구조

xml양식으로 작성된 웹 브라우저들이 이해하는 표준 문서이다.

simple.html

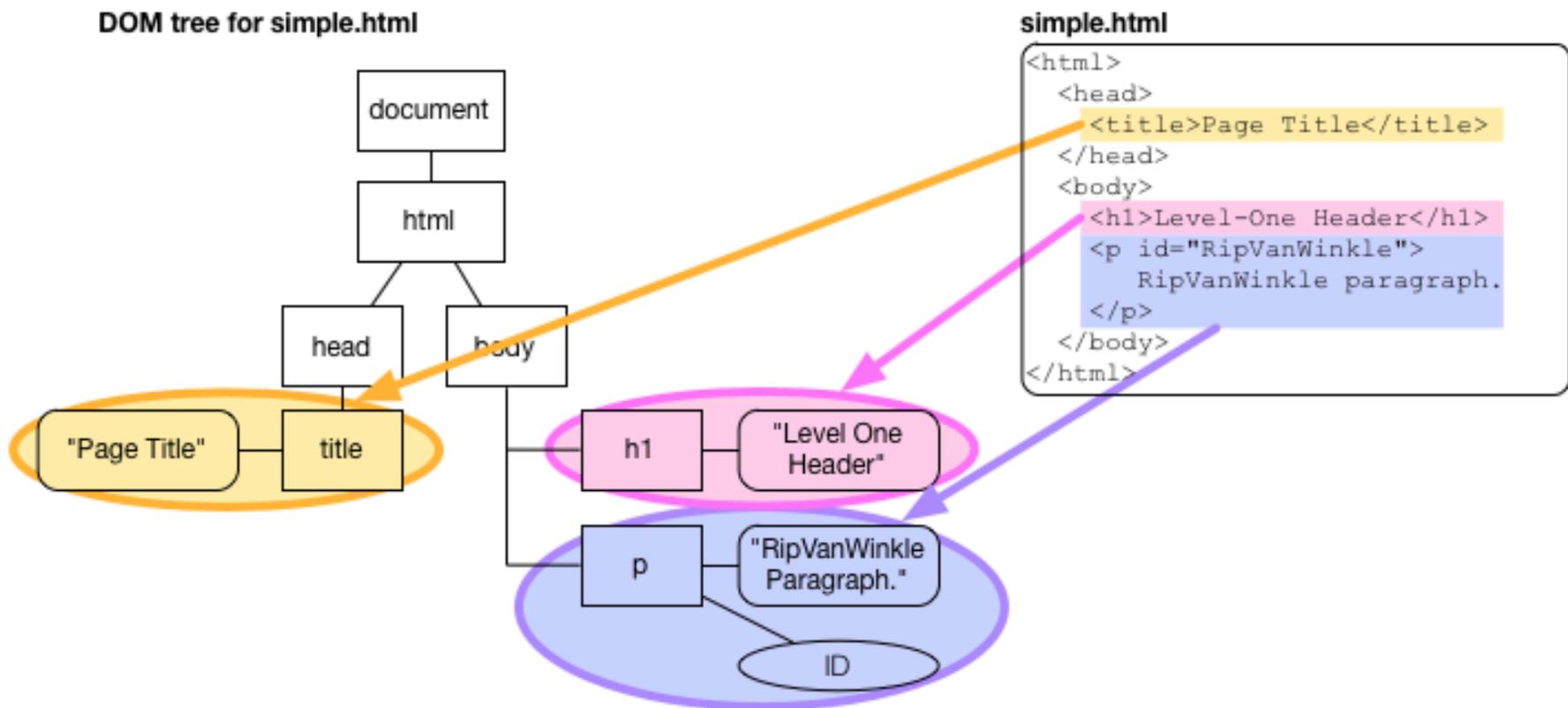
```
<html>
  <head>
    <title>Page Title</title>
  </head>
  <body>
    <h1>Level-One Header</h1>
    <p id="RipVanWinkle">
      RipVanWinkle paragraph.
    </p>
  </body>
</html>
```

A Browser Window



웹 사이트의 구조

xml양식으로 작성된 웹 브라우저들이 이해하는 표준 문서이다.



웹 크롤링 수행하기

웹 크롤링 수행하기

1. 정규표현식을 활용한 웹크롤링

데이터 정리 및 패키지 설치

```
> #memory & previous_works
> gc()
      used   (Mb) gc trigger  (Mb) limit (Mb) max used   (Mb)
Ncells 3180984 169.9     6125541 327.2        NA 4237698 226.4
Vcells 6085608  46.5    12425283  94.8        32768 12419957  94.8
> rm(list=ls())
> #Encoding_mac
> Sys.setlocale(category = "LC_CTYPE", locale = "ko_KR.UTF-8")
[1] "ko_KR.UTF-8"
```

웹 크롤링 수행하기

1. 정규표현식을 활용한 웹크롤링

데이터 정리 및 패키지 설치

```
> install.packages("RCurl")
trying URL 'https://cran.rstudio.com/bin/macosx/el-capitan/contrib/3.5/RCurl_1.95-4.11.tgz'
Content type 'application/x-gzip' length 985623 bytes (962 KB)
=====
downloaded 962 KB
```

The downloaded binary packages are in
/var/folders/h4/hwr_1_hn7tz9zrwpjn9b6xcc0000gn/T//Rtmp05XntC downloaded_packages

```
> install.packages("stringr")
trying URL 'https://cran.rstudio.com/bin/macosx/el-capitan/contrib/3.5/stringr_1.3.1.tgz'
Content type 'application/x-gzip' length 186353 bytes (181 KB)
=====
downloaded 181 KB
```

The downloaded binary packages are in
/var/folders/h4/hwr_1_hn7tz9zrwpjn9b6xcc0000gn/T//Rtmp05XntC downloaded_packages

```
> library(stringr)
> library(RCurl)
Loading required package: bitops
```

웹 크롤링 수행하기

웹 사이트 소스코드 가져오기

웹 크롤링 수행하기

웹 사이트 소스를 클래스 명을 기준으로 정제하기

웹 크롤링 수행하기

정규표현식을 사용해서 데이터 추출하기

웹 크롤링 수행하기

정규표현식을 사용해서 데이터 추출하기

```
> #영화 리스트
> for(i in 1:5){
+   print(str_match(hold_3[[1]][i], "[가-힣'\\s]+\\<\\/\\w+\\>"))
+ }
[,1]
[1,] NA
[,1]
[1,] "베일리 어게인</strong>"
[,1]
[1,] "신비한 동물들과 그린델왈드의 범죄</strong>"
[,1]
[1,] "완벽한 타인</strong>"
[,1]
[1,] "보헤미안 랩소디</strong>"
```

웹 크롤링 수행하기

정규표현식을 사용해서 데이터 추출하기

```
> #영화 리스트_정제
> movielist <- NULL
> for(i in 1:5){
+   before <- str_match(hold_3[[1]][i],"[가-힣'\\s]+\\<\\/\\w+\\>")
+   print(gsub("</strong>","",before))
+   target <- gsub("</strong>","",before)
+   if(!is.na(target)){
+     movielist <- c(movielist,target)
+   }
+ }
[,1]
[1,] NA
[,1]
[1,] "베일리 어게인"
[,1]
[1,] "신비한 동물들과 그린델왈드의 범죄"
[,1]
[1,] "완벽한 타인"
[,1]
[1,] "보헤미안 랩소디"
> movielist
[1] "베일리 어게인" "신비한 동물들과 그린델왈드의 범죄" "완벽한 타인" "보헤미안 랩소디"
```

2. GET요청 방식을 활용한 웹크롤링

데이터 정리 및 패키지 설치

```
> #memory & previous_works  
> gc()  
      used   (Mb) gc trigger  (Mb) limit (Mb) max used   (Mb)  
Ncells 3180984 169.9    6125541 327.2        NA 4237698 226.4  
Vcells 6085608  46.5    12425283 94.8        32768 12419957  94.8  
> rm(list=ls())  
> #Encoding_mac  
> Sys.setlocale(category = "LC_CTYPE", locale = "ko_KR.UTF-8")  
[1] "ko_KR.UTF-8"
```

웹 크롤링 수행하기

데이터 정리 및 패키지 설치

```
> install.packages('rvest')
trying URL 'https://cran.rstudio.com/bin/macosx/el-capitan/contrib/3.5/rvest_0.3.2.tgz'
Content type 'application/x-gzip' length 880013 bytes (859 KB)
=====
downloaded 859 KB
```

The downloaded binary packages are in

```
/var/folders/h4/hwr_1_hn7tz9zrwpjn9b6xcc0000gn/T//Rtmp05XntC downloaded_packages
> #install.packages('rvest')
> library(rvest)
Loading required package: xml2
```

웹 크롤링 수행하기

문서 제목 크롤링하기

```
> url <- "http://terms.naver.com/entry.nhn?docId=1623779&categoryId=49629&cid=49629"
> Gorea <- read_html(url)
> GoreaTitle <- html_nodes(Gorea, css='.headword_title')
> head(GoreaTitle)
{xml_nodeset (1)}
[1] <div class="headword_title">\n\t\t\t\t<p class="cite">\n\t\t\t\t\t<a href="list.nhn?categoryId=62161&so=st4.asc" onclick="clickcr(thi ...
> str(GoreaTitle)
List of 1
$ :List of 2
..$ node:<externalptr>
..$ doc :<externalptr>
..- attr(*, "class")= chr "xml_node"
- attr(*, "class")= chr "xml_nodeset"
> GoreaTitle[1]
{xml_nodeset (1)}
[1] <div class="headword_title">\n\t\t\t\t<p class="cite">\n\t\t\t\t\t<a href="list.nhn?categoryId=62161&so=st4.asc" onclick="clickcr(thi ...
> GoreaTitle[1] %>% html_nodes('h2') %>% html_text()
[1] "태조 원년(918) 무인년"
> pagetitle<-GoreaTitle[1] %>% html_nodes('h2') %>% html_text()
> pagetitle
[1] "태조 원년(918) 무인년"
```

웹 크롤링 수행하기

텍스트 정보 크롤링 하기

```
> #텍스트 크롤링하기
> GoreaInfos <- html_nodes(Gorea, css='.size_ct_v3')
> head(GoreaInfos)
{xml_nodeSet (1)}
[1] <div id="size_ct" class="size_ct_v3">\n\n\t\t\t\t\t<script type="text/javascript" src="https://audioapi.nmv.naver.com/re
> str(GoreaInfos)
List of 1
$ :List of 2
..$ node:<externalptr>
..$ doc :<externalptr>
.. - attr(*, "class")= chr "xml_node"
- attr(*, "class")= chr "xml_nodeSet"
> GoreaInfos[1]
{xml_nodeSet (1)}
[1] <div id="size_ct" class="size_ct_v3">\n\n\t\t\t\t\t<script type="text/javascript" src="https://audioapi.nmv.naver.com/re
> GoreaInfos[1] %>% html_nodes('.txt') %>% html_text()
[1] "● 원년 여름 6월병진일. 태조가 포정전(布政殿)1)에서 즉위하여 국호를 고려(高麗)2)라 하고 연호를 고쳐 천수(天授)라고 했다. 정사일. 다음과 같은 조서를 내렸다.“? 붕괴할 때에 도적의 무리들을 제거하고 점차로 영토를 넓혀갔다. 그러나 천하를 아우르기도 전에 갑자기 잔혹한 폭정으로 백성들을 다스렸으며 간사함을 가장 옳은 것으로 여기고 단으로 삼았다. 요역과 부세가 번거롭고 과중하여 인구는 줄어들고 농토는 텅 비게 되었다. 그런데도 오히려 궁실만은 크고 으리으리하며 옛 제도를 준수하지 않고 힘든 부역은 어나게 된 것이다. 더군다나 함부로 연호를 정하고 황제를 칭했으며 처자를 살육한 죄는 천지간에 용납되지 못할 일이며 귀신과 사람이 함께 노할 일로서 왕업의 기반을 송두리째 은 공들이 추대하는 마음에 힘입어 가장 높은 자리에 올랐으니 낡은 풍속을 고쳐 모든 것을 다함께 새롭게 만들려 한다. 마땅히 법도와 규범을 혁신하는 길4)을 쫓을 것이며 가계로 삼으리라. 임금과 신하는 물과 물고기처럼 서로 어울려 즐거움[魚水之歡6)]을 같이 할 것이며 온 천하는 태평시대의 경사[晏清之慶7)]를 함께 누릴지니 나라의 모든 백성! 신하들이 절을 올리고 사례했다.“전 임금의 통치 기간에는 선량한 사람들이 악독한 피해를 입고 죄 없는 사람들이 잔혹한 학대를 받는 통에 남녀노소가 모두 불만에 싸여 원한을
```

웹 크롤링 수행하기

텍스트 정보 크롤링 하기

```
> pagetext<-GoreaInfos[1] %>% html_nodes('.txt') %>% html_text()  
> pagetext
```

[1] "● 원년 여름 6월병진일. 태조가 포정전(布政殿)1)에서 즉위하여 국호를 고려(高麗)2)라 하고 연호를 고쳐 천수(天授)라고 봉고할 때에 도적의 무리들을 제거하고 점차로 영토를 넓혀갔다. 그러나 천하를 아우르기도 전에 갑자기 잔혹한 폭정으로 백성들을 단으로 삼았다. 요역과 부세가 번거롭고 과중하여 인구는 줄어들고 농토는 텅 비게 되었다. 그런데도 오히려 궁실만은 크고 으리으리나게 된 것이다. 더군다나 함부로 연호를 정하고 황제를 칭했으며 처자를 살육한 죄는 천지간에 용납되지 못할 일이며 귀신과 사는 공들이 추대하는 마음에 힘입어 가장 높은 자리에 올랐으니 낡은 풍속을 고쳐 모든 것을 다함께 새롭게 만들려 한다. 마땅히 법계로 삼으리라. 임금과 신하는 물과 물고기처럼 서로 어울려 즐거움[魚水之歡6)]을 같이 할 것이며 온 천하는 태평시대의 경사[晏] 신하들이 절을 올리고 사례했다.“전 임금의 통치 기간에는 선량한 사람들이 악독한 피해를 입고 죄 없는 사람들이 잔혹한 학대를 받히 목숨을 보전하여 성스럽고 현명한 임금을 만날 수 있게 되었으니 어찌 힘을 다하여 성은에 보답하지 않겠습니까?”무오일. 왕이 , 경의 고향 청주(青州)는 땅이 기름지고 호걸이 많았기 때문에 변란을 일으킬까 우려한 나머지 그 곳 사람들의 씨를 말리려 했다. 가 없는데도 형구를 찬 채 끌려오고 있으니 경은 빨리 가서 그들8)을 고향으로 돌려보내도록 하라.”경신일. 마군장군(馬軍將軍) 혼직을 설치하고 직책을 분담하는 일에는 유능한 사람을 임명하는 것이 중요하며, 세상을 이롭게 하고 백성을 평안하게 하는 일에는 해질 까닭이 없는 것이다. 짐이 외람되게도 천명[景命9)]을 받아 큰 계획을 품어하려니, 높은 지위를 차지하게 되면 편안하기 어렵로 알지 못하고 관리를 선발함에 실수가 많아 어진 사람을 누락시켰다는 탄식을 야기시키고 선비를 얻는 도리에 어긋남이 있을까 걱신료들이 모두 그 임무를 잘 감당할 수 있으면 현재 훌륭한 치적을 이룩할 수 있을 뿐 아니라 후대의 칭송까지 받을 수 있는 것�이 뽑아 적재적소에 배치해야 할 것이다. 온 나라 사람들은 모두 짐의 뜻을 헤아릴지어다.”이에 따라 한찬(韓粲) 김행도(金行灝)11)

```
> str(pagetext)  
chr [1:2] "● 원년 여름 6월병진일. 태조가 포정전(布政殿)1)에서 즉위하여 국호를 고려(高麗)2)라 하고 연호를 고쳐 천수(天授)라" | __truncated__ ...
```

3. POST요청 방식을 활용한 웹크롤링

데이터 정리 및 패키지 설치

```
> #memory & previous_works  
> gc()  
      used   (Mb) gc trigger  (Mb) limit (Mb) max used   (Mb)  
Ncells 3180984 169.9    6125541 327.2        NA 4237698 226.4  
Vcells 6085608  46.5    12425283 94.8        32768 12419957  94.8  
> rm(list=ls())  
> #Encoding_mac  
> Sys.setlocale(category = "LC_CTYPE", locale = "ko_KR.UTF-8")  
[1] "ko_KR.UTF-8"
```

웹 크롤링 수행하기

데이터 정리 및 패키지 설치

```
> install.packages("httr")
trying URL 'https://cran.rstudio.com/bin/macosx/el-capitan/contrib/3.5/httr_1.3.1.tgz'
Content type 'application/x-gzip' length 447943 bytes (437 KB)
=====
downloaded 437 KB
```

The downloaded binary packages are in

```
/var/folders/h4/hwr_1_hn7tz9zrwpjn9b6xcc0000gn/T//Rtmpou1C3a downloaded_packages
> install.packages("rvest")
trying URL 'https://cran.rstudio.com/bin/macosx/el-capitan/contrib/3.5/rvest_0.3.2.tgz'
Content type 'application/x-gzip' length 880013 bytes (859 KB)
=====
downloaded 859 KB
```

The downloaded binary packages are in

```
/var/folders/h4/hwr_1_hn7tz9zrwpjn9b6xcc0000gn/T//Rtmpou1C3a downloaded_packages
> library(httr)
> library(rvest)
Loading required package: xml2
```

웹 크롤링 수행하기

웹 사이트 소스코드 가져오기

```
> comment.url <- 'https://www.kinds.or.kr/news/newsResult.do'
> comment <- POST(comment.url,
+                     body=list(
+                         pageInfo= "bksMain",
+                         indexName= "news",
+                         keyword= "도서관",
+                         searchScope= "1",
+                         searchFtr= "1",
+                         startDate= "2018-09-02",
+                         endDate= "2018-12-02",
+                         sortMethod= "date",
+                         contentLength= "100",
+                         startNo= "1",
+                         resultNumber= "10",
+                         period= "3month"))
> outresult = content(comment, "text")
> outresult
[1] "{\"keywordJson\":null,\"sessionUSID\":null,\"topmenuoff\":null,\"resultState\":\"detailSe
ams\":{\"indexName\":\"news\",\"query\":null,\"keyword\":\"도서관\",\"byLine\":null,\"categoryCc
2\",\"endDate\":\"2018-12-02\",\"incidentCode\":null,\"dateCode\":null,\"providerCode\":null,\n
rue\",\"searchFtr\":\"1\",\"sortMethod\":\"date\",\"searchScope\":\"1\",\"contentLength\":\"10
30\",\"maxSentences\":\"5\",\"maxNewsCount\":\"30\",\"keyWindowSize\":\"50\",\"edgeWindowSize\l
off\":null,\"researchQuery\":null,\"listMode\":null,\"categoryTab\":null,\"newsId\":null,\"de
\"quotationKeyword1\":null,\"quotationKeyword2\":null,\"quotationKeyword3\":null,\"filterProvi
de\":null,\"filterDateCode\":null,\"filterAnalysisCode\":null,\"providerLinkPage\":null,\"prin
inTodayPersonYn\":null,\"period\":\"3month\",\"sectionDiv\":null,\"nation\":null,\"keywordJson
```

웹 크롤링 수행하기

관련 정보 가져오기

```
> #관련정보들 가져오기 (json)
> jsonlite::fromJSON(outresult)
```

```
$keywordJson
```

```
NULL
```

```
$sessionUSID
```

```
NULL
```

```
$topmenuoff
```

```
NULL
```

```
$resultState
```

```
[1] "detailSearch"
```

```
$sessionUUID
```

```
NULL
```

```
$pageInfo
```

```
[1] "newsResult"
```

```
$params
```

```
$params$indexName
```

```
[1] "news"
```

웹 크롤링 수행하기

관련 정보 가져오기

```
> data <- jsonlite::fromJSON(outresult)
> summary(data)
      Length Class Mode
keywordJson        0   -none- NULL
sessionUSID        0   -none- NULL
topmenuoff         0   -none- NULL
resultState        1   -none- character
sessionUUID        0   -none- NULL
pageInfo           1   -none- character
params              59  -none- list
keywordFilterJson  0   -none- NULL
resultSet          13  -none- list
bigkindsMode       1   -none- character
> data$resultSet$resultList$CONTENT
```

[1] "지난 30일 사망한 조지 허버트 워커 부시 전 대통령의 장례식은 제럴드 포드 전 대통령의 사망 이후 미국에서 11년만의 하지 않았지만, 5일 워싱턴 국립 대성당에서 장례식을 진행한 후 부인 바버라 부시 여사가 묻힌 텍사스 칼리지스테이션 조지 부시 게 된다."

[2] "극단 '쨍'이 오는 16일 안양시 석수도서관에서 낭독극 'L 의견에 간다'는 프랑스 비평가가 뽑은 '주목받지 못..공연은 일요일인 16일 오후 1시, 오후 3시 총 2회 안양시립석수다. 도서관 홈페이지와 방문·전화 등을 통해 사전 접수 할 수 있

[3] "지난달 30일 북구 오토밸리컨벤션에서 열린 북구 도서관>도서관 우수자원봉사자와 작은도서관 운영

[4] "자유한국당 이채익(남구갑?사진) 의원은 지난 30일 국회 산업통상자원중소벤처기업위원회 회의실에서 열린 에너지특별위원회 13년 조합 이사장에 취임한 허인회 이사장은 2015년 2월까지 발효현미 상품 판매, SH공사의 작은 <font color='고 밝혔다."

지정된 사이트에서 정보 가져오기

지정된 사이트에서 정보 가져오기

- 사이트 주소 : <http://sjw.history.go.kr/main.do>

The screenshot shows the homepage of the 'Chungyeon Ilgi' website (<http://sjw.history.go.kr/main.do>). The page features a large image of historical documents at the bottom, with text from the documents visible. Above the image, there's a search bar with placeholder text '인기검색어' and a green '검색' button. To the left, there are sections for '인조 - 순종' and '기타'. The main content area displays a grid of historical records. On the right side of the screen, the Chrome DevTools Network tab is open, showing a list of resources loaded by the page. The Network tab includes a timeline at the top and a detailed table below it. The table lists various files with their status, type, initiator, size, time taken, and a waterfall chart. The total number of requests is 44, and the total transferred data is 9.0 MB.

Name	Status	Type	Initiator	Size	Time	Waterfall
main_visual.png	200	png	main.do	416 KB	562 ms	
logo_mow.png	200	png	main.do	3.3 KB	238 ms	
bar01.gif	200	gif	main.do	296 B	26 ms	
NanumGothicBold.woff	200	font	main.do	2.3 MB	2.69 s	
arrow_green.png	200	png	main.do	3.4 KB	35 ms	
list_line.png	200	png	main.do	3.0 KB	74 ms	
bullet_dot.gif	200	gif	main.do	295 B	137 ms	
ko.js	200	script	tiny_mce_src.js:8715	10.1 KB	17 ms	
editor_template_src.js	200	script	tiny_mce_src.js:8715	2.3 KB	19 ms	
editor_plugin_src.js	200	script	tiny_mce_src.js:8715	30.4 KB	55 ms	
ui.js	200	script	tiny_mce_src.js:8715	10.1 KB	34 ms	
ui.css	200	stylesheet	tiny_mce_src.js:4218	922 B	65 ms	
content.css	200	stylesheet	tiny_mce_src.js:117...	3.2 KB	21 ms	
NanumGothic.woff	200	font	main.do	1.2 MB	1.56 s	
NanumBarunGothic-YetHangul.woff	200	font	main.do	2.6 MB	3.06 s	

44 requests | 9.0 MB transferred | Finish: 5.46 s | DOMContentLoaded: 845 ms | Load: 5.53 s

Console What's New X

Highlights from the Chrome 70 update

Live Expressions in the Console
Pin expressions to the top of the Console to monitor their values in real-time.

Highlight DOM nodes during Eager Evaluation
Type an expression that evaluates to a node to highlight that node in the viewport.

Autocomplete Conditional Breakpoints
Type expressions quickly and accurately.

Performance panel optimizations
Faster loading and processing of Performance recordings.

More reliable debugging
Bug fixes for sourcemaps and blackboxing.

Debug Node.js apps with ndb
Detect and attach to child processes, place breakpoints before modules are required, edit files within DevTools, and more.

Learn more Close