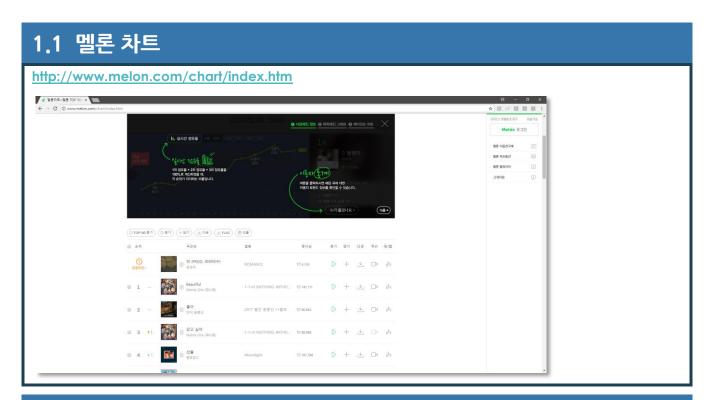






# 1. 크롤링 할 웹사이트 확보



### 1.2 HTML 소스 확인

검색을 통해 차트 정보가 있는 곳 확인





## 2. 웹크롤링

#### 2.1 모듈 설치 2.1.1 간단한 웹서버에 요청 모듈 명령프롬프트에서 pip install requests 5 명령 프롬프트 Microsoft Windows [Version 10.0.15063] (c) 2017 Microsoft Corporation. All rights reserved. C:\Users\Ki-ung>pip install requests Collecting requests Using cached requests-2.18.4-py2.py3-none-any.whl Requirement already satisfied: chardet<3.1.0,>=3.0.2 in c:\users\ki-ung\appdata\ local\programs\python\python36\lib\site-packages (from requests) Requirement already satisfied: idna<2.7,>=2.5 in c:\u00c4users\u00fcki-ung\u00fcappdata\u00fclocal\u00fcp rograms\u00fcpython\u00fcpython36\u00fclib\u00fcsite-packages (from requests) Requirement already satisfied: certifi>=2017.4.17 in c:\u00fcusers\u00fcki-ung\u00fcappdata\u00fclocal\u00fcap al\mprograms\mpython\mpython36\mathrid3 \text{ucitivity}=2011.4.11 \text{in c.\masers\mathrid3} \text{dig\mappaata\mappa} \text{100} \\ al\mappappams\mappata\mappaat Installing collected packages: requests Successfully installed requests-2.18.4 C:\Users\Ki-ung> 2.1.2 HTML Parser 모듈 명령프롬프트에서 pip install bs4 명령 프롬프트 × Microsoft Windows [Version 10.0.15063] (c) 2017 Microsoft Corporation. All rights reserved. C:\Users\Ki-ung>pip install bs4 Collecting bs4 Using cached bs4-0.0.1.tar.gz Collecting beautifulsoup4 (from bs4) Using cached beautifulsoup4-4.6.0-py3-none-any.whl Installing collected packages: beautifulsoup4, bs4 Running setup.py install for bs4 ... done Successfully installed beautifulsoup4-4.6.0 bs4-0.0.1 C:\Users\Ki-ung>\_





### 2.2 HTML 다운받기

```
### Content - Type | Content = "text/html; charset = utf-8"/>

### weta http-equiv="X-UA-Compatible" content = "IE = edge, chrome = 1" />

### weta name="keywords" content = "IE = edge, chrome = 1" />

### weta name="keywords" content = "IE = edge, chrome = 1" />

### weta name = "IE = edge, chrome = 1" />

### weta name = "IE = edge, chrome = 1" />

### weta name = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwords" content = "IE = edge, chrome = 1" />

### weta name = "Newwor
```

#### 2.3 HTML Parsing

