

## Project: Google Play Store

Google Play Store is a rich source of information about the various aspects of mobile apps. This includes information about customers, businesses, and technical information in the form of permission details, descriptions, version information and reviews. On July 24th, 2013, Google announced that they had crossed one million applications, and estimates now state about one and a half million apps and billions of downloads. This rich source of information has been largely untapped due to there being no public access to this data. This data is a potential source of information about customer sentiment and perception of apps. App developers can use such information to build better apps and respond to customer complaints faster. A few companies like ‘AppAnnie’ and ‘Distimo’ provide limited access to such data, and a few researchers have written APIs to scrape some selected amount of data to suit their needs. However, such a selective sampling of data using the custom written APIs can lead to an analytical bias.

The operation of every API is based on the “*docid*” parameter.

The “*docid*” is the unique identifier for each application. This docid is available in the URL of every application page when you visit Google Play Store. For example, for the Facebook application the URL is

“https://play.google.com/store/apps/details?id=com.facebook.katana” and the docid is “com.facebook.katana”.

\*\*\*\*\*

## Scope of the Project

The project has two major folds:

**Task-1:** Extracting the Name, Docid and the app URL/Link from the Google Play Store.

**Task-2:** Extracting the Name, Docid and rank of the application for any of the website of your choice that rank the Google Play Store apps on either weekly basis or monthly basis. I need to approve the site of your choice before you start with the this task

\*\*\*\*\*

## Environment

**Language:** Python

**API:** Any, provided it is free and had a decent limit on download quota. Before using an API, please consult with me.

**Database:** MySQL Version (8.0.12), community edition.

**This is not the final draft of the project.**

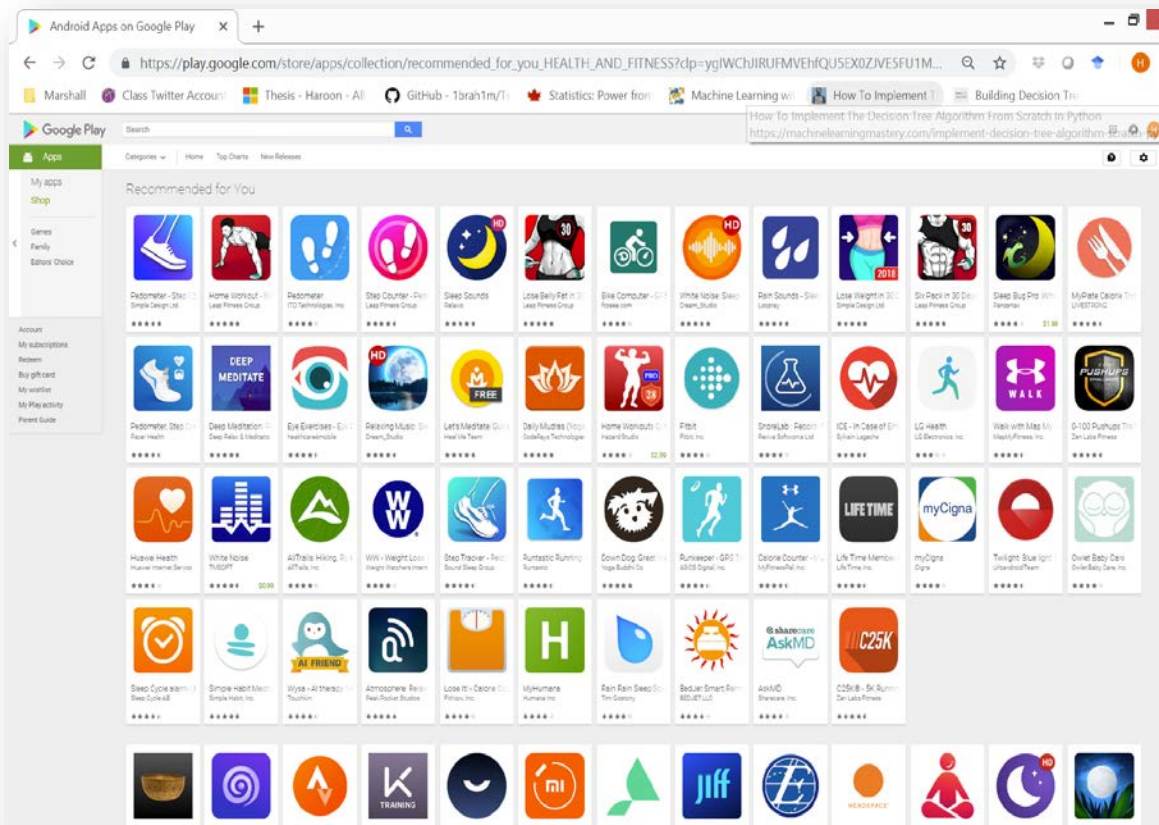
## Task-1

You need to produce a command line crawler, that will take at the minimum, a URL as an argument and fetch the App name, its corresponding Docid, its respective URL and save it into the DB.

For example, if the following URL is provide to the crawler:

[https://play.google.com/store/apps/collection/recommended\\_for\\_you\\_HEALTH\\_AND\\_FITNESS?clp=ygIWChJIRUFMVEhfQU5EX0ZJVE5FU1MQAAQ%3D%3D%3AS%3AANO1ljKVgyA&hl=en](https://play.google.com/store/apps/collection/recommended_for_you_HEALTH_AND_FITNESS?clp=ygIWChJIRUFMVEhfQU5EX0ZJVE5FU1MQAAQ%3D%3D%3AS%3AANO1ljKVgyA&hl=en) (Page shown below)

The crawler should retrieve all the Apps name, Doc ID and app URLs listed on the page (minimum of 30).



**Hint:** if you view the page source, i.e., view-

source:[https://play.google.com/store/apps/collection/recommended\\_for\\_you\\_HEALTH\\_AND\\_FITNESS?clp=ygIWChJIRUFMVEhfQU5EX0ZJVE5FU1MQAAQ%3D%3D%3AS%3AANO1ljKVgyA&hl=en](https://play.google.com/store/apps/collection/recommended_for_you_HEALTH_AND_FITNESS?clp=ygIWChJIRUFMVEhfQU5EX0ZJVE5FU1MQAAQ%3D%3D%3AS%3AANO1ljKVgyA&hl=en) (Page shown below)

**This is not the final draft of the project.**

You should have access to all the required information to get the job done.

```
view-source:https://play.google.com/store/apps/collection/recommended_for_you_HEALTH_AND_FITNESS?cp=ygiWChJIRUFMVEhfQ...
style="display:none"></div> <div js1="$x 7;$t t-fIi-cv3Fi0I;$x 0; track-impression id-action-bar-settings-button" aria-label="Settings" data-uitype="111" title="Settings"
jsan="t-fIi-cv3Fi0I,7.action-bar-link,7.id-track-click,7.id-track-impression,7.id-action-bar-settings-
button,0.aria-label,0.data-uitype,0.title"> <a href="/settings/?authuser=0" title="Settings" js1="$x 1;"> <div
class="action-bar-dropdown-top"> <span class="dropdown-icon"></span> </div> </a> </div> </div> </div> </div>
</div> </div><div class="body-content-loading-overlay" style="display:none"><div class="body-content-loading-
spinner"></div></div><div class="body-content" id="body-content" role="main"><div class="outer-container"><div
class="inner-container"><div class="main-content"><div class="id-cluster-container"><div class="id-cluster-
container"><div class="cluster id-track-impression normal square-cover apps show-all id-show-
cover-type-css="square-cover" data-fetch-start="49" data-layout="NORMAL" data-orig
track-impression normal square-cover apps show-all" data-server-
cookie="CAMiGAgBogETCPDvR0KNkN4CFdQ1yQod3QkNUg==" data-short-classes="cluster id-t
cover apps show-all" data-uitype="400"><div class="cluster-heading"><h2>Recommend
class="subtitle" data-server-cookie="CAMiGAgBogETCPDvR0KNkN4CFdQ1yQod3QkNUg==" dat
<div class="id-card-list card-list two-cards"><div class="card no-rationale square
docid="pedometer.stepcounter.calorieburner.pedometerforwalking" data-original-clas
cover apps small" data-short-classes="card no-rationale square-cover apps small" d
rationale square-cover apps small"> <div class="card-content id-track-click id-tra
docid="pedometer.stepcounter.calorieburner.pedometerforwalking" data-server-
cookie="CAIaWAO/Ej0KN3B1ZG9tZXR1ci5zdGVwY291bnR1ci5jYWxvcml1YnVybmVyLnB1ZG9tZXR1cm
3gIViDXJCh3dCQ1S0gA=" data-uitype="500"> <a class="card-click-target" data-server
cookie="CAIaWAO/Ej0KN3B1ZG9tZXR1ci5zdGVwY291bnR1ci5jYWxvcml1YnVybmVyLnB1ZG9tZXR1cm
3gIViDXJCh3dCQ1S0gA=" data-uitype="500" href="/store/apps/details?
id=pedometer.stepcounter.calorieburner.pedometerforwalking" aria-hidden="true" tabindex="-1"></a> <div
class="cover"> <div class="cover-image-container"> <div class="cover-outer-align"> <div class="cover-inner-
align">  </div> </div> </div> <a class="card-click-target"
href="/store/apps/details?id=pedometer.stepcounter.calorieburner.pedometerforwalking" aria-label=" Pedometer
```

URL:  
[href="/store/apps/details?id=pedometer.stepcounter.calorieburner.pedometerforwalking"](/store/apps/details?id=pedometer.stepcounter.calorieburner.pedometerforwalking)  
Note that the link store in the href tag leads to the app page. Also, after /store/apps/details? You can find the docid, besides other tags .

By simply knowing the appropriate tags, you can scrape the corresponding values for App name, URL and Docid.

\*\*\*\*\*

## Task-2

Pickup any website of your choice that maintains the Top-Chart(s) of google Play Store (Andriod) apps and scrape the App Name, its Docid, its rank and store them into the database.