



درس پایگاه داده پیشرفته

گزارش پروژه

تحلیل داده‌های Google Play Store

نام استاد: دکتر نعمت بخش

نام دانشجو: علی ابراهیمی

شماره دانشجویی: ۴۰۳۳۶۴۴۰۰۱

بهمن ۱۴۰۳

بخش اول: جمع‌آوری و آماده‌سازی داده‌ها

در این پروژه، داده‌های مربوط به Google Play Store از Kaggle دانلود و سپس برای استفاده در پایگاه داده PostgreSQL، پیش پردازش و آماده شدند.

۱. دانلود مجموعه داده‌ها از Kaggle

ابتدا، مجموعه داده‌ای از سایت Kaggle دانلود شد که شامل اطلاعات اپلیکیشن‌های موجود در Google Play Store بود.

۲. تجزیه و تحلیل داده‌های گم‌شده

در این مرحله، ابتدا بررسی شد که کدام ستون‌ها دارای مقادیر گم‌شده هستند و میزان گم‌شدن داده‌ها برای هر ستون محاسبه شد. این اطلاعات به ما کمک کرد تا تصمیم بگیریم که کدام داده‌ها باید پر شوند و کدام‌ها باید حذف شوند. پس از شناسایی داده‌های گم‌شده، داده‌هایی که برای تحلیل‌های اصلی ضروری بودند، اصلاح شدند.

۳. حذف رکوردهای تکراری یا نامعتبر

رکوردهای تکراری که ممکن بود چندین بار برای یک اپلیکیشن ثبت شده باشند، شناسایی و حذف شدند. همچنین رکوردهایی که شامل داده‌های نادرست یا ناقص در فیلدهای حیاتی مانند نام اپلیکیشن، شناسه توسعه‌دهنده، تعداد نصب‌ها و حجم بودند، حذف شدند.

۴. استانداردسازی فرمت داده‌ها

در این مرحله، فرمت داده‌ها به‌طور یکنواخت استاندارد شد تا وارد پایگاه داده PostgreSQL شوند. مهم‌ترین تغییرات شامل:

تبدیل تاریخ‌ها: مقادیر تاریخ‌های مربوط به تاریخ انتشار و آخرین به‌روزرسانی به فرمت استاندارد YYYY-MM-DD تبدیل شدند.

تبدیل داده‌های عددی: مقادیر فیلدهایی مانند تعداد نصب‌ها و قیمت که به صورت متنی ذخیره شده بودند، به مقادیر عددی تبدیل شدند.

حجم اپلیکیشن‌ها: مقادیر موجود در فیلد حجم که واحدهای مختلفی داشتند، به مگابایت تبدیل شدند.

تبدیل مقادیر گم‌شده: مقادیر گم‌شده در فیلدهای خاص مانند امتیاز اپلیکیشن‌ها، تعداد نصب‌ها، حجم و تاریخ‌ها پر شدند.

۵. اصلاح داده‌های خاص

در این مرحله، برخی از داده‌های خاص که شامل مشکلات فرمت بودند، اصلاح شدند. به عنوان مثال، نام اپلیکیشن‌ها و شناسه‌های توسعه‌دهندگان که حاوی کوتیشن‌های اضافی بودند، اصلاح شدند.

۶. ذخیره‌سازی داده‌های تمیز شده

پس از انجام عملیات پاکسازی و استانداردسازی داده‌ها، اطلاعات تمیز شده در قالب فایل‌های CSV ذخیره شدند. داده‌ها به سه فایل جداگانه تقسیم شدند:

فایل دسته‌بندی‌ها: شامل اطلاعات مربوط به تمامی دسته‌بندی‌ها.

فایل توسعه‌دهندگان: شامل اطلاعات مربوط به توسعه‌دهندگان اپلیکیشن‌ها.

فایل اپلیکیشن‌ها: شامل تمامی ویژگی‌های اپلیکیشن‌ها همراه با شناسه‌های دسته‌بندی‌ها و توسعه‌دهندگان.

بخش دوم: ایجاد پایگاه داده

در این بخش، مراحل طراحی پایگاه داده برای ذخیره سازی داده ها توضیح داده می شود.

۱. طراحی اسکیمای

جدول apps که شامل اطلاعات اپلیکیشن هاست و به جداول دسته بندی ها و توسعه دهندگان مرتبط می شود.

جدول categories که برای ذخیره اطلاعات دسته بندی ها است.

جدول developers که اطلاعات توسعه دهندگان را نگهداری می کند.

۲. وارد کردن داده

بعد از طراحی اسکیمای و ایجاد جداول، داده های تمیز شده که در قالب فایل های CSV ذخیره شده بودند، وارد پایگاه داده شدند. داده ها شامل اطلاعات اپلیکیشن ها، دسته بندی ها و توسعه دهندگان بودند.

بخش سوم: توسعه API های سمت بک اند

۱. فیلترها و جستجو

GET /filters	استخراج مقادیر پیشفرض فیلترها از پایگاه داده
GET /apps	جستجو و صفحه بندی اپلیکیشن ها بر اساس فیلترها

۲. آمار و تحلیل های داده

GET /statistics/rating_distribution	محاسبه توزیع امتیازات اپلیکیشن ها
GET /statistics/release_trend	محاسبه روند انتشار اپلیکیشن ها
GET /statistics/update_trend	محاسبه روند به روزرسانی اپلیکیشن ها
GET /statistics/average_rating	محاسبه میانگین امتیاز اپلیکیشن ها

۳. مدیریت دسته بندی ها

GET /categories	دریافت لیست تمامی دسته بندی ها
GET /categories/{category_id}	دریافت اطلاعات یک دسته بندی خاص
POST /categories	ایجاد دسته بندی جدید
PUT /categories/{category_id}	ویرایش دسته بندی موجود
DELETE /categories/{category_id}	حذف دسته بندی خاص

۴. مدیریت توسعه دهندگان

GET /developers	دریافت لیست توسعه دهندگان با صفحه بندی
GET /developers/{developer_id}	دریافت اطلاعات یک توسعه دهنده خاص
POST /developers	ایجاد توسعه دهنده جدید
PUT /developers/{developer_id}	ویرایش توسعه دهنده موجود
DELETE /developers/{developer_id}	حذف توسعه دهنده خاص

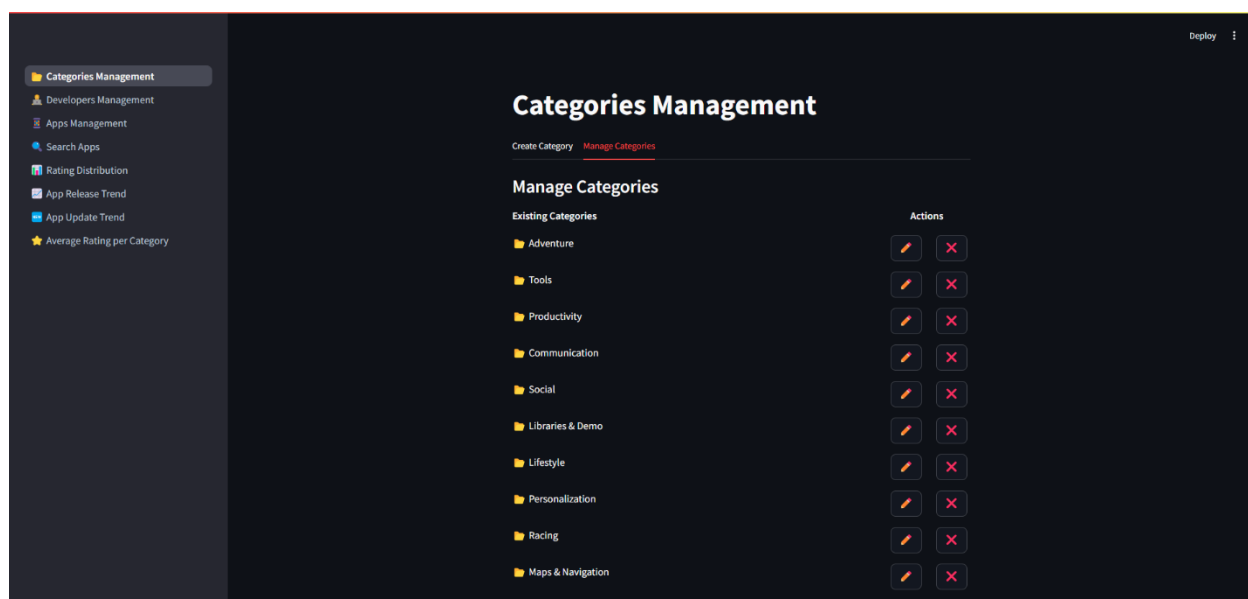
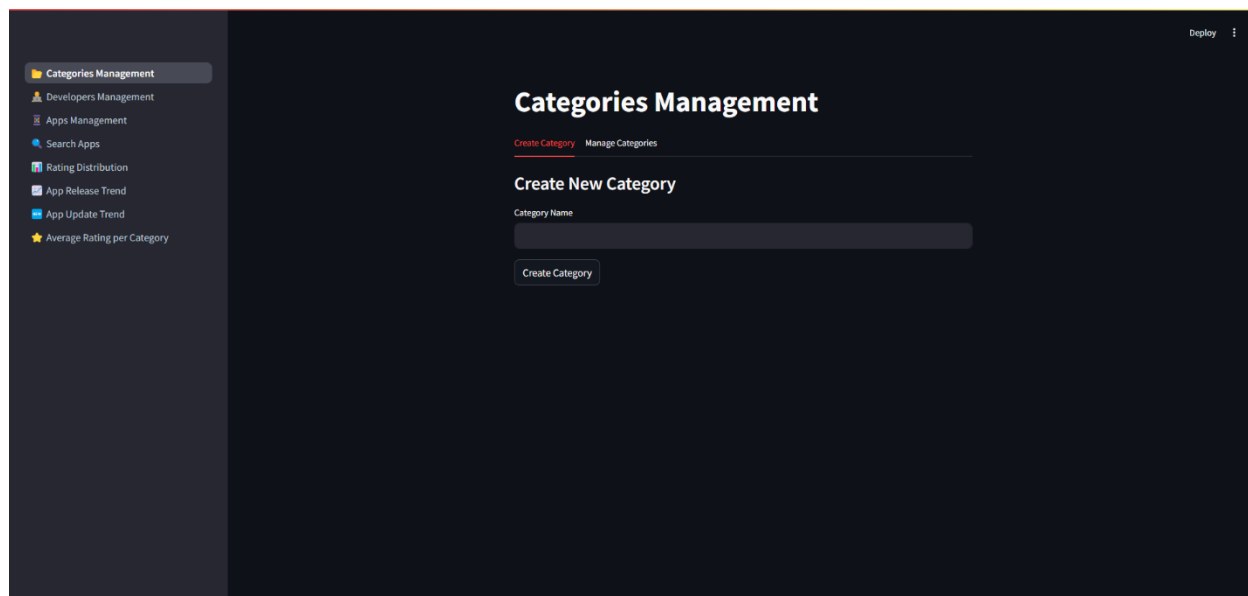
۵. مدیریت اپلیکیشن ها

GET /apps/{app_id}	دریافت اطلاعات یک اپلیکیشن خاص
POST /apps	ایجاد اپلیکیشن جدید
PUT /apps/{app_id}	ویرایش اپلیکیشن موجود
DELETE /apps/{app_id}	حذف اپلیکیشن خاص

بخش چهارم: توسعه داشبورد با Streamlit

داشبورد Streamlit شامل چندین صفحه است که به کاربران امکان مدیریت داده‌ها، جستجوی اپلیکیشن‌ها و مشاهده تحلیل‌های مختلف را می‌دهد. در ادامه، صفحات مختلف و قابلیت‌های هر یک توضیح داده شده است.

۱. مدیریت دسته‌بندی‌ها



قابلیت‌ها:

- نمایش لیست دسته‌بندی‌ها
- افزودن دسته‌بندی جدید
- ویرایش نام دسته‌بندی
- حذف دسته‌بندی

۲. مدیریت توسعه‌دهندگان

The screenshot shows the 'Developers Management' section with the 'Create New Developer' form. The sidebar on the left contains links to 'Categories Management', 'Developers Management' (active), 'Apps Management', 'Search Apps', 'Rating Distribution', 'App Release Trend', 'App Update Trend', and 'Average Rating per Category'. The main content area has a 'Deploy' button in the top right. Below the 'Developers Management' header, there are links for 'Create Developer' (active) and 'Manage Developers'. The form includes fields for 'Developer Name' and 'Developer Email', and a 'Create Developer' button.

The screenshot shows the 'Manage Developers' table. The sidebar and header are identical to the previous screenshot. The table lists existing developers with their names and email addresses, and provides actions for each (edit and delete).

Existing Developers	Actions
MrScratch (mithalaarush@gmail.com)	
Eqra Tech (ehiyassat@eqratech.com)	
Jie App (jie.myapp@gmail.com)	
has.to.be gmbh (support@has-to-be.com)	
WayCall (waycallcol@gmail.com)	
pietechsolution (info@ionicfirebaseapp.com)	
Hotels Attitude (info@hotels-attitude.com)	
TK Applications (tkapplications001@gmail.com)	
Grit Technology (info@redmondcompany)	
Human Droid Apps (arukabdlillah4@gmail.com)	

فیلترها:

- صفحه‌بندی (شماره صفحه، تعداد نتایج در هر صفحه)

قابلیت‌ها:

- نمایش لیست توسعه‌دهندگان
- افزودن توسعه‌دهنده جدید
- ویرایش نام و ایمیل توسعه‌دهنده
- حذف توسعه‌دهنده

۳. مدیریت اپلیکیشن‌ها

The screenshot shows the 'Create New App' form in the 'Apps Management' dashboard. The left sidebar contains navigation links: Categories Management, Developers Management, Apps Management (selected), Search Apps, Rating Distribution, App Release Trend, App Update Trend, and Average Rating per Category. The main content area has a 'Deploy' button in the top right. Below the 'Apps Management' header, there are tabs for 'Create App' (active) and 'Manage Apps'. The form fields include: App Name (text input), App ID (text input), Category ID (dropdown menu with '1' selected), Developer ID (dropdown menu with '1' selected), Rating (text input with '0.00' and a range selector), a 'Free' checkbox, and a 'Create App' button.

The screenshot shows the 'Manage Apps' table in the 'Apps Management' dashboard. The left sidebar is identical to the previous screenshot. The main content area has a 'Deploy' button in the top right. Below the 'Apps Management' header, there are tabs for 'Create App' and 'Manage Apps' (active). The table lists existing apps with their names and IDs, and provides actions for each (edit and delete). The table has two columns: 'Existing Apps' and 'Actions'.

Existing Apps	Actions
Rainbow Cup Launcher Theme (rainbow.cup.themes)	
Bound Strike (com.jfigames.gbbs)	
Multi League: Soccer/ Football Live Scores Results (com.appaso.live.score.football)	
Pipa Coloring Book (com.nolimit.pipacoloring)	
St Patricks Raquets Club (app.activitypro.stpatrick)	
TicTacToe (com.chincotectictactoe.games)	
Wedding Salon - Bride Princess (alr.LilDressUpGames.BridePrincessWeddingSalon)	
Piano Clicker (appinventor.ai._adam_zhakenov.dfgjkl)	
Push All My Buttons (com.PxelGameHouseLTD.PushAllMyButtons)	
افغان (afghanistan.afg.rg)	

فیلترها:

- صفحه‌بندی (شماره صفحه، تعداد نتایج در هر صفحه)

قابلیت‌ها:

- نمایش لیست اپلیکیشن‌ها
- افزودن اپلیکیشن جدید
- ویرایش مشخصات اپلیکیشن
- حذف اپلیکیشن

۴. جستجوی اپلیکیشن‌ها

The screenshot shows a 'Search Apps' interface. On the left is a sidebar with various filters: 'Select Category' (set to 'All'), 'Rating Range' (0.00 to 5.00), 'Price Range (\$)' (0.00 to 400.00), 'Minimum Installs' (0), 'Maximum Installs' (1,000,000,000), 'Content Rating' (set to 'All'), and checkboxes for 'Show Only Free Apps', 'Apps with Ads', 'Apps with In-App Purchases', and 'Editors' Choice'. The main area is titled 'Search Apps' and displays a table of results. The table has columns for 'id', 'app_id', 'app_name', and 'category_id'. It shows 10 results per page, with the first page displaying 10 items. Below the table, it says 'Showing page 1 of 231258 (Total results: 2312573)'.

id	app_id	app_name	category_id
0	3,701	rainbow.cup.themes	Rainbow Cup Launcher Theme
1	3,702	com.jlgames.gbts	Bound Strike
2	3,703	com.appaso.live.score.football	Multi League: Soccer/ Football Li
3	3,704	com.nolimit.pipacoloring	Pipe Coloring Book
4	3,705	app.activitypro.stpatrick	St Patricks Raquets Club
5	3,706	com.chincotec.tictactoe.games	TicTacToe
6	3,707	ali.LiliDressUpGames.BridePrincessWeddingSalon	Wedding Salon - Bride Princess
7	3,708	appinventor.ai._adam_zhakenov.dfgghjk	Piano Clicker
8	3,709	com.PeelGameHouseLTD.PushAllMyButtons	Push All My Buttons
9	3,710	afghanistan.afg.rg	افغانستان

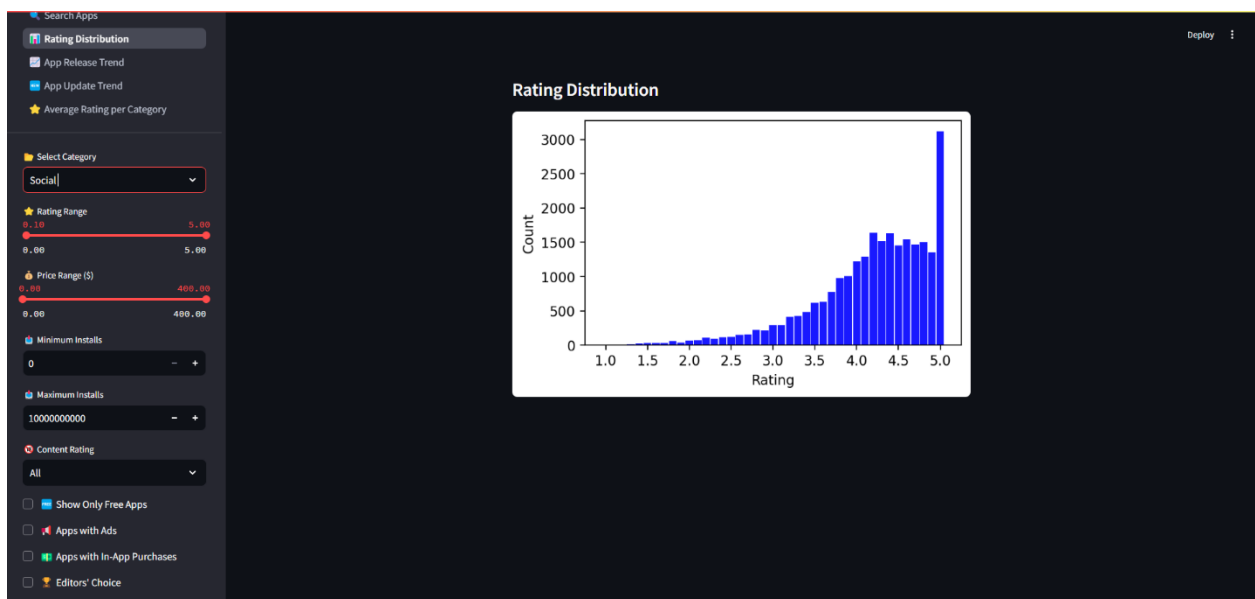
فیلترها:

- دسته‌بندی
- محدوده امتیاز
- محدوده قیمت
- محدوده تعداد نصب‌ها
- رده‌بندی سنی
- وضعیت رایگان بودن، پشتیبانی از تبلیغات، خریدهای درون‌برنامه‌ای، انتخاب سردبیر
- صفحه‌بندی (شماره صفحه، تعداد نتایج در هر صفحه)

قابلیت‌ها:

- جستجو و نمایش اپلیکیشن‌ها بر اساس فیلترهای انتخابی
- نمایش تعداد کل نتایج و تعداد صفحات

۵. توزیع امتیازات



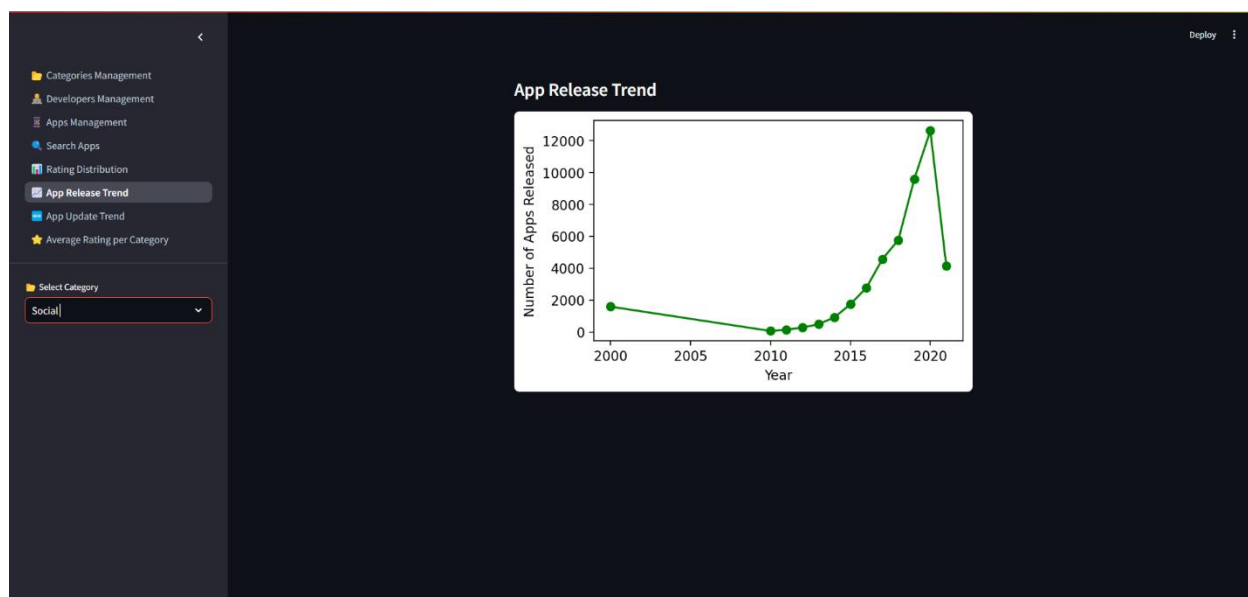
فیلترها:

- دسته‌بندی
- محدوده امتیاز
- محدوده قیمت
- محدوده تعداد نصب‌ها
- رده‌بندی سنی
- وضعیت رایگان بودن، پشتیبانی از تبلیغات، خریدهای درون‌برنامه‌ای، انتخاب سردبیر

نمایش:

- نمودار میله‌ای توزیع تعداد اپلیکیشن‌ها بر اساس امتیاز

۶. روند انتشار اپلیکیشن‌ها



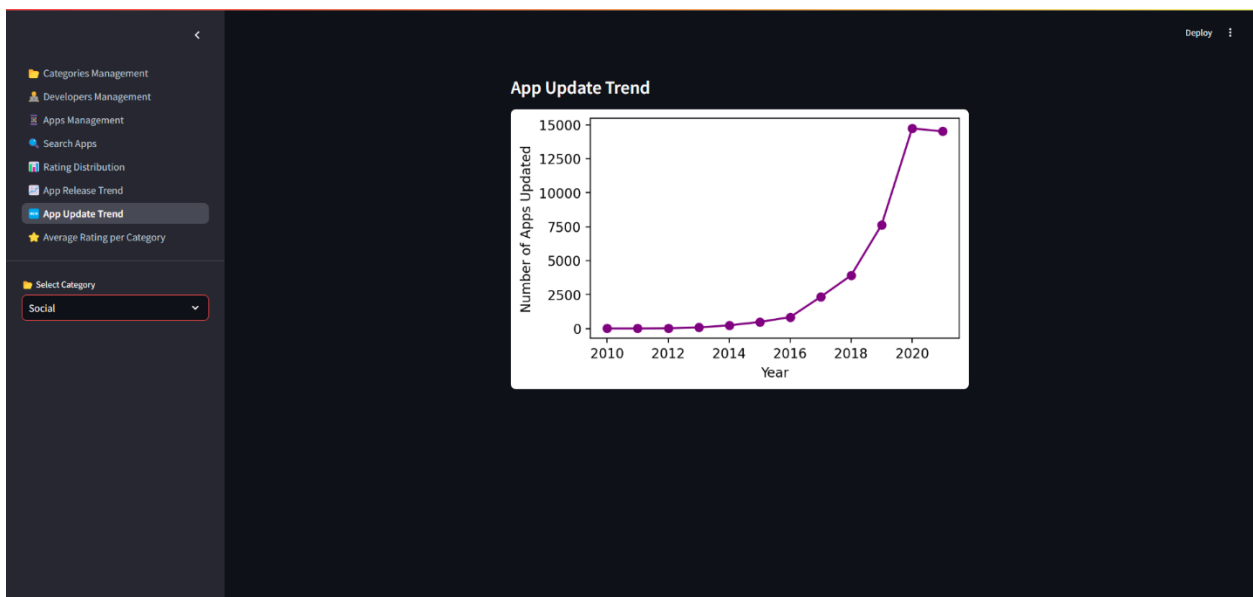
فیلترها:

• دسته‌بندی

نمایش:

• نمودار خطی تعداد اپلیکیشن‌های منتشر شده در هر سال

۷. روند به‌روزرسانی اپلیکیشن‌ها



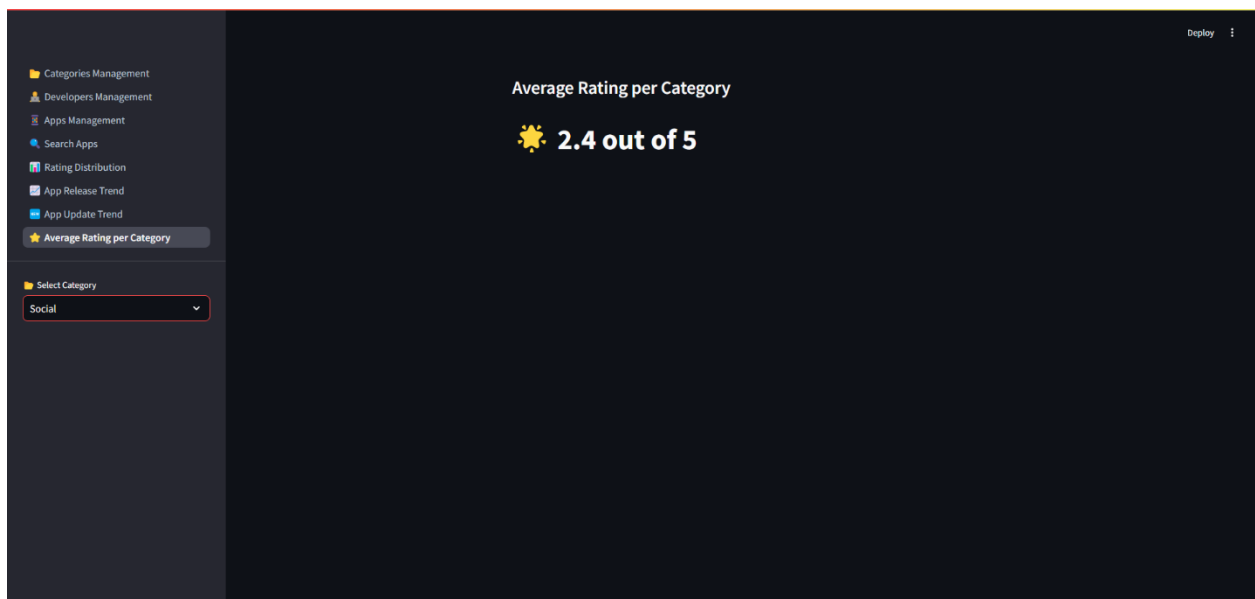
فیلترها:

• دسته‌بندی

نمایش:

• نمودار خطی تعداد اپلیکیشن‌هایی که در هر سال به‌روزرسانی شده‌اند

۸. میانگین امتیاز اپلیکیشن‌ها



فیلترها:

- دسته‌بندی

نمایش:

- مقدار میانگین امتیاز اپلیکیشن‌های یک دسته‌بندی خاص یا همه دسته‌بندی‌ها

بخش پنجم: بهینه‌سازی پایگاه داده

در این بخش، ایندکس‌های ایجادشده در پایگاه داده PostgreSQL و تاثیر آن‌ها بر عملکرد جستجوها بررسی می‌شود. هر ایندکس برای بهبود سرعت فیلترها، مرتب‌سازی‌ها و نمودارهای تحلیلی در داشبورد Streamlit طراحی شده است.

۱. ایندکس‌های تک‌ستونی

idx_apps_category_id (category_id)

- بهینه‌سازی جستجو و فیلتر اپلیکیشن‌ها بر اساس دسته‌بندی
- بهبود فیلتر دسته‌بندی در جستجوی اپلیکیشن‌ها، توزیع امتیازات، روند انتشار، روند به‌روزرسانی، میانگین امتیاز

idx_apps_rating (rating)

- افزایش سرعت مرتب‌سازی و فیلتر اپلیکیشن‌ها بر اساس امتیاز
- بهبود فیلتر امتیاز در جستجوی اپلیکیشن‌ها، توزیع امتیازات، میانگین امتیاز

idx_apps_price (price)

- بهبود جستجو و فیلتر اپلیکیشن‌ها بر اساس محدوده قیمت
- بهبود فیلتر قیمت در جستجوی اپلیکیشن‌ها

idx_apps_installs (installs)

- افزایش سرعت مرتب‌سازی و جستجوی اپلیکیشن‌ها بر اساس تعداد نصب‌ها.
- بهبود فیلتر نصب‌ها در جستجوی اپلیکیشن‌ها

idx_apps_content_rating (content_rating)

- بهینه‌سازی جستجو بر اساس امتیاز محتوا (رده‌بندی سنی)
- بهبود فیلتر امتیاز محتوا در جستجوی اپلیکیشن‌ها

idx_apps_free (free)

- افزایش سرعت فیلتر اپلیکیشن‌های رایگان یا پولی
- بهبود فیلتر اپلیکیشن‌های رایگان در جستجوی اپلیکیشن‌ها

idx_apps_ad_supported (ad_supported)

- بهینه‌سازی فیلتر اپلیکیشن‌هایی که تبلیغات دارند
- بهبود فیلتر تبلیغات در جستجوی اپلیکیشن‌ها

idx_apps_in_app_purchases (in_app_purchases)

- افزایش سرعت فیلتر اپلیکیشن‌های دارای خرید درون برنامه‌ای
- بهبود فیلتر خرید درون برنامه‌ای در جستجوی اپلیکیشن‌ها

idx_apps_editors_choice (editors_choice)

- بهینه‌سازی جستجو برای نمایش اپلیکیشن‌های منتخب سردبیر
- بهبود فیلتر انتخاب سردبیر در جستجوی اپلیکیشن‌ها

۲. ایندکس‌های چندستونی

idx_apps_category_released (category_id, released)

- افزایش سرعت جستجو بر اساس ترکیب دسته‌بندی و تاریخ انتشار
- بهبود روند انتشار اپلیکیشن‌ها

idx_apps_category_last_updated (category_id, last_updated)

- بهبود جستجو برای ترکیب دسته‌بندی و تاریخ آخرین به‌روزرسانی
- بهبود روند به‌روزرسانی اپلیکیشن‌ها

idx_apps_category_rating (category_id, rating)

- بهینه‌سازی مرتب‌سازی و فیلتر اپلیکیشن‌ها بر اساس دسته‌بندی و امتیاز
- بهبود جستجوی اپلیکیشن‌ها، میانگین امتیاز

۳. نتایج ایندکس کردن داده‌ها

در ویدیوی دمو، لاگ زمان اجرای کوئری‌های مختلف روی داده‌های ایندکس شده و همچنین بعد از حذف ایندکس‌ها بررسی شده‌است. با استفاده از ایندکس روی ستون‌ها، میانگین زمان اجرای کوئری‌ها از ۲ تا ۱۰ ثانیه به حدود ۰.۰۰۱ تا ۰.۱ ثانیه رسیده است.

نمونه‌ها:

۱. کوئری جستجوی اپلیکیشن‌ها در دسته‌بندی Social با رده‌بندی سنی Teen

- قبل از اعمال ایندکس

```
DEBUG:database:Executing SQL Query: SELECT count(*) AS count_1
FROM (SELECT apps.id AS apps_id, apps.app_id AS apps_app_id, apps.app_name AS
apps_app_name, apps.category_id AS apps_category_id, apps.developer_id AS
apps_developer_id, apps.rating AS apps_rating, apps.rating_count AS
apps_rating_count, apps.installs AS apps_installs, apps.min_installs AS
apps_min_installs, apps.max_installs AS apps_max_installs, apps.free AS
apps_free, apps.price AS apps_price, apps.currency AS apps_currency, apps.size
AS apps_size, apps.min_android AS apps_min_android, apps.released AS
apps_released, apps.last_updated AS apps_last_updated, apps.content_rating AS
apps_content_rating, apps.ad_supported AS apps_ad_supported,
apps.in_app_purchases AS apps_in_app_purchases, apps.editors_choice AS
apps_editors_choice, apps.scraped_time AS apps_scraped_time
FROM apps
WHERE apps.category_id = 17 AND apps.content_rating = Teen) AS anon_1
DEBUG:database:Query executed in 5.5680 seconds
INFO: 127.0.0.1:59142 - "GET
/apps?category=Arcade&content_rating=Teen&free=False&ad_supported=False&in_app
_purchases=False&editors_choice=False&page=1&per_page=10 HTTP/1.1" 200 OK
```

- بعد از اعمال ایندکس

```
DEBUG:database:Executing SQL Query: SELECT count(*) AS count_1
FROM (SELECT apps.id AS apps_id, apps.app_id AS apps_app_id, apps.app_name AS
apps_app_name, apps.category_id AS apps_category_id, apps.developer_id AS
apps_developer_id, apps.rating AS apps_rating, apps.rating_count AS
apps_rating_count, apps.installs AS apps_installs, apps.min_installs AS
apps_min_installs, apps.max_installs AS apps_max_installs, apps.free AS
apps_free, apps.price AS apps_price, apps.currency AS apps_currency, apps.size
```



```
AS apps_size, apps.min_android AS apps_min_android, apps.released AS
apps_released, apps.last_updated AS apps_last_updated, apps.content_rating AS
apps_content_rating, apps.ad_supported AS apps_ad_supported,
apps.in_app_purchases AS apps_in_app_purchases, apps.editors_choice AS
apps_editors_choice, apps.scraped_time AS apps_scraped_time
FROM apps
WHERE apps.category_id = 17 AND apps.content_rating = Teen) AS anon_1
DEBUG:database:Query executed in 0.0030 seconds
INFO: 127.0.0.1:59142 - "GET
/apps?category=Arcade&content_rating=Teen&free=False&ad_supported=False&in_app
_purchases=False&editors_choice=False&page=1&per_page=10 HTTP/1.1" 200 OK
```

زمان اجرا با بیش از ۱۸۰۰ درصد بهبود، از ۵.۵۶۸۰ ثانیه به ۰.۰۰۳۰ ثانیه رسیده است.

۲. کوئری میانگین امتیاز دسته بندی

- قبل از اعمال ایندکس

```
DEBUG:database:Executing SQL Query: SELECT avg(apps.rating) AS avg_1
FROM apps
WHERE apps.category_id = 5
DEBUG:database:Query executed in 6.5326 seconds
INFO: 127.0.0.1:64691 - "GET
/statistics/average_rating/?category_name=Social HTTP/1.1" 200 OK
```

- بعد از اعمال ایندکس

```
DEBUG:database:Executing SQL Query: SELECT avg(apps.rating) AS avg_1
FROM apps
WHERE apps.category_id = 5
DEBUG:database:Query executed in 0.0040 seconds
INFO: 127.0.0.1:64691 - "GET
/statistics/average_rating/?category_name=Social HTTP/1.1" 200 OK
```

زمان اجرا با بیش از ۱۶۰۰ درصد بهبود، از ۶.۵۳۲۶ ثانیه به ۰.۰۰۴۰ ثانیه رسیده است.

بخش ششم: اطلاعات تکمیلی

۱. لینک ویدیوی دمو:

<https://drive.google.com/file/d/1r8GfgGVyj56qlhsO5HK8IE80V5kjl4d/view?usp=sharing>

۲. لینک ریپوی گیتهاب:

<https://github.com/AliSK81/google-playstore-analysis>