



دانشکده مهندسی کامپیوتر

گزارش فاز اول

باوان دیوانی آذر ۹۸۵۲۲۱۱

نیم سال دوم

سال تحصیلی ۱۴۰۱-۰۲

اطلاعات جمع آوری شده از سایت زیر است:

<https://store.steampowered.com/>

با استفاده از کتابخانه BeautifulSoup اطلاعات مورد نیاز از سایت crawl شده است. به دلیل تعداد بالای رکوست‌ها و پرشدن حافظه‌ی RAM هر ۵ هزار داده در یک فایل ذخیره شده‌اند. اما در نهایت، اطلاعات مربوط به بازی‌های یک گروه، در یک فایل قرار داده شده‌اند.

روش کار به این صورت است که در ابتدا، لیست app_id مربوط به بازی‌های یک گروه را استخراج کرده، به کمک تابع add_information به تک تک app_idها، رکوست می‌فرستیم تا اطلاعات (نام و توضیحات) مربوط به بازی، به دست آید.
کد:

```
categories = ['action', 'adventure', 'rpg', 'strategy', 'simulation', 'sports_and_racing']
data = []
for category in categories:
    print(category)
    start = 5000
    while True:
        response = requests.get(f"https://store.steampowered.com/saleaction/ajaxgetaledyna")
        app_ids = response.json()["appids"]
        response.close()
        start += 100
        print(start)
        if len(app_ids) == 0:
            break
        add_information(app_ids, category, data)
        if start % 5000 == 0:
            df = pd.DataFrame(data, columns=['id', 'name', 'category', 'about'])
            filename = f'gdrive/MyDrive/NLP/Dataset/{category}{start//5000}.csv'
            df.to_csv(filename, index=False, quoting=csv.QUOTE_ALL)
            data = []
        if start % 5000 != 0:
            df = pd.DataFrame(data, columns=['id', 'name', 'category', 'about'])
            filename = f'gdrive/MyDrive/NLP/Dataset/{category}{(start//5000)+1}.csv'
            df.to_csv(filename, index=False, quoting=csv.QUOTE_ALL)
            data = []
```

```
def add_information(app_ids, category, data):
    for i, app_id in enumerate(app_ids):
        response = requests.get(f"https://store.steampowered.com/app/{app_id}")
        soup = BeautifulSoup(response.content, 'lxml')
        response.close()
        try:
            name = soup.title.string
            info = soup.find("div", {"id": "aboutThisGame"}).get_text("\n")
            info = re.sub('\r\n|\r|\n+', '\n', info).replace("\t", "")
            data.append([app_id, name, category, info])
        except:
            print(f"An exception occurred app_id = {app_id}")
            print(f"https://store.steampowered.com/app/{app_id}")
```

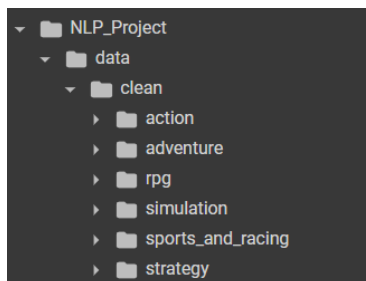
```
categories = ['action', 'adventure', 'rpg', 'strategy', 'simulation', 'sports_and_racing']
for category in categories:
    PATH = 'gdrive/MyDrive/NLP/Dataset/'
    df1 = pd.read_csv(f'{PATH}{category}1.csv')
    df2 = pd.read_csv(f'{PATH}{category}2.csv')
    df3 = pd.read_csv(f'{PATH}{category}3.csv')

    dict1 = df1.to_dict('dict')
    dict2 = df2.to_dict('dict')
    dict3 = df3.to_dict('dict')

    combined_dict = []
    for i in range(len(dict1['id'])):
        combined_dict.append([dict1['id'][i], dict1['name'][i], dict1['category'][i], dict1['about'][i]])
    for i in range(len(dict2['id'])):
        combined_dict.append([dict2['id'][i], dict2['name'][i], dict2['category'][i], dict2['about'][i]])
    for i in range(len(dict3['id'])):
        combined_dict.append([dict3['id'][i], dict3['name'][i], dict3['category'][i], dict3['about'][i]])

    directory = f'gdrive/MyDrive/NLP/github/data/raw/{category}/'
    df = pd.DataFrame(combined_dict, columns=['id', 'name', 'category', 'about'])
    if not os.path.exists(directory):
        os.makedirs(directory)
    df.to_csv(f'{directory}{category}.csv', index=False, quoting=csv.QUOTE_ALL)
```

لیست اطلاعات بازی های مربوط به ۶ گروه درآورده شده و اطلاعات هر گروه در فولدر مربوطه قرار گرفته شده است. فرمت داده ها به شکل فایل CSV هست.



برای تفکیک جملات و کلمات، از متدهای موجود در کتابخانه nltk استفاده کردم. این متدها شامل توابعی مانند `sent_tokenize()` برای تفکیک جملات و `word_tokenize()` برای تفکیک کلمات هستند.

```
for i in range(len(data['clean_description'])):
    id = data['id'][i]
    description = data['clean_description'][i]
    sentences = nltk.sent_tokenize(description)
    words = word_tokenize(description)
    sentences_data.append([id, sentences])
    words_data.append([id, words])
```

۳

برای بخش clean کردن داده، ابتدا عبارت "About the game" را از ابتدای متن حذف کرده و سپس بررسی کرده‌ایم که آیا توضیحات بازی کمتر از ده حرف هستند؟ در صورتی که چنین باشد، آن‌ها را نیز حذف می‌کنیم.

```
for i in range(len(data['about'])):
    id = data['id'][i]
    description = data['about'][i]
    clean_description = description.replace('\nAbout This Game\n', '')
    data_count+=1
    if(len(clean_description)>10):
        clean_data_count +=1
        clean_data.append([id, clean_description])
```

تعداد داده قبل و بعد از clean کردن :

```
before cleaning data : 78719
after cleaning data: 78665
```

واحد برچسب گذاری برای هر متن بازی ، ژانر خود بازی هست. برای مثال اگر ژانر بازی اکشن باشند برچسب توضیحات آن ، اکشن می‌شود.

آمار داده به تفکیک برچسب در قالب جدول «و» نمودار
أ. تعداد «واحد» داده

"action"	"adventure"	"rpg"	"strategy"	"simulation"	"sports_and_racing"
"14983"	"14982"	"14987"	"11513"	"14994"	"7206"

ب. تعداد جملات

"action"	"adventure"	"rpg"	"strategy"	"simulation"	"sports_and_racing"
"167376"	"170770"	"168316"	"155396"	"172196"	"76045"

ج. تعداد کلمات

"action"	"adventure"	"rpg"	"strategy"	"simulation"	"sports_and_racing"
"3287666"	"3306827"	"3458951"	"3236588"	"3563133"	"1497667"

د. تعداد کلمات منحصر به فرد

"category"	"unique words"	"common words"	"uncommon"
"action"	"103295"	"85947"	"17348"
"adventure"	"106669"	"90025"	"16644"
"rpg"	"112026"	"95099"	"16927"
"strategy"	"97843"	"77818"	"20025"
"simulation"	"115101"	"86454"	"28647"
"sports_and_racing"	"62055"	"52599"	"9456"

ه. تعداد کلمات منحصر به فرد مشترک و غیر مشترک بین پرچسب ها

"category"	"unique words"	"common words"	"uncommon"
"action"	"103295"	"85947"	"17348"
"adventure"	"106669"	"90025"	"16644"
"rpg"	"112026"	"95099"	"16927"
"strategy"	"97843"	"77818"	"20025"
"simulation"	"115101"	"86454"	"28647"
"sports_and_racing"	"62055"	"52599"	"9456"

و. ۱۰ کلمه پرتکرار غیر مشترک هر برجسب

"category"	"word1"	"word2"	"word3"	"word4"	"word5"	"word6"	"word7"	"word8"	"word9"	"word10"
"action"	"Uchiban"	"Touken"	"CA2S"	"Ranbu"	"Honmaru"	"Warzone "	"honmaru"	"Jutsu"	"Paintjob"	"GameGuru"
"adventure"	"mini-games/puzzles"	"hidden-object"	"HDPs"	"Deponia"	"•Never"	"Log1"	"Mimpi"	"ScumVM"	"Putt-Putts"	"Pinecreek"
"rpg"	"Hajimari"	"Battlers"	"Lumena"	"x2000"	"S.F.A"	"Zelam"	"Sevin"	"WML2"	" Music"	" Special"
"strategy"	"3.3.5"	"3.3.4"	"Catan"	"ENnie"	"eXtra"	"non-mythic"	"Freecell"	"ChessBase"	"2.9.9"	"Senet"
"simulation"	"Trainz"	"CSX"	"Subdivision"	"TRS19"	"Quinnimont"	"QJ"	"doubledecker"	"couplers"	"Railfan"	"Bahn"
"sports_and_racing"	"WHEELS "	"Hellfest"	"•Maximum"	"Riskers"	"Ragnarock"	"Betelgeuze"	"TrackMania"	"Ablman"	"Gloryhammer"	" "

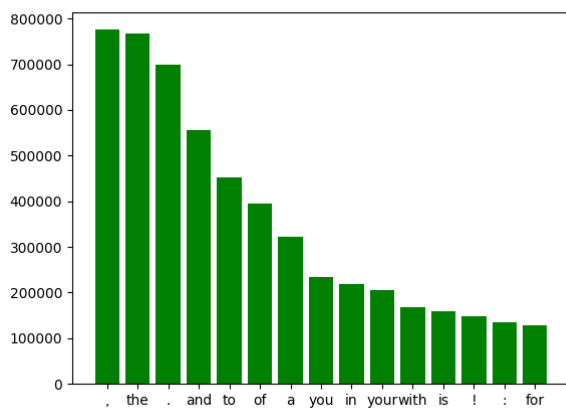
ز. ۱۰ کلمه مشترک برتر هر برجسب نسبت به برجسبهای دیگر بر اساس معیار زیر

"category"	"word1"	"word2"	"word3"	"word4"	"word5"	"word6"	"word7"	"word8"	"word9"	"word10"
"action"	"resized"	"locomotives"	"parcels"	"pre-placed"	"Grounds"	"freight"	"Pathfinder"	"Converted"	"Adapted"	"Railroad"
"adventure"	"Starfinder"	"resized"	"pre-placed"	"handouts"	"km/h"	"parcels"	"mm"	"BR"	"DB"	"Arma"
"rpg"	"liveries"	"locomotives"	"Airport"	"Arma"	"RC"	"brake"	"GT"	"Livery"	"cab"	"BR"
"strategy"	"x10"	"mm"	"RC"	"brake"	"Stella"	"PACK"	"XR"	"Tyre"	"Dutys"	"off-road"
"simulation"	"Grounds"	"Pathfinder"	"pre-placed"	"handouts"	"ruleset"	"JRPG"	"Adapted"	"resized"	"Conversion"	"Terms"
"sports_and_racing"	"ruleset"	"Grounds"	"Pathfinder"	"Adapted"	"freight"	"Savage"	"subscription"	"Terms"	"erotic"	"Male"

ح. ده کلمه برتر بر اساس $TF_IDF(W)$

"category"	"word1"	"word2"	"word3"	"word4"	"word5"	"word6"	"word7"	"word8"	"word9"	"word10"
"action"	"resized"	"locomotives"	"parcels"	"pre-placed"	"Grounds"	"freight"	"Pathfinder"	"Converted"	"Adapted"	"Railroad"
"adventure"	"Starfinder"	"resized"	"pre-placed"	"handouts"	"km/h"	"parcels"	"mm"	"BR"	"DB"	"Arma"
"rpg"	"liveries"	"locomotives"	"Airport"	"Arma"	"RC"	"brake"	"GT"	"Livery"	"cab"	"BR"
"strategy"	"x10"	"mm"	"RC"	"brake"	"Stella"	"PACK"	"XR"	"Tyre"	"Dutys"	"off-road"
"simulation"	"Grounds"	"Pathfinder"	"pre-placed"	"handouts"	"ruleset"	"JRPG"	"Adapted"	"resized"	"Conversion"	"Terms"
"sports_and_racing"	"ruleset"	"Grounds"	"Pathfinder"	"Adapted"	"freight"	"Savage"	"subscription"	"Terms"	"erotic"	"Male"

ط. هیستوگرام تعداد تکرار هر کلمه منحصر به فرد به ترتیب از فرکانس بالا به پایین



لینک github :

https://github.com/bavanDA/NLP_Project

لینک huggingface :

https://huggingface.co/datasets/Bavanda/Steam_DG