

도서관 빅데이터를 활용한 20대 성별에 따른 대출 트렌드 분석

202044069 배예진

<https://github.com/daraeyaa/data-analysis>

목차

1. 데이터 수집 방법
2. 데이터 전처리 - 20대 성별에 따른 대출 트렌드 분석
3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석
4. 데이터 전처리 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석
5. 데이터 시각화 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석
6. 프로젝트 리뷰

1. 데이터 수집 방법

사용 데이터 : 도서관 정보나루 제공 데이터
<https://data4library.kr/>

수집 방법 : 도서관 정보나루 사이트의 '인기대출도서' 데이터를 수집한다.

조건 1) 기간(금년), 성별(여성, 남성), 연령(20대)

조건 2) 기간(작년, 분기별 설정), 성별(여성, 남성), 연령(20대)

⇒ 여성과 남성의 데이터를 따로 받아야 하므로
동시에 선택하지 않고 각각 선택하여 csv 파일을 다운 받는다.



참여 도서관 목록 장서/대출데이터 **인기대출도서** 도서별 이용분석 대출 급상승 도서 지역별 비교분석 이달의 키워드

검색 조건들을 선택한 후 데이터를 다운로드 하거나, API 이용신청을 할 수 있습니다.

데이터 제공방식	<input checked="" type="radio"/> 전체	<input type="radio"/> 월별	<input type="radio"/> 연도별			
대출 기간	<input checked="" type="radio"/> 금년	<input type="radio"/> 금월	<input type="radio"/> 금주			
	기간설정	2022-01-01 ~ 2022-12-10				
성별	<input type="radio"/> 전체	<input type="radio"/> 남성	<input checked="" type="radio"/> 여성	<input type="radio"/> 미상		
연령	<input checked="" type="radio"/> 전체	<input type="radio"/> 영유아(0~5세)	<input type="radio"/> 유아(6~7세)	<input type="radio"/> 초등(8~13세)	<input type="radio"/> 청소년(14~19세)	<input type="radio"/> 20대
	<input type="radio"/> 30대	<input type="radio"/> 40대	<input type="radio"/> 50대	<input type="radio"/> 60세 이상	<input type="radio"/> 미상	
지역 <input type="checkbox"/> 세부지역	<input checked="" type="radio"/> 전체	<input type="radio"/> 서울	<input type="radio"/> 부산	<input type="radio"/> 대구	<input type="radio"/> 인천	<input type="radio"/> 광주
	<input type="radio"/> 대전	<input type="radio"/> 울산	<input type="radio"/> 세종	<input type="radio"/> 경기	<input type="radio"/> 강원	<input type="radio"/> 충북
	<input type="radio"/> 충남	<input type="radio"/> 전북	<input type="radio"/> 전남	<input type="radio"/> 경북	<input type="radio"/> 경남	<input type="radio"/> 제주
ISBN 부가기호 <input type="checkbox"/> 학습참고서2(초등)	<input checked="" type="radio"/> 전체	<input type="radio"/> 고양	<input type="radio"/> 실용	<input type="radio"/> 여성	<input type="radio"/> 청소년	<input type="radio"/> 학습참고서1(중고)
	<input type="radio"/> 학습참고서2(초등)	<input type="radio"/> 아동	<input type="radio"/> 전문			
주제 <input type="checkbox"/> 세부주제	<input checked="" type="radio"/> 전체	<input type="radio"/> 종류	<input type="radio"/> 철학	<input type="radio"/> 종교	<input type="radio"/> 사회과학	<input type="radio"/> 자연과학
	<input type="radio"/> 기술과학	<input type="radio"/> 예술	<input type="radio"/> 언어	<input type="radio"/> 문학	<input type="radio"/> 역사	
결과건수	<input checked="" type="radio"/> 200건	<input type="radio"/> 500건	<input type="radio"/> 1000건			

<https://www.data4library.kr/loanDataL>

1. 데이터 수집 방법

[최종 사용 데이터]

- 1) 20대 성별에 따른 대출 트렌드 분석용 : 2022년_20대_남성_인기도서
2022년_20대_여성_인기도서
- 2) 20대 성별, 분기별 대출 트렌드 분석용 : 2021_n분기_남성, 2021_n분기_여성

[csv 파일 구성 내용]

순위, 서명, 저자, 출판사, 출판년도, 권, ISBN, ISBN부가기호, KDC, 대출건수 (파일당 200건)

	A	B	C	D	E	F	G	H	I	J	
1	순위	서명	저자	출판사	출판년도	권	ISBN	ISBN부가기호	KDC	대출건수	
2		1 지구 끝의	지은아: 김	Giant Boo	2021		9.79E+12	3810	813.7	5176	
3		2 달러구트	지은아: 이	팩토리나인	2020		9.79E+12	3810	813.7	5028	
4		3 시선으로	지은아: 정	문학동네	2020		9.79E+12	3810	813.7	4891	
5		4 우리가 빛	지은아: 김	허블	2019		9.79E+12	3810	813.7	4763	

- 2021_1분기_남성.csv
- 2021_1분기_여성.csv
- 2021_2분기_남성.csv
- 2021_2분기_여성.csv
- 2021_3분기_남성.csv
- 2021_3분기_여성.csv
- 2021_4분기_남성.csv
- 2021_4분기_여성.csv
- 2022년_20대_남성_인기도서.csv
- 2022년_20대_여성_인기도서.csv

2. 데이터 전처리

1) 한글 사용을 위한 폰트 적용

```
[ ] # 단계 1: 폰트 설치
import matplotlib.font_manager as fm

!apt-get -qq -y install fonts-nanum > /dev/null
fontpath = '/usr/share/fonts/truetype/nanum/NanumBarunGothic.ttf'
font = fm.FontProperties(fname=fontpath, size=9)
fm._rebuild()
```

```
[ ] # 단계 2: 런타임 재시작
import os
os.kill(os.getpid(), 9)
```

```
[1] # 단계 3: 한글 폰트 설정
import matplotlib.pyplot as plt
import matplotlib as mpl
import matplotlib.font_manager as fm

# 마이너스 표시 문제
mpl.rcParams['axes.unicode_minus'] = False

# 한글 폰트 설정
path = '/usr/share/fonts/truetype/nanum/NanumGothicBold.ttf'
font_name = fm.FontProperties(fname=path, size=18).get_name()
plt.rc('font', family=font_name)
fm._rebuild()
```

2) 라이브러리 불러오기

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

데이터 가공 : pandas, numpy
데이터 시각화 : seaborn, matplotlib

2. 데이터 전처리 - 20대 성별에 따른 대출 트렌드 분석

3) csv -> 데이터프레임 변환

```
# 파일 경로
file_path_man = '/content/2022년_20대_남성_인기도서.csv'
file_path_woman = '/content/2022년_20대_여성_인기도서.csv'

# 데이터프레임 변환(인코딩 에러 -> encoding 파라미터 활용)
df_man = pd.read_csv(file_path_man, encoding='cp949')
df_woman = pd.read_csv(file_path_woman, encoding='cp949')
```

4) 성별 열 추가

```
#성별 열 추가
df_man['성별'] = '남성'
print('2022년_20대_남성_인기도서')
print(df_man.head())

df_woman['성별'] = '여성'
print('2022년_20대_여성_인기도서')
print(df_woman.head())
```

남녀 csv를 각각 따로 다운 받아서 합칠 예정이므로,
구분을 위해 성별 열을 추가한다.

2. 데이터 전처리 - 20대 성별에 따른 대출 트렌드 분석

5) 데이터프레임 병합, 가공

```
# 데이터프레임 합치기(행 방향으로 합치기)
df = pd.concat([df_man, df_woman])
print(df.info()) #400row 확인

#합친 데이터 프레임 가공
#권,ISBN 부가기호 열 삭제
df.drop(['권', 'ISBN부가기호'], axis=1, inplace=True)

#KDC 결측치 처리(396 non-null, 4건의 NaN)
df['KDC'].fillna(method='bfill', inplace=True)
print(df.info())
```

- 남녀 데이터프레임을 행 방향으로 합친다.

각 데이터가 200건이므로 concat 이후에 info()로 400건의 데이터를 확인한다.

- 필요한 데이터만 사용하기 위해 '권', 'ISBN부가기호' 열을 삭제한다.
- KDC(한국십진분류법) 열에 결측값이 있어 결측값을 제거한다.

2. 데이터 전처리 - 20대 성별에 따른 대출 트렌드 분석

6) 데이터 전처리 전후 비교

#	Column	Non-Null Count	Dtype
0	순위	400 non-null	int64
1	서명	400 non-null	object
2	저자	400 non-null	object
3	출판사	400 non-null	object
4	출판년도	400 non-null	object
5	권	41 non-null	float64
6	ISBN	400 non-null	float64
7	ISBN부가기호	394 non-null	float64
8	KDC	396 non-null	float64
9	대출건수	400 non-null	int64
10	성별	400 non-null	object



#	Column	Non-Null Count	Dtype
0	순위	400 non-null	int64
1	서명	400 non-null	object
2	저자	400 non-null	object
3	출판사	400 non-null	object
4	출판년도	400 non-null	object
5	ISBN	400 non-null	float64
6	KDC	400 non-null	float64
7	대출건수	400 non-null	int64
8	성별	400 non-null	object

열이 삭제되고, 결측값이 제거된 것을 확인할 수 있다.

2. 데이터 전처리 - 20대 성별에 따른 대출 트렌드 분석

7) KDC(한국십진분류법) 기반 분류를 위해 함수를 적용한다.

```
def kdcFunc(x):  
    if (x > 0 and x < 100):  
        return '종류'  
    elif (x >= 100 and x < 200):  
        return '철학'  
    elif (x >= 200 and x < 300):  
        return '종교'  
    elif (x >= 300 and x < 400):  
        return '사회과학'  
    elif (x >= 400 and x < 500):  
        return '자연과학'  
    elif (x >= 500 and x < 600):  
        return '기술과학'  
    elif (x >= 600 and x < 700):  
        return '예술'  
    elif (x >= 700 and x < 800):  
        return '언어'  
    elif (x >= 800 and x < 900):  
        return '문학'  
    elif (x >= 900 and x < 1000):  
        return '역사'  
    else:  
        return '미상'
```

```
df['sorted_KDC'] = df.apply(lambda x:kdcFunc(x['KDC']), axis=1)
```

000 종류	100 철학	200 종교	300 사회과학	400 자연과학
010 도서학, 서지학	110 형이상학	210 비교종교	310 통계학	410 수 학
020 문헌정보학	120 인식론, 인과론, 인간학	220 불 교	320 경제학	420 물리학
030 백과사전	130 철학의 체계	230 기독교	330 사회학, 사회문제	430 화 학
040 강연집, 수필집, 연설문집	140 경 학	240 도 교	340 정치학	440 천문학
050 일반연속간행물	150 동양철학, 사상	250 천도교	350 행정학	450 지 학
060 일반학회, 단체, 협회, 기관	160 서양철학	260 신 도	360 법 학	460 광물학
070 신문, 언론, 저널리즘	170 논리학	270 힌두교, 브라만교	370 교육학	470 생명과학
080 일반전집, 총서	180 심리학	280 이슬람교(회교)	380 풍속, 예절, 민속학	480 식물학
090 향토자료	190 윤리학, 도덕철학	290 기타 제종교	390 국방, 군사학	490 동물학
500 기술과학	600 예술	700 언어	800 문학	900 역사
510 의 학	610 건축물	710 한국어	810 한국문학	910 아시아
520 농업, 농학	620 조각, 조형예술	720 중국어	820 중국문학	920 유럽
530 공학, 공업일반, 토목공학, 환경공학	630 공예, 장식미술	730 일본어, 기타아시아제어	830 일본문학, 기타아시아문학	930 아프리카
540 건축공학	640 서 예	740 영 어	840 영미문학	940 북아메리카
550 기계공학	650 회화, 도화	750 독일어	850 독일문학	950 남아메리카
560 전기공학, 전자공학	660 사진예술	760 프랑스어	860 프랑스문학	960 오세아니아
570 화학공학	670 음 악	770 스페인어, 포르투갈어	870 스페인, 포르투갈문학	970 양극지방
580 제조업	680 공연예술, 매체예술	780 이탈리아어	880 이탈리아문학	980 지 리
590 생활과학	690 오락, 스포츠	790 기타제어	890 기타제문학	990 전 기

csv에는 KDC가 숫자로 입력되어 있어
한글로 변환한 새로운 열을 추가하기 위해 함수를 적용한다.

2. 데이터 전처리 - 20대 성별에 따른 대출 트렌드 분석

8) 함수 적용 후 데이터 확인

```
# Column      Non-Null Count  Dtype
---  -
0 순위        400 non-null    int64
1 서명        400 non-null    object
2 저자        400 non-null    object
3 출판사      400 non-null    object
4 출판년도    400 non-null    object
5 ISBN        400 non-null    float64
6 KDC         400 non-null    float64
7 대출건수    400 non-null    int64
8 성별        400 non-null    object
9 sorted_KDC  400 non-null    object
dtypes: float64(2), int64(2), object(6)
```

```
문학      220
사회과학  103
철학       37
종류       13
자연과학   9
기술과학   9
역사        6
예술        3
Name: sorted_KDC, dtype: int64
```

df.info()를 통해 sorted_KDC가 추가된 것을 확인.

df['sorted_KDC'].value_counts()를 통해 고유값 확인

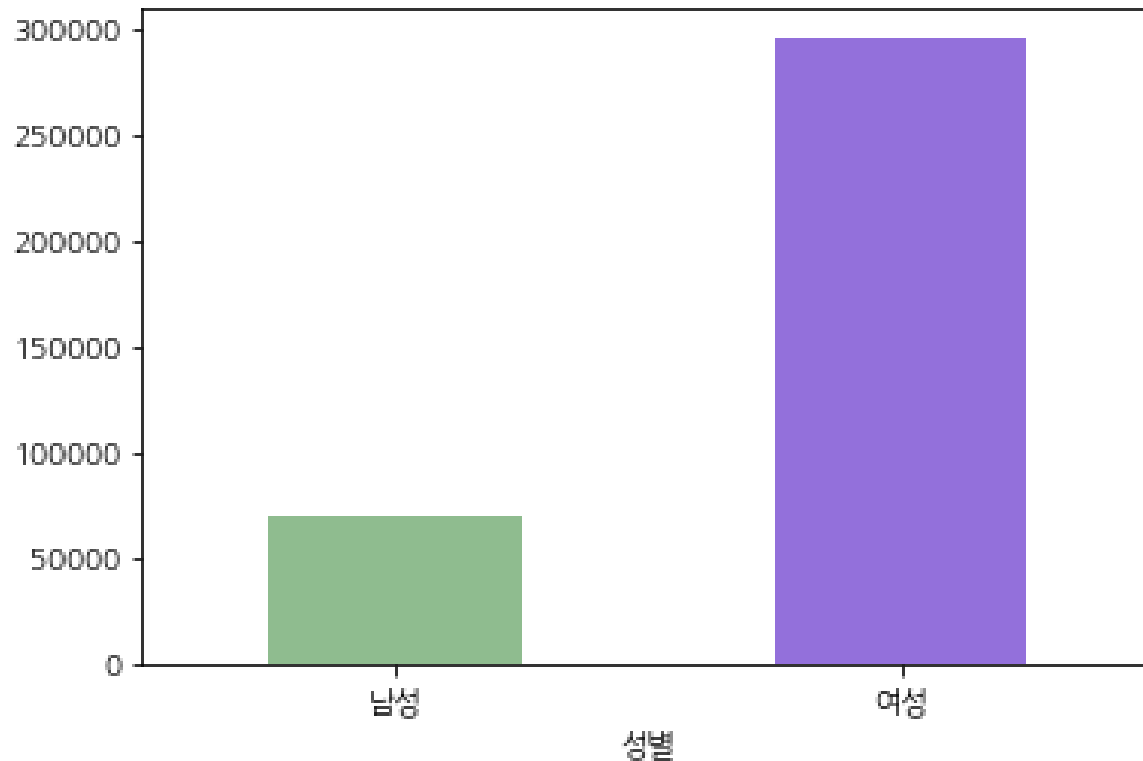
```
출판사  출판년도      ISBN      KDC  대출건수  성별  sorted_KDC
0      김영사   2015  9.788930e+12  909.00   973  남성      역사
1      미래엔   2020  9.791160e+12  332.60   879  남성      사회과학
2      현대문학  2012  9.788970e+12  833.60   842  남성      문학
3      팩토리나인 2020  9.791170e+12  813.70   783  남성      문학
4  Snowfox(스노우폭스북스) 2020  9.791190e+12  327.04   730  남성      사회과학
..      ...      ...      ...      ...      ...
195     해냄   2018  9.790000e+12  189.00   761  여성      철학
196     문학동네 2001  9.790000e+12  879.00   760  여성      문학
197     소미미디어 2020  9.790000e+12  833.60   759  여성      문학
198     은행나무 2014  9.790000e+12  813.70   757  여성      문학
199     문학동네 2017  9.790000e+12  813.62   756  여성      문학
```

[400 rows x 10 columns]

KDC 가 함수를 통해 sorted_KDC 열에
한글로 입력되어 있음을 확인.

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

1) 성별에 따른 대출건수 시각화



```
#성별별 대출건수 합산을 위한 groupby 연산
grouped_sex = df.groupby(df['성별'])['대출건수'].sum()

#성별별 대출건수 시각화
colors=['darkseagreen', 'mediumpurple']
grouped_sex.plot.bar(color=colors, rot=0)
```

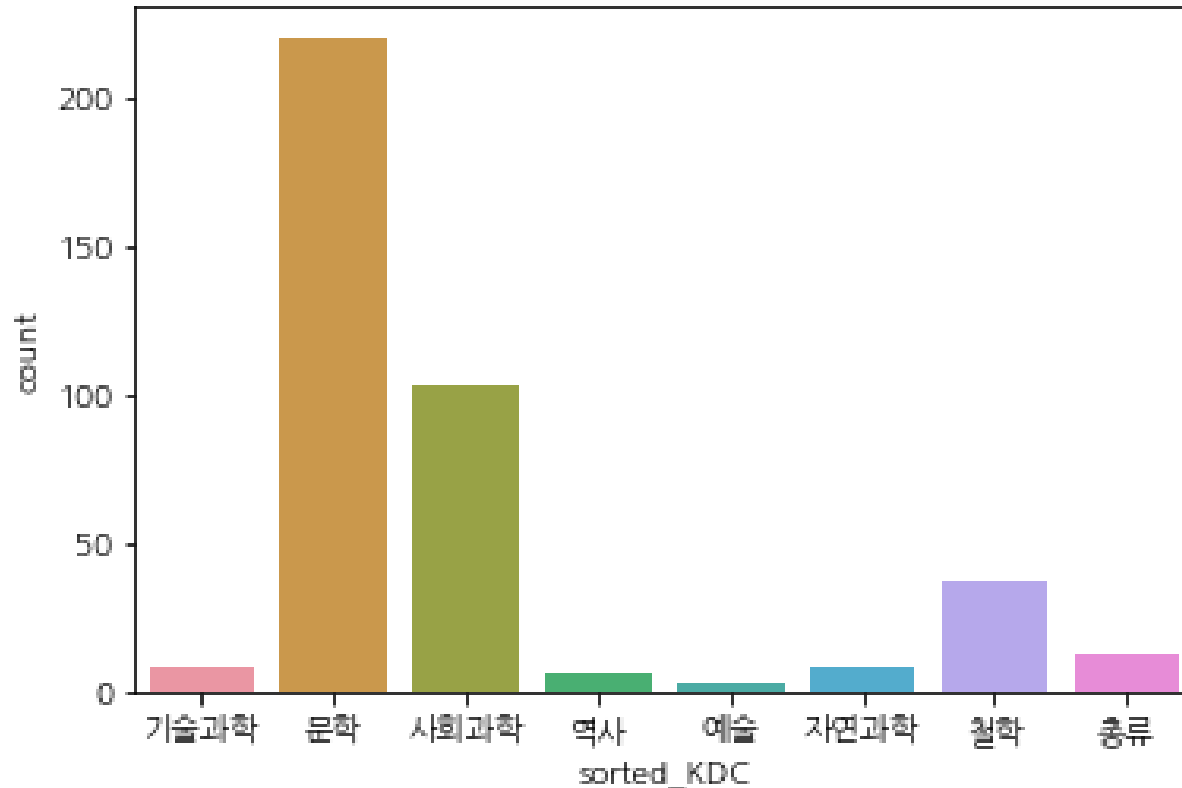
[2022년 12월 10일까지 남녀 대출 건수 비교]

남성 70985
여성 295929

⇒ 여성이 남성의 4배 이상 대출건수가 많은 것을 확인할 수 있다.

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

2) 인기대출도서의 KDC 항목별 개수 (성별 통합)



#대출 순위 200위 내 KDC 항목별 개수 그래프

```
sns.set_palette("hls")  
df.sort_values(by=['sorted_KDC'], inplace=True) #가나다순 정렬  
sns.countplot(data=df, x='sorted_KDC')
```

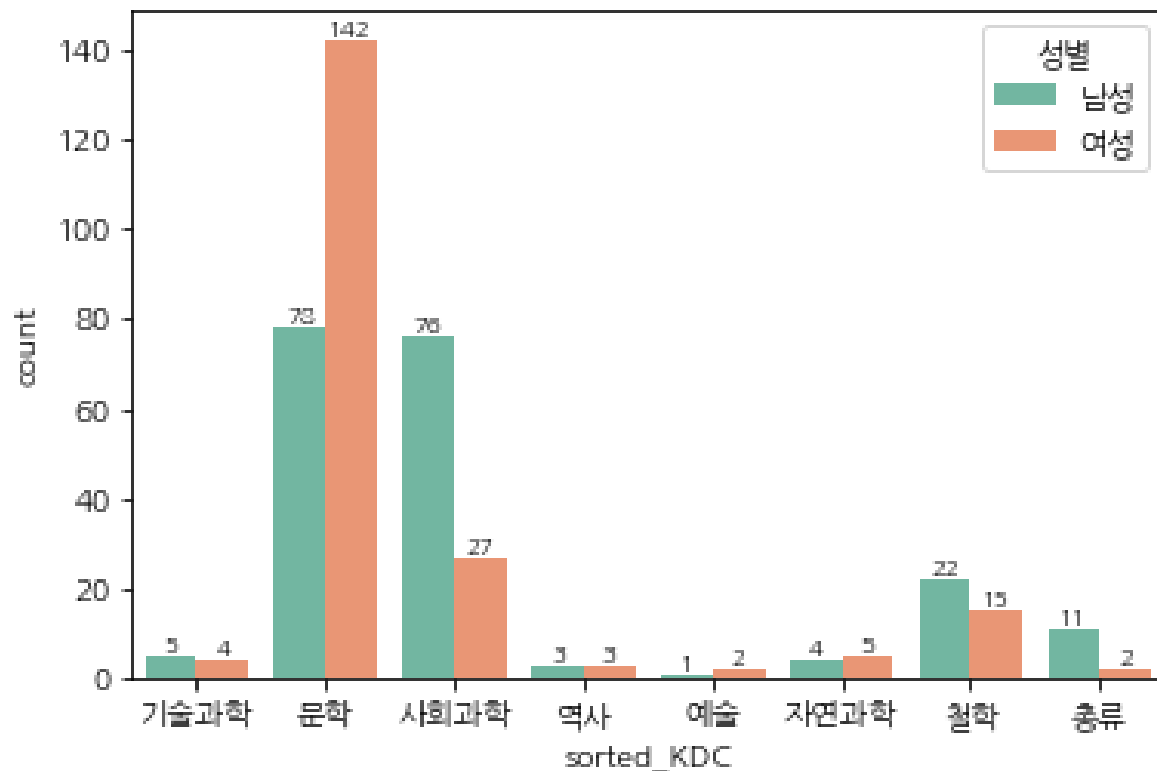
[인기대출도서의 KDC 순위(총 400건)]

문학 -> 사회과학 -> 철학 -> 총류 -> ...

=> 문학과 사회과학이 차지하는 비중이 크다.

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

3) 인기대출도서의 KDC 항목별 개수 (성별 기준)



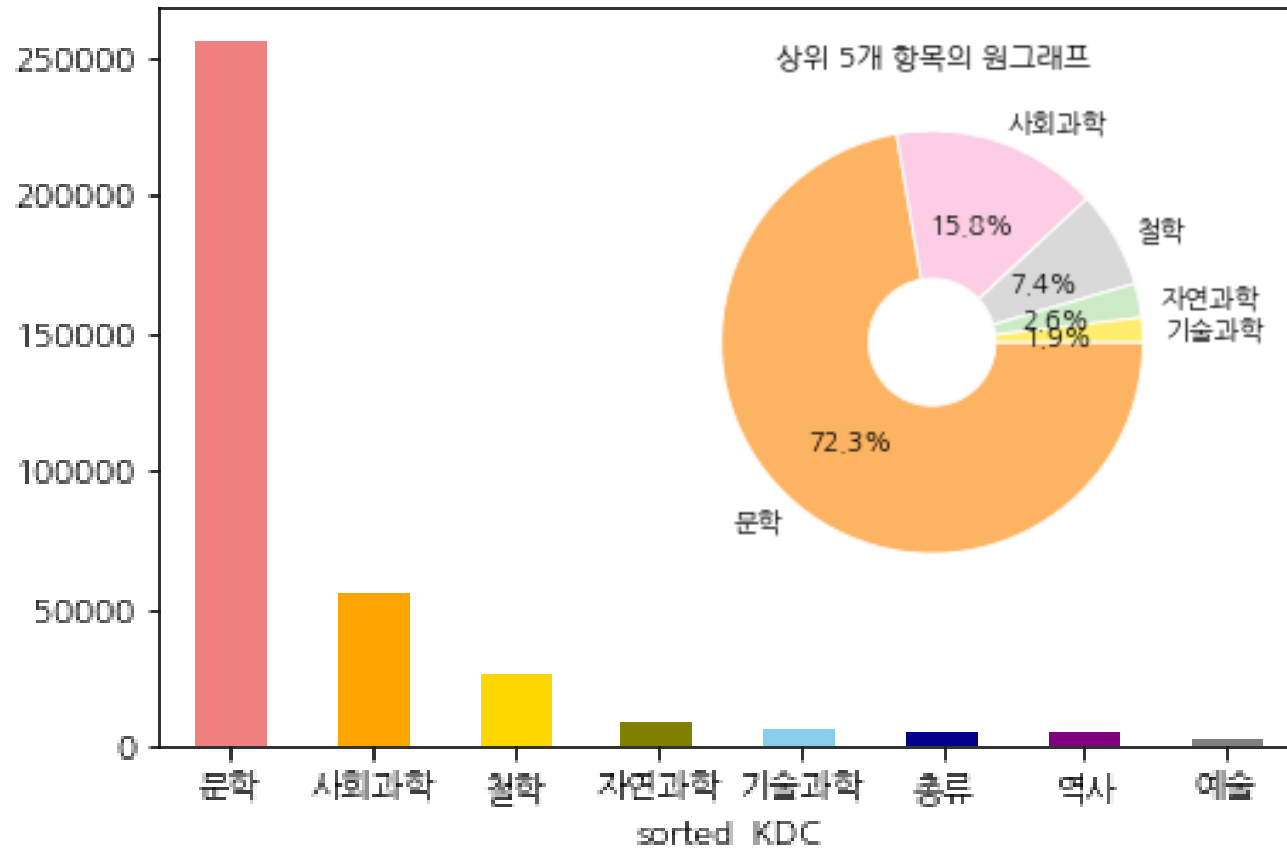
```
sns.set_palette("Set2")
ax = sns.countplot(data=df, x='sorted_KDC', hue='성별')

# countplot에 값 표시
for p in ax.patches:
    height = p.get_height()
    ax.text(p.get_x() + p.get_width() / 2,
            height + 1, height, ha = 'center', size = 9)
```

여성의 경우 인기대출도서 200권 중 문학과 사회과학 도서의 차이가 크지만, 남성의 경우 두 분류가 비등함을 그래프로 확인할 수 있다.

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

4) KDC별 대출건수 시각화 (성별 통합)



[KDC별 대출건수 순위 (성별 통합)]

문학 -> 사회과학 -> 철학 -> 자연과학 ->
기술과학 -> 총류 -> 역사 -> 예술

상위 5개 항목의 원그래프로 비율을 비교했을 때, 문학의 비율이 72.3% 그 다음으로 많은 사회과학이 15.8%로 **문학의 총 대출건수가 압도적**임을 확인할 수 있다.

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

4) KDC별 대출건수 시각화 (성별 통합)

20대 전체 기준 KDC별 대출건수 시각화 (bar chart)

```
sumKDC = df.groupby(df['sorted_KDC'])['대출건수'].sum().sort_values(ascending=False)
print(sumKDC)
colors=['lightcoral', 'orange', 'gold', 'olive', 'skyblue', 'darkblue', 'purple', 'grey']
sumKDC.plot.bar(color=colors, rot=0)
```

20대 전체 기준 KDC별 대출건수 시각화(상위 5개 항목의 pie chart)

```
sumKDC2 = sumKDC.head() #상위 5개 행 추출
print(sumKDC2)
```

```
plt.axis('equal')
plt.title('상위 5개 항목의 원그래프')
color_list = plt.cm.Set3(np.linspace(0.45, 1, 5))
wedgeprops={'width': 0.7, 'edgecolor': 'w', 'linewidth': 1}
```

autopct 숫자 소수점 한자리 표현 #shadow=True

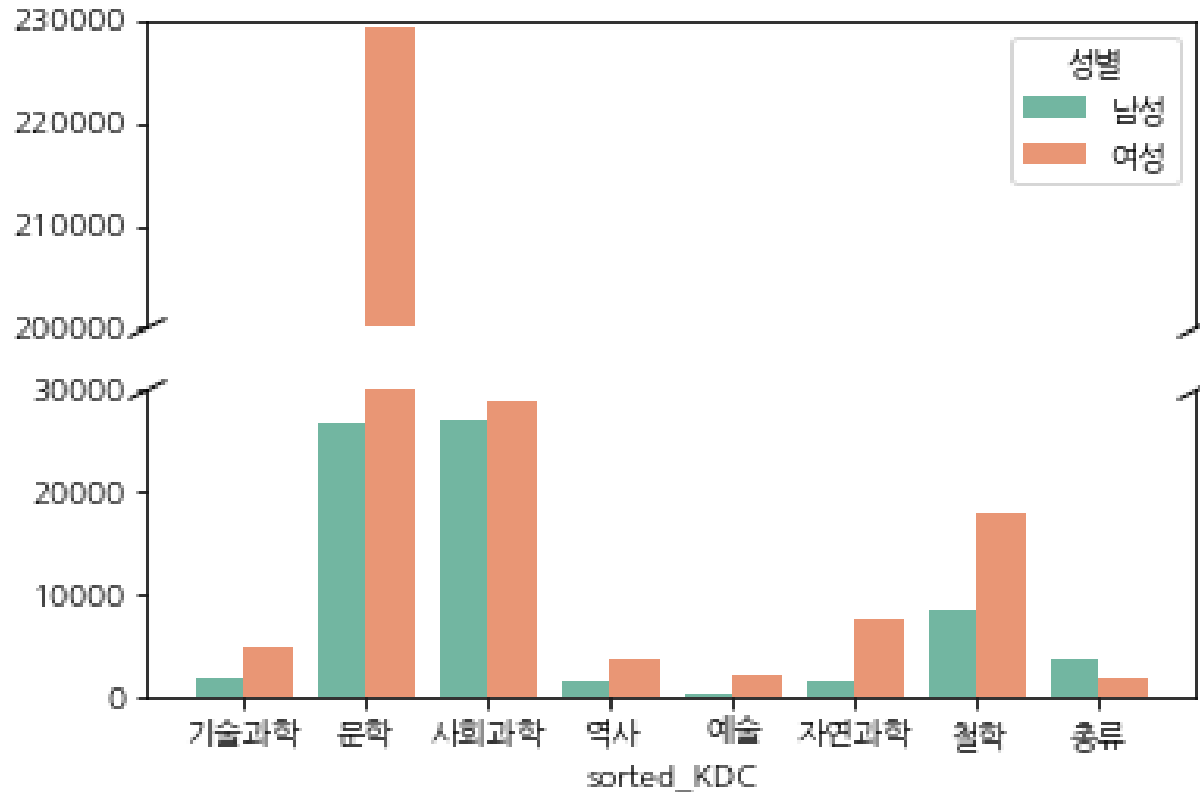
```
plt.pie(sumKDC2, labels=sumKDC2.index, colors=color_list, autopct='%.1f%%', counterclock=False,
        wedgeprops=wedgeprops, textprops={'fontsize': 11}, )
plt.show()
```

```
sorted_KDC
문학      255819
사회과학  55997
철학      26224
자연과학  9199
기술과학  6616
총류      5450
역사      5304
예술      2305
Name: 대출건수, dtype: int64
```

```
sorted_KDC
문학      255819
사회과학  55997
철학      26224
자연과학  9199
기술과학  6616
Name: 대출건수, dtype: int64
```

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

5) KDC별 대출건수 시각화 (성별 기준)



[KDC별 대출건수 순위 (성별 기준)]

남녀 총 대출건수 차이로 인해 대부분의 항목들이 여성의 대출건수가 많지만,

사회과학의 경우 여성 28,995건, 남성 27,002건으로 비슷하며 총류의 경우 여성 1,870, 남성 3,580건으로 남성의 대출건수가 높다는 특징이 있다.

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

5) KDC별 대출건수 시각화 (성별 기준)

```
# 20대 전체 기준 KDC별 대출건수 시각화 (bar chart)

sumKDC = df.groupby(['성별','sorted_KDC'])['대출건수'].sum()
sumKDC = sumKDC.to_frame().reset_index()
print(sumKDC)

colors=['lightcoral', 'orange', 'gold', 'olive', 'skyblue', 'darkblue', 'purple', 'grey']

ax1 = plt.subplot(2, 1, 1)
ax1 = sns.barplot(data=sumKDC, x='sorted_KDC', y='대출건수', hue='성별')
plt.xticks(visible=False)

ax2 = plt.subplot(2, 1, 2, sharex=ax1)
ax2 = sns.barplot(data=sumKDC, x='sorted_KDC', y='대출건수', hue='성별')

ax1.set_ylim(200000, 230000)
ax2.set_ylim(0, 30000)
```

groupby 이후 to_frame으로 시리즈를 데이터프레임으로 변환하고, reset_index로 인덱스를 재설정한다.

```
# 그래프 사이의 경계선 제거
ax1.spines['bottom'].set_visible(False)
ax2.spines['top'].set_visible(False)
ax1.xaxis.tick_top()
ax1.tick_params(labeltop=False)
ax2.xaxis.tick_bottom()

ax2.get_legend().remove()
ax1.axes.xaxis.set_visible(False)

# y축에 물결선 표시
kwargs = dict(marker=[(-1, -0.5), (1, 0.5)], markersize=12,
                linestyle="none", color='k', mec='k', mew=1, clip_on=False)
ax1.plot([0, 1], [0, 0], transform=ax1.transAxes, **kwargs)
ax2.plot([0, 1], [1, 1], transform=ax2.transAxes, **kwargs)

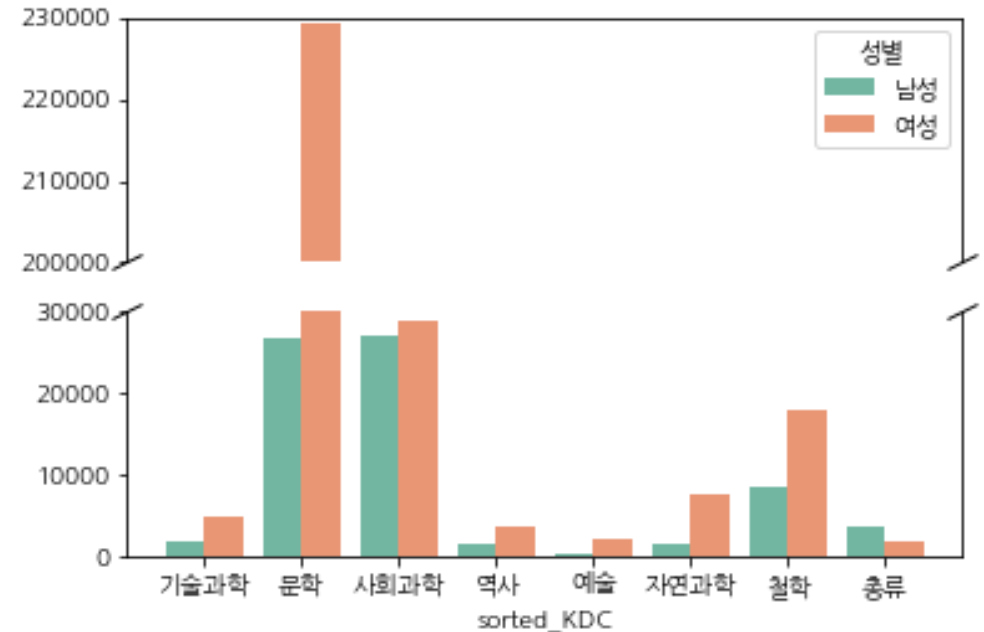
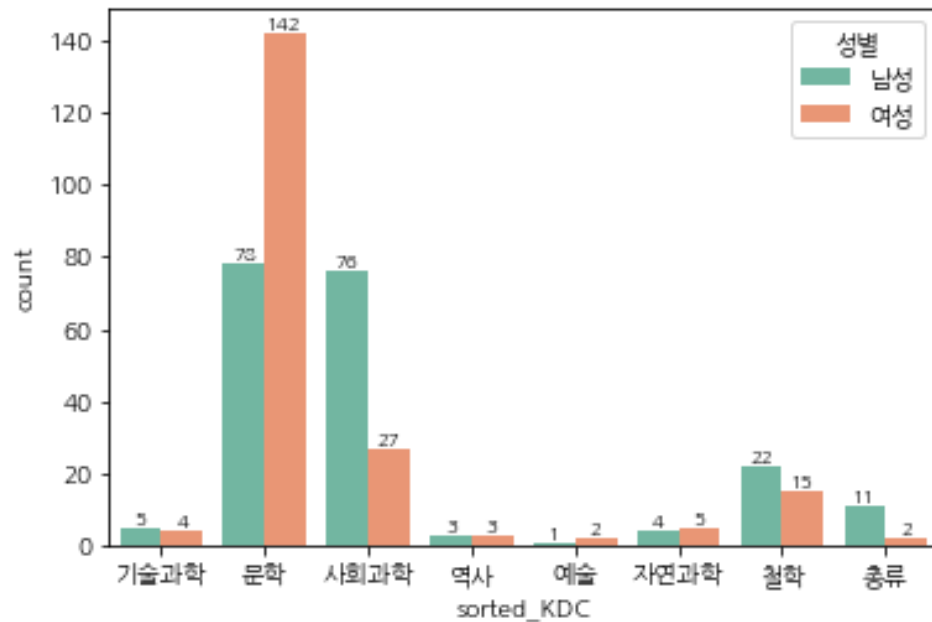
plt.show()
```

	성별	sorted_KDC	대출건수
0	남성	기술과학	1812
1	남성	문학	26648
2	남성	사회과학	27002
3	남성	역사	1566
4	남성	예술	267
5	남성	자연과학	1678
6	남성	철학	8432
7	남성	종류	3580
8	여성	기술과학	4804
9	여성	문학	229171
10	여성	사회과학	28995
11	여성	역사	3738
12	여성	예술	2038
13	여성	자연과학	7521
14	여성	철학	17792
15	여성	종류	1870

여성의 문학 대출건수(229,170건)가 다른 항목의 값에 비해 7배 이상 높아서 그래프를 2개 그려 합친다.

3. 데이터 시각화 - 20대 성별에 따른 대출 트렌드 분석

6) 인기대출도서의 KDC 항목별 개수와 KDC 항목별 대출건수 비교



⇒ 사회과학의 경우 여성이 27권으로 남성(76권)보다 49권 적지만 대출건수는 높다.

⇒ 총류의 경우 컴퓨터과학, 프로그래밍 도서, 지식 및 학문 일반 도서로 인해 남성이 높은 수치를 보였다.

4. 데이터 전처리 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석

1) csv -> 데이터프레임 변환

```
# csv 불러오기
file_path_man1 = '/content/2021_1분기_남성.csv'
file_path_man2 = '/content/2021_2분기_남성.csv'
file_path_man3 = '/content/2021_3분기_남성.csv'
file_path_man4 = '/content/2021_4분기_남성.csv'

file_path_woman1 = '/content/2021_1분기_여성.csv'
file_path_woman2 = '/content/2021_2분기_여성.csv'
file_path_woman3 = '/content/2021_3분기_여성.csv'
file_path_woman4 = '/content/2021_4분기_여성.csv'
```

8개의 CSV를 앞서 변환한 방법과 같은 방식으로 데이터 프레임으로 변환한다.

4. 데이터 전처리 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석

2) 성별, 분기 열 추가

```
# 성별, 분기 열 추가
man_list = [df_man1, df_man2, df_man3, df_man4]
woman_list = [df_woman1, df_woman2, df_woman3, df_woman4]
num_list = [1, 2, 3, 4]

for i, j in zip(man_list, num_list):
    i['성별'] = '남성'
    i['분기'] = str(j) + '분기'

for i, j in zip(woman_list, num_list):
    i['성별'] = '여성'
    i['분기'] = str(j) + '분기'
```

분기별 분석을 위해, 각 csv 별로 분기 열을 추가한다.

3) 데이터프레임 병합, 가공

8개의 데이터프레임을 병합하고 불필요한 행을 삭제한 뒤, KDC 결측치를 처리한다.
(기존 2-5번과 동일)

이후 KDC 처리 함수를 통해 한글로 된 KDC 열을 추가한다. (기존 2-7번과 동일)

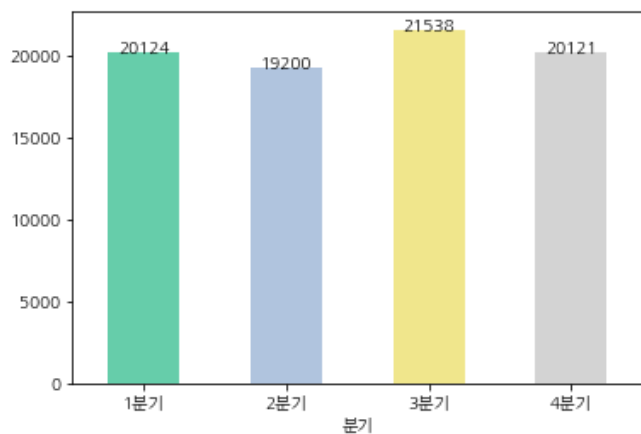
#	Column	Non-Null Count	Dtype
0	순위	1600 non-null	int64
1	서명	1600 non-null	object
2	저자	1600 non-null	object
3	출판사	1600 non-null	object
4	출판년도	1600 non-null	object
5	KDC	1600 non-null	float64
6	대출건수	1600 non-null	int64
7	성별	1600 non-null	object
8	분기	1600 non-null	object
9	sorted_KDC	1600 non-null	object

dtypes: float64(1), int64(2), object(7)

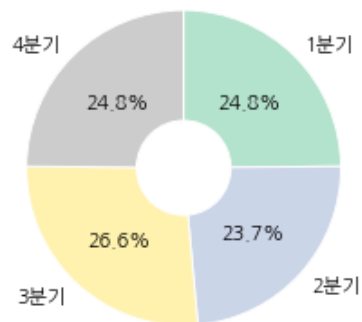
5. 데이터 시각화 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석

1) 분기별 대출건수 비교 (bar, pie chart)

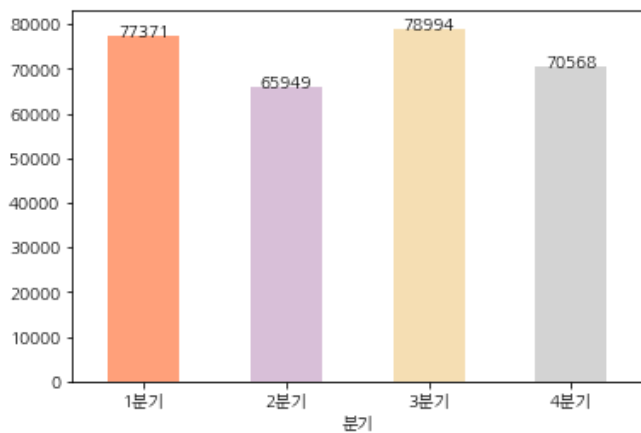
남성



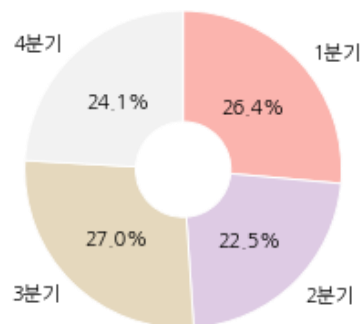
남성의 분기별 대출비율



여성



여성의 분기별 대출비율



[20대 남녀의 공통적인 특징]

⇒ 1, 3분기 대출량 증가

⇒ 2, 4분기 대출량 감소

5. 데이터 시각화 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석

1) 분기별 대출건수 비교 (bar, pie chart)

```
# 여성
quater_woman = df_quater[df_quater['성별']=='여성']
grouped_woman = quater_woman.groupby(['분기'])['대출건수'].sum()
print(grouped_woman)
```

```
# pie chart
plt.axis('equal')
plt.title('여성의 분기별 대출비율')
color_list = plt.cm.Pastel1((np.linspace(0, 1, 4)))
wedgeprops={'width': 0.7, 'edgecolor': 'w', 'linewidth': 1}
```

```
# autopct 숫자 소수점 한자리 표현
plt.pie(grouped_woman, labels=grouped_woman.index, colors=color_list, autopct='%0.1f%%',
        counterclock=False, wedgeprops=wedgeprops, textprops={'fontsize': 11}, startangle=90)
plt.show()
```

```
# bar chart
colors=['lightsalmon', 'thistle', 'wheat', 'lightgrey']
ax = grouped_woman.plot.bar(color=colors, rot=0)
```

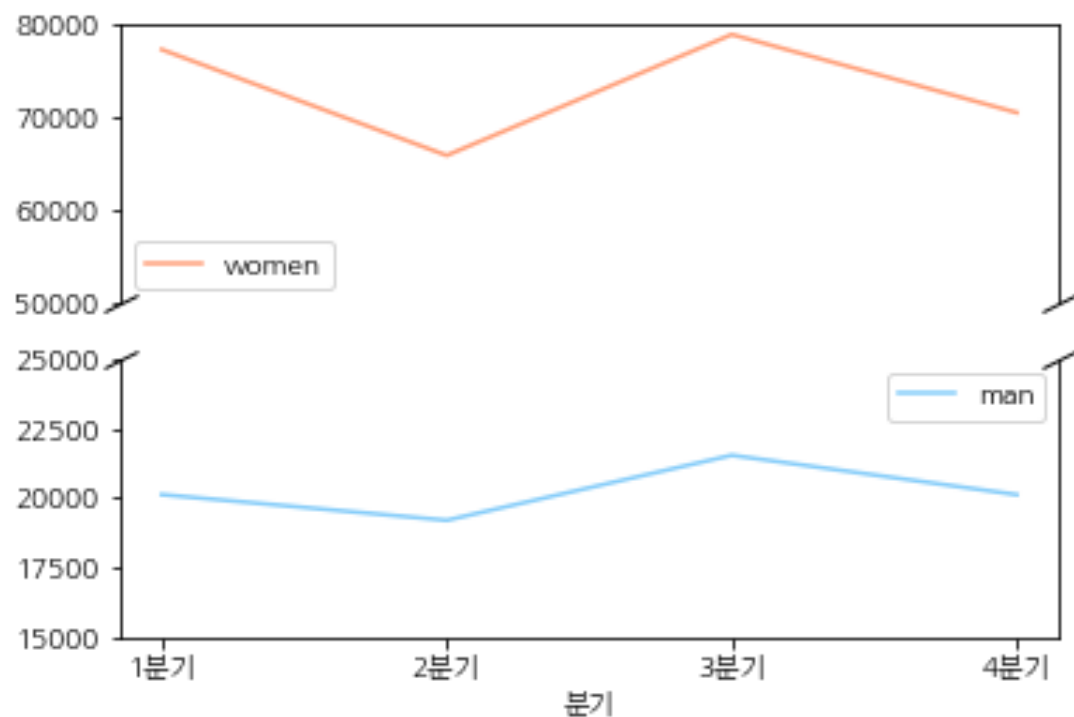
```
for p in ax.patches:
    height = p.get_height()
    ax.text(p.get_x() + p.get_width() / 2, height + 5, height, ha = 'center', size = 10)
```

```
#남성
quater_man = df_quater[df_quater['성별']=='남성']
grouped_man = quater_man.groupby(['분기'])['대출건수'].sum()
```

Bar chart에 값을 표시하여 가독성을 높이고
분기별 비율을 알기 위해 Pie chart를 사용함.

5. 데이터 시각화 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석

1) 분기별 대출건수 비교 (line chart)



```
#남녀 라인플롯 합치기
ax1 = plt.subplot(2, 1, 1)
ax1 = sns.lineplot(data=grouped_woman, label='women', color='lightsalmon')
plt.xticks(visible=False)

ax2 = plt.subplot(2, 1, 2, sharex=ax1)
ax2 = sns.lineplot(data=grouped_man, label='man', color='lightskyblue')

ax1.set_ylim(50000, 80000)
ax2.set_ylim(15000, 25000)

# 그래프 사이의 경계선 제거
ax1.spines['bottom'].set_visible(False)
ax2.spines['top'].set_visible(False)
ax1.xaxis.tick_top()
ax1.tick_params(labeltop=False)
ax2.xaxis.tick_bottom()

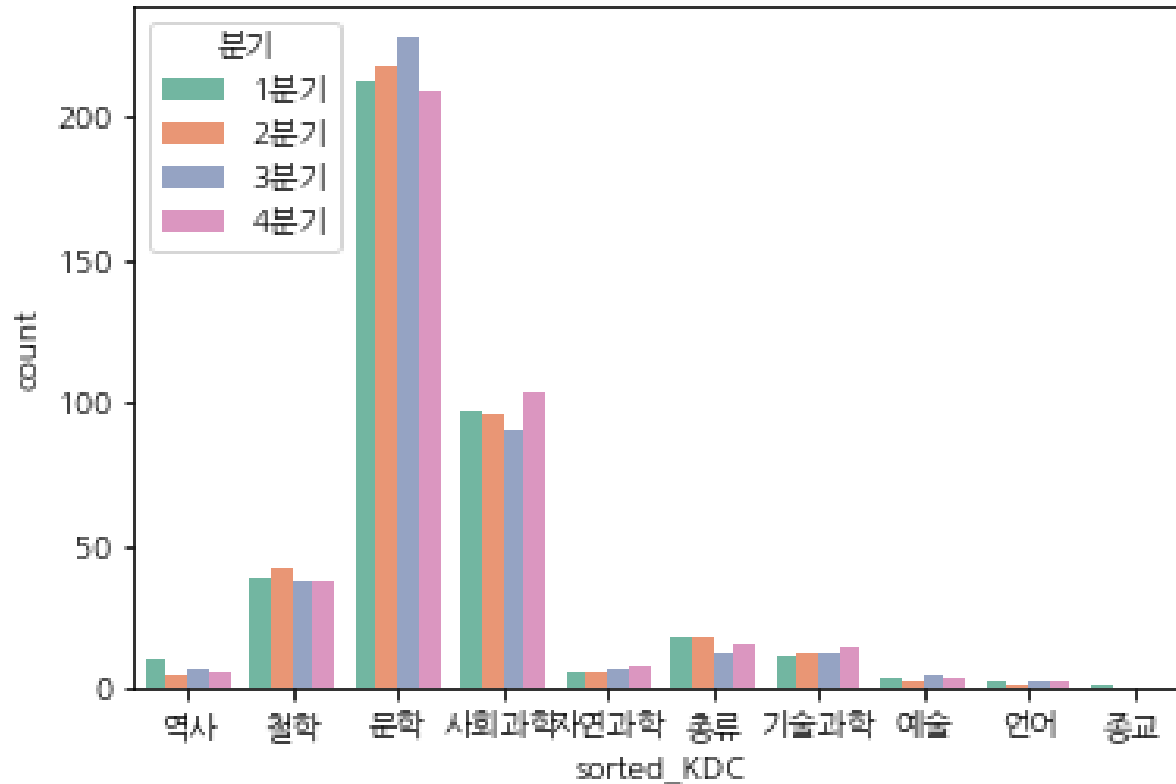
# ax2.get_legend().remove()
ax1.axes.xaxis.set_visible(False)

# y축에 물결선 표시
kwargs = dict(marker=[(-1, -0.5), (1, 0.5)], markersize=12,
                linestyle="none", color='k', mec='k', mew=1, clip_on=False)
ax1.plot([0, 1], [0, 0], transform=ax1.transAxes, **kwargs)
ax2.plot([0, 1], [1, 1], transform=ax2.transAxes, **kwargs)

plt.show()
```

5. 데이터 시각화 - 2021년 기준, 20대 성별 · 분기별 대출건수와 KDC 분석

2) 분기별 대출도서의 KDC 비교 (성별 통합)



```
# 분기별 대출도서 유형 분류
sns.set_palette("Set2")
sns.countplot(data=df_quater, x='sorted_KDC', hue='분기')
```

분기별 대출 도서의 KDC 건수를 비교한다.

문학은 1~3분기에 증가 추세를 보이다가 4분기에 대출건수가 감소하여 분기 중 가장 적은 대출건수를 보였고, 사회과학은 4분기의 대출건수가 가장 많은 것을 확인할 수 있다.

6. 프로젝트 리뷰

도서관 빅데이터를 활용해서 20대 성별에 따른 대출 트렌드를 파악할 수 있었다.
한 학기 동안 배운 내용을 통해 직접 데이터를 수집하여 분석한 뜻깊은 경험이었다.

[주요한 특징]

- 20대 여성은 20대 남성에 비해 4배 이상 대출량이 많다.
- 20대 여성의 대출량 1위는 문학이며, 이는 2위인 사회과학의 7배 이상의 대출량이다.
- 20대 남성의 대출량 1, 2위인 문학과 사회과학은 매우 근소한 차이이다.
- 20대 남성은 20대 여성에 비해 총류에 속하는 도서에 관심이 많다.
- 도서관의 대출량은 1,3분기가 2,4분기에 비해 높다.

[보완할 점]

- 20대에 한정해서 데이터를 분석해 보았는데, 전 연령층을 대상으로 분석하지 못한 아쉬움이 남는다.
- 전 연령층을 대상으로 데이터를 분석하면 도서관 운영 시 도움이 되는 정보가 도출될 것이라고 생각한다.
- 도서관 정보나루의 데이터와 국가통계포털의 전자책 관련 데이터를 비교해 보고 싶다는 생각이 들었다.
ex) 도서관 대출량 추이와 전자책 독서량 추이 비교 · 예측, 도서관 대출 도서와 전자책 선호도서 비교 등