

```
In [5]: ## 数据操作
import numpy as np
import pandas as pd
pd.set_option('display.max_rows', 100)
pd.set_option('display.unicode.ambiguous_as_wide', True)
pd.set_option('display.unicode.east_asian_width', True)

from tqdm import tqdm
```

```
In [4]: # 读取文件
data = pd.read_csv('E:\爬虫\jjwx\无差别爬取数据\无差别爬取.csv')
data
```

D:\anaconda2022\lib\site-packages\IPython\core\interactiveshell.py:3444: DtypeWarning: Columns (0,3,4,5,6,9,13,14,15,16,17,19,20,24) have mixed types.Specify dtype option on import or set low_memory=False.

```
exec(code_obj, self.user_global_ns, self.user_ns)
```

Out[4]:

	a_作品 字数	a_作品 状态	a_作 者	a_作 者所 发送 红包 数	a_小 说 完 本 数	a_小 说 暂 停 数	a_小 说 连 载 数	a_最 后 更 新 时 间	a_最 近 更 新 作 品	a_被 收 藏 数	...	b_总 点 击 量	b_文 章 名 称	b_文 章 当 前 被 收 藏 数	b_文 章 积 分	b_文 章 分 类
0	157033	连载	挖坑不填	40	121	25	173	2022-02-20 22:07:26	《[主名柯]关于我有三话存稿但还是被催稿是不是哪里搞错了》...	25490	...	39657	[综]尤莉卡的秘密恋爱关系	865	98595808	行生情史
1	60744	连载	一宫兴子	20	18	0	25	2020-12-13 15:34:56	《[综英美]而那个成功的反派却想要回到过去》	1530	...	22671	[综英美]而那个成功的反派却想要回到过去	591	9550995	行生情史

	a_作品 字数	a_作品 状态	a_作者	a_作 者所 发送 红包 数	a_小 说 完 本 数	a_小 说 暂 停 数	a_小 说 连 载 数	a_最后 更新 时间	a_最近 更新 作品	a_被 收藏 数	...	b_总点 击量	b_文章 名称	b_文 章 当 前 被 收 藏 数	b_文章 积 分	b_...
2	60744	连载	一宫 兴子	20	18	0	25	2020- 12-13 15:34:56	《[综 英美] 而那个成功的反派却想要回到过去》	1530	...	123219	[TSN/ME] 当你的脑子里多了一个人的正确解决方法	904	15546017	行 生 结 爱 的 代 言 行 人
3	157033	连载	挖坑 不填	40	121	25	173	2022- 02-20 22:07:26	《[主 名柯] 关于我有三话存稿但还是被催稿是不是哪里搞错了》...	25490	...	10135	[猫鼠游 戏]圆桌 骑士和亚瑟王的圣杯	110	90343528	行 生 情 迷 月 史 行 人
4	157033	连载	挖坑 不填	40	121	25	173	2022- 02-20 22:07:26	《[主 名柯] 关于我有三话存稿但还是被催稿是不是哪里搞错了》...	25490	...	17097	[奇幻贵 公子/恶 灵猎人] 生当复来归	263	93998864	行 生 情 迷 月 史 行 人
...
320276	146993	连载	写文 字的 静一	109	1	0	1	2021- 06-07 10:40:42	《替 嫁假 千金 每天都想休夫》	52	...	13368	醉扶青山	309	16599650	月 创 情 迷 月 史 行 人

In [14]:

```
## 字段, nunique等基本信息
def get_info(df, head=10):
    print(df.shape)
    print('\n', '-----各列特征信息如下-----')
    stats = []
    for col in df.columns:
        stats.append(
            (col, df[col].nunique(),
             round(df[col].isnull().sum()*100 / df.shape[0], 3),
             round(df[col].value_counts(normalize=True, dropna=False).values[0]*100, 3),
             df[col].dtype)
        )
    stats_df = pd.DataFrame(stats, columns=['特征', '属性个性', '缺失值占比', '最大属性占比', '特征类型'])
    print(stats_df.sort_values('缺失值占比', ascending=False).head(head))
```

In [15]:

```
# 获取数据的基本信息
get_info(data, 100)
```

(318676, 27)

```
-----各列特征信息如下-----
      特征  属性个性  缺失值占比  最大属性占比  特征类型
11      b_作品视角          7      0.291      55.443  object
12      b_作品风格          7      0.026      58.647  object
18      b_文章名称    305010      0.003       0.061  object
13      b_全文字数    187112      0.001      19.996  object
15      b_总书评数    19960      0.001      20.294  object
25      b_霸王票全站排行  127180      0.001      47.866  object
24      b_营养液数    22149      0.001      31.990  object
23      b_签约状态          2      0.001      86.713  object
22      b_文章进度          3      0.001      49.290  object
21      b_文章类型        452      0.001      21.247  object
20      b_文章积分    277016      0.001       0.072  object
19      b_文章当前被收藏数  33585      0.001       2.338  object
17      b_总点击量    170228      0.001      19.961  object
16      b_总共地雷数量    9512      0.001      37.744  object
26      piece_url    318663      0.001       0.001  object
14      b_前进一名所需地雷数  379      0.001      70.969  object
1      a_作品状态          4      0.000      69.971  object
10      author_url    38072      0.000       0.074  object
9      a_被收藏数    6731      0.000       0.340  object
8      a_最近更新作品    37173      0.000       0.074  object
7      a_最后更新时间    35491      0.000       1.057  object
6      a_小说连载数     150      0.000       7.697  object
5      a_小说暂停数      54      0.000      50.026  object
4      a_小说完本数     185      0.000       6.400  object
3      a_作者所发送红包数    4287      0.000      21.450  object
2      a_作者        38072      0.000       0.074  object
0      a_作品字数    38013      0.000       0.104  object
```

7种作品视角, 305010部小说, 作品视角缺失率占比0.291, 特征类型全部为object,需要合理的转换, 比如总书评数, 应该为int,整型

In [67]:

```
# 数据转换
data['b_霸王票全站排行'] = data['b_霸王票全站排行'].replace(r'暂无', np.nan)
int_cols = ['a_作品字数', 'a_作者所发送红包数', 'a_小说完本数', 'a_小说暂停数', 'a_小说连载数', 'a_总书评数', 'b_总书评数', 'b_总共地雷数量', 'b_总点击量', 'b_文章当前被收藏数', 'b_文章积分', 'b_营养液数']
dtypes_ = {i: np.int64 for i in int_cols}
```

```
data = data.astype(dtypes_)
data['a_最后更新时间'] = pd.to_datetime(data['a_最后更新时间'])
```

In [74]: data.dtypes

Out[74]:

a_作品字数	int64
a_作品状态	object
a_作者	object
a_作者所发送红包数	int64
a_小说完本数	int64
a_小说暂停数	int64
a_小说连载数	int64
a_最后更新时间	datetime64[ns]
a_最近更新作品	object
a_被收藏数	int64
author_url	object
b_作品视角	object
b_作品风格	object
b_全文字数	int64
b_前进一名所需地雷数	int64
b_总书评数	int64
b_总共地雷数量	int64
b_总点击量	int64
b_文章名称	object
b_文章当前被收藏数	int64
b_文章积分	int64
b_文章类型	object
b_文章进度	object
b_签约状态	object
b_营养液数	int64
b_霸王票全站排行	object
piece_url	object
dtype:	object

所有字段被转换成了正确的类型，正确的类型有助于数据的分析

In [78]: data.head()

Out[78]:

	a_作品 字数	a_作品 状态	a_作者	a_作者 所发送 红包数	a_小说 完本数	a_小说 暂停数	a_小说 连载数	a_最后 更新时 间	a_最 近更 新作品	a_被 收藏 数	...	b_总点 击量	b_文章 名称	b_文 章当 前被 收藏 数	b_文章 积分	b_文 章类 型	b_文 章进 度	b_签 约状 态	b_营养 液数
0	157033	连载	挖坑不填	40	121	25	173	2022-02-20 22:07:26	《[主名柯]关于我有三话存稿但还是被催稿是不是哪里搞错了》...	25490	...	39657	[综]尤莉卡的秘密恋爱关系	865	98595808	衍生-言情-架空历史-东方衍生	连载	已签约	6

	a_作品 字数	a_作品 状态	a_作者	a_作者所 发送红包数	a_小说 完本数	a_小说 暂停数	a_小说 连载数	a_最后 更新时间	a_最近 更新作品	a_被 收藏数	...	b_总点 击量	b_文章 名称	b_文章 当前被收 藏数	b_文章 积分	b_文章 类型	b_文章 进度	b_签 约状态	b_理 身评数
1	60744	连载	一宫 兴子	20	18	0	25	2020- 12-13 15:34:56	《[综 英美] 而那个 成功的 反派却 想要回 到过去》	1530	...	22671	[综英美] 而那个成 功的反派 却想要回 到过去	591	9550995	衍生-纯 爱-架空 历史-西 方衍生	连载	已签约	43
2	60744	连载	一宫 兴子	20	18	0	25	2020- 12-13 15:34:56	《[综 英美] 而那个 成功的 反派却 想要回 到过去》	1530	...	123219	[TSN/ME] 当你的脑 子里多了 一个人的 正确解决 方法	904	15546017	衍生-纯 爱-近代 现代-西 方衍生	完结	已签约	29
3	157033	连载	挖坑不 填	40	121	25	173	2022- 02-20 22:07:26	《[主 名柯] 关于我 有三话 存稿但 还是被 催稿是 不是哪 里搞错 了》...	25490	...	10135	[猫鼠游 戏]圆桌 骑士和亚 瑟王的圣 杯	110	90343528	衍生-言 情-架空 历史-西 方衍生	完结	已签约	1
4	157033	连载	挖坑不 填	40	121	25	173	2022- 02-20 22:07:26	《[主 名柯] 关于我 有三话 存稿但 还是被 催稿是 不是哪 里搞错 了》...	25490	...	17097	[奇幻贵 公子/恶 灵猎人] 生当复来 归	263	93998864	衍生-言 情-架空 历史-东 方衍生	完结	已签约	4

5 rows × 27 columns

In [76]:

```
# 将结果保存, 用于数据分析, 并覆盖源数据  
data.to_csv('E:\爬虫\jjwx\无差别爬取数据\无差别爬取.csv', index=False)
```