

STAR WARS SURVEY DATA ANALYSIS



A PRESENTATION FROM A GALAXY FAR, FAR AWAY



SURVEY DATA BACKGROUND

- The data was taken from a survey conducted via SurveyMonkey on FiveThirtyEight.com
- There were 1186 total respondents
- The data includes questions about people's opinion and attitude about Star Wars films:
 - Whether or not respondents have seen any of the original six movies
 - Whether or not respondents consider themselves to be fans of the franchise
 - Which of the original six movies they had seen
 - How they ranked each of the six movies
- The data also includes how favorably they ranked the major characters
- Each respondent also provides demographic data about themselves, such as their age, gender, location, household income and education



DATA PREPARATION





DATA PREPARATION

- Changed the column names, and had separate columns for the sub-columns (check boxes)
- Filled all the columns of whether respondents had seen each of the six movies with 'Yes' or 'No' if they had checked the box for that movie

DATA PREPARATION

TYPE CORRECTION

- Looped through each column and printed out the unique values for that column.

```
expanded_universe_familiarity: ['Yes' nan 'No']  
expanded_universe_fan: ['No' nan 'Yes' 'Yess']  
star_trek_fan: ['No' 'Yes' nan 'Noo' 'yes' 'no ']  
gender: ['Male' nan 'Female' 'F' 'female' 'male']
```

- Made necessary corrections to these typos:

Unique value counts in the column star_trek_fan

No	639
Yes	426
yes	1
Noo	1
no	1

BEFORE

Unique value counts in the column star_trek_fan

No	641
Yes	427

AFTER



DATA PREPARATION

CONVERSION TO UPPERCASE AND REMOVAL OF EXTRA WHITESPACE

- Removed extra whitespaces, so our exploration functions don't treat a response like 'no ' and 'no' differently since they mean the same thing
- Converted all text values to uppercase for parity

```
star_trek_fan: ['No' 'Yes' nan 'Noo' 'yes' 'no ']  
gender: ['Male' nan 'Female' 'F' 'female' 'male']
```

BEFORE

```
star_trek_fan: ['NO' 'YES' nan]  
gender: ['MALE' nan 'FEMALE']
```

AFTER

DATA PREPARATION

SANITY CHECKS

- Looped through each column and printed out the unique values for that column.

```
age: ['18-29' nan '500' '30-44' '> 60' '45-60']
```

```
household_income: [nan '$0 - $24,999' '$100,000 - $149,999' '$25,000 - $49,999'  
'$50,000 - $99,999' '$150,000+']
```

```
education: ['High school degree' 'Bachelor degree' 'Some college or Associate degree'  
nan 'Graduate degree' 'Less than high school degree']
```

- Made necessary corrections:

Unique value counts in the column 'age'

45-60	291
> 60	269
30-44	268
18-29	217
500	1

BEFORE

Unique value counts in the column 'age'

45-60	291
> 60	269
30-44	268
18-29	217
-1	1

AFTER



DATA PREPARATION

MISSING VALUES

- Missing values in the columns were dealt with by setting them all to (-1)
- Putting a -1 allows us to easily identify occurrences of missing values in each column via indexing in the code
- We did not take the column-wise mean because it doesn't work for most columns, as most columns aren't numeric, and the ones that are happened to represent ordinal data
- Taking the column-wise median is a possibility, however, it may not accurately reflect the truth since some columns are interdependent.

DATA PREPARATION

Jupyter assignment1_presentation Last Checkpoint: Yesterday at 12:43 AM (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

Out[30]:

	respondent_ID	seen_any_movies	consider_themselves_fan	seen_episode_I	seen_episode_II	seen_episode_III	seen_episode_IV	seen_episode_V	seen_episode_VI
0	3292879998	YES	YES	YES	YES	YES	YES	YES	YES
1	3292879538	NO	-1	NO	NO	NO	NO	NO	NO
2	3292765271	YES	NO	YES	YES	YES	NO	NO	NO
3	3292763116	YES	YES	YES	YES	YES	YES	YES	YES
4	3292731220	YES	YES	YES	YES	YES	YES	YES	YES
...
1181	3288388730	YES	YES	YES	YES	YES	YES	YES	YES
1182	3288378779	YES	YES	YES	YES	YES	YES	YES	YES
1183	3288375286	NO	-1	NO	NO	NO	NO	NO	NO
1184	3288373068	YES	YES	YES	YES	YES	YES	YES	YES
1185	3288372923	YES	NO	YES	YES	NO	NO	NO	YES

1186 rows x 38 columns



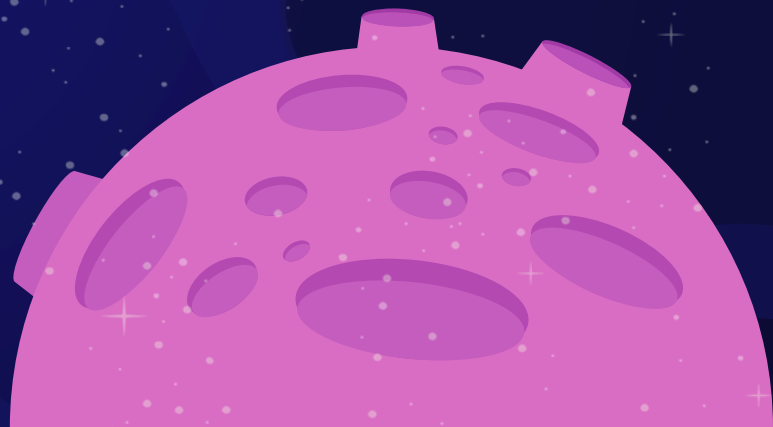
HOW DO PEOPLE RANK STAR WARS MOVIES?

DATA EXPLORATION

First checked if people had ranked the episodes, even if they hadn't seen them.

As an example, for Episode I, around 162 respondents marked that they had not seen this episode yet ranked it.

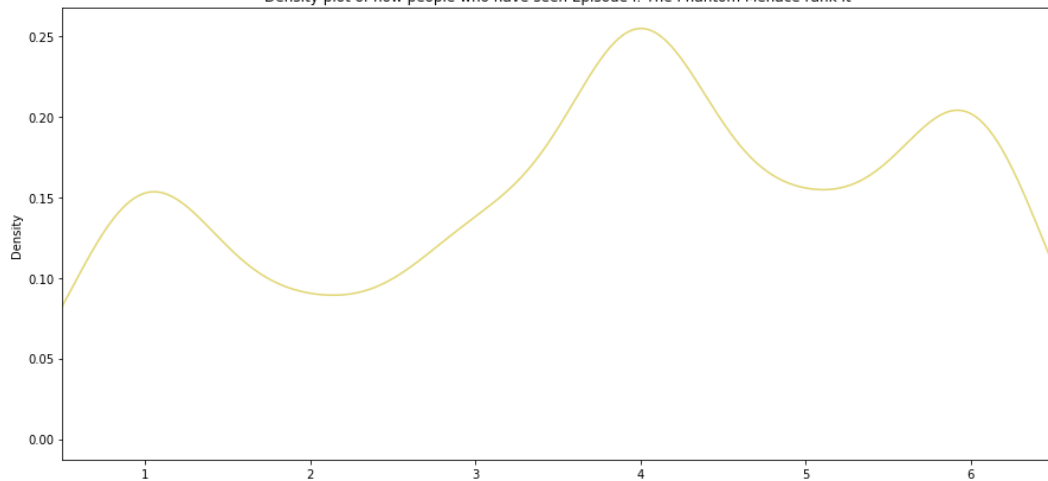
Thus we do not consider these rows when doing our analysis of the rankings as people who haven't seen the movie shouldn't be allowed to rate it.



HOW PEOPLE RANK

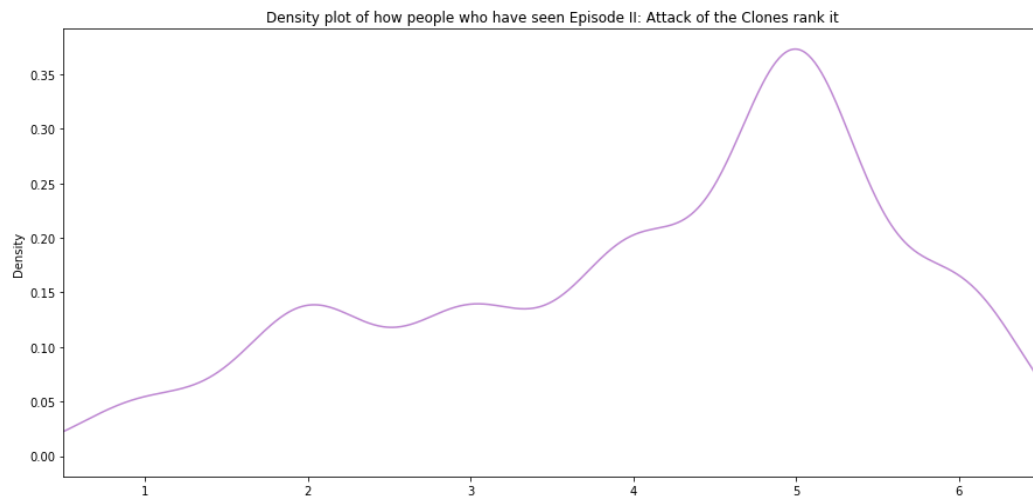
EPISODE I: THE PHANTOM MENACE

Density plot of how people who have seen Episode I: The Phantom Menace rank it

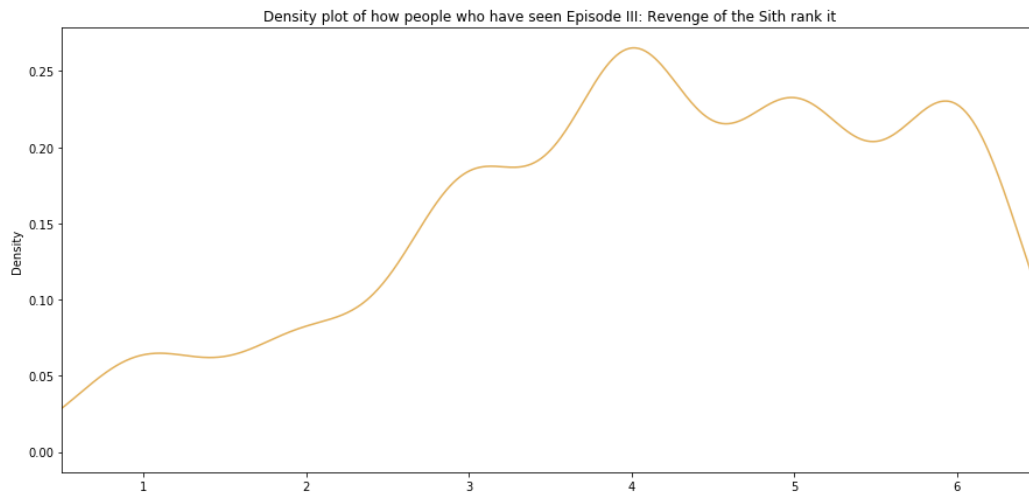


HOW PEOPLE RANK

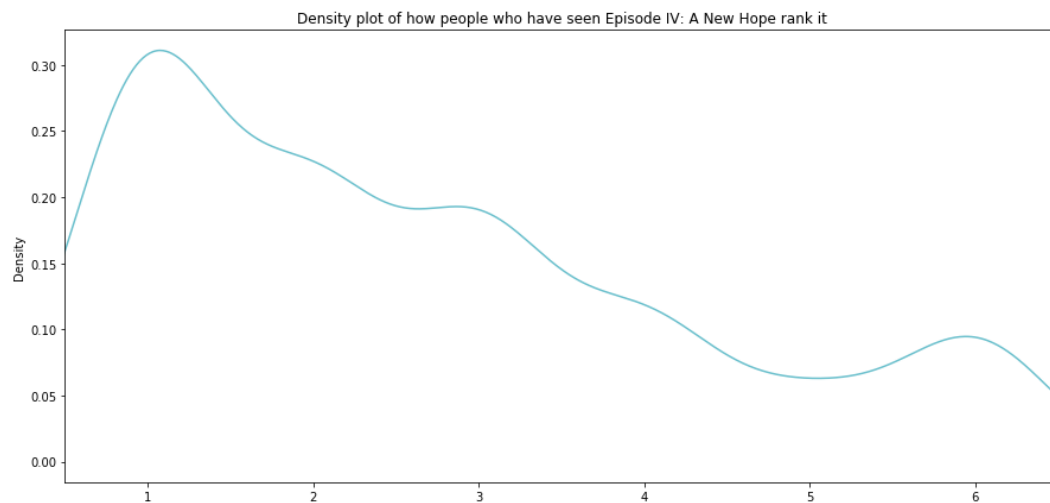
EPISODE II: ATTACK OF THE CLONES



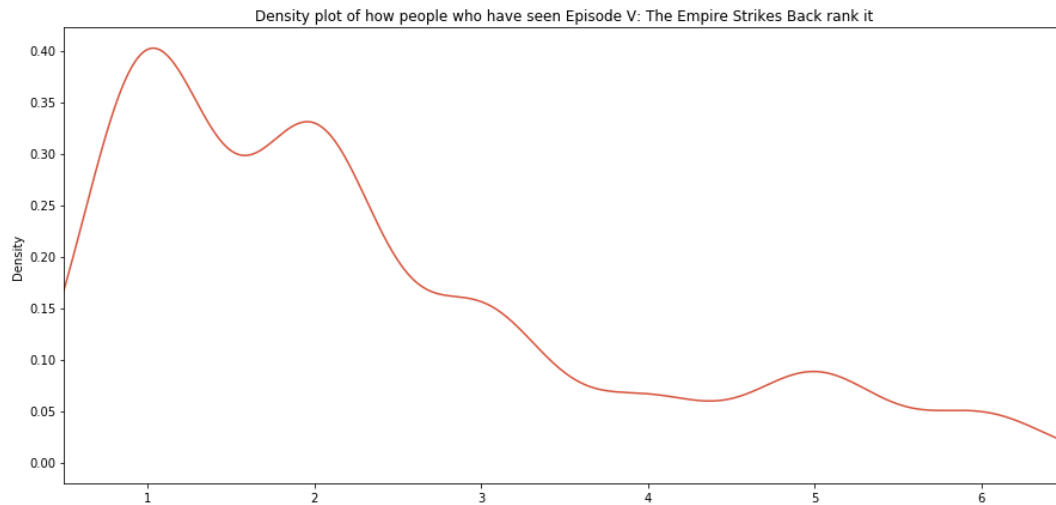
HOW PEOPLE RANK EPISODE III: REVENGE OF THE SITH



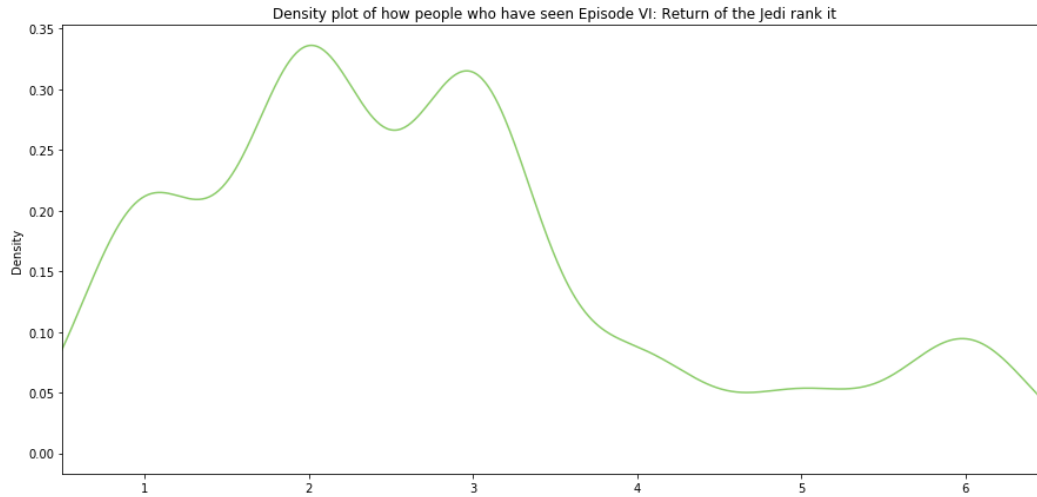
HOW PEOPLE RANK EPISODE IV: A NEW HOPE



HOW PEOPLE RANK EPISODE V: THE EMPIRE STRIKES BACK



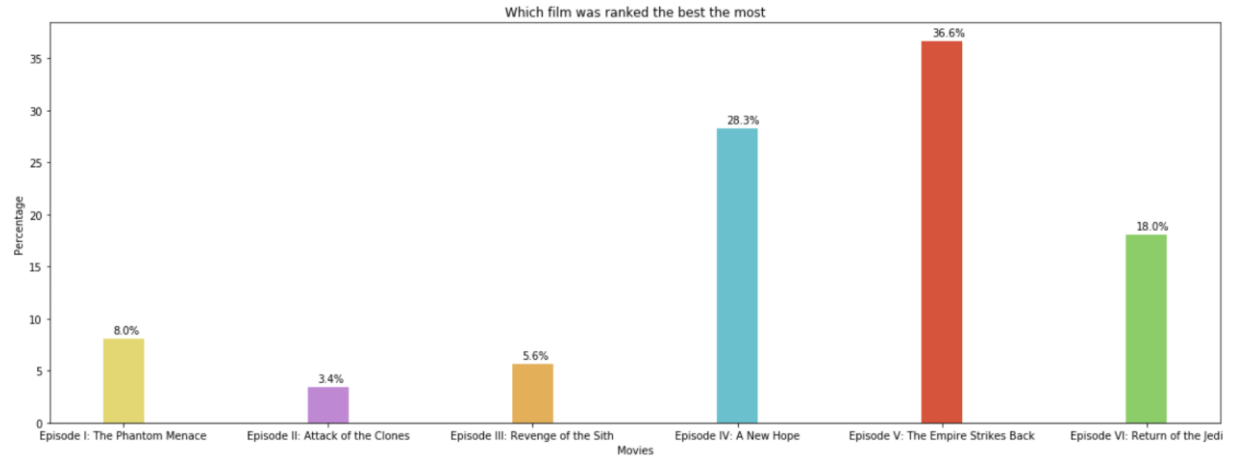
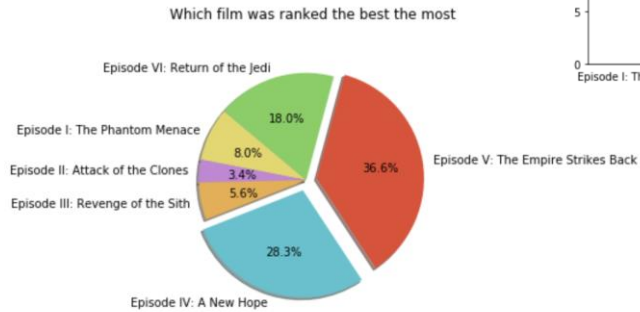
HOW PEOPLE RANK EPISODE VI: RETURN OF THE JEDI





BEST STAR WARS MOVIE?

WHICH FILM WAS RANKED THE BEST THE MOST

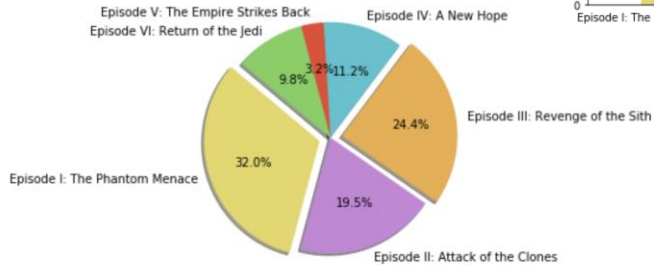




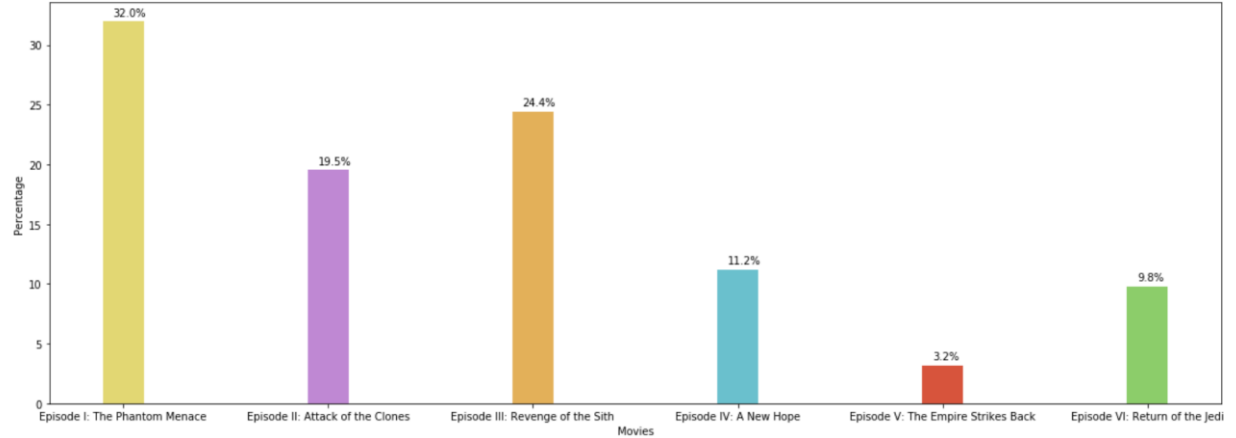
**WORST STAR WARS
MOVIE?**

WHICH FILM WAS RANKED THE WORST THE MOST

Which film was ranked the worst the most



Which film was ranked the worst the most





RELATIONSHIP BETWEEN COLUMNS



EXPLORATION PAIR 1

RELATIONSHIP BETWEEN RESPONDENT'S
AGE AND HOW THEY RANK

**EPISODE I: THE PHANTOM
MENACE**

BACKGROUND

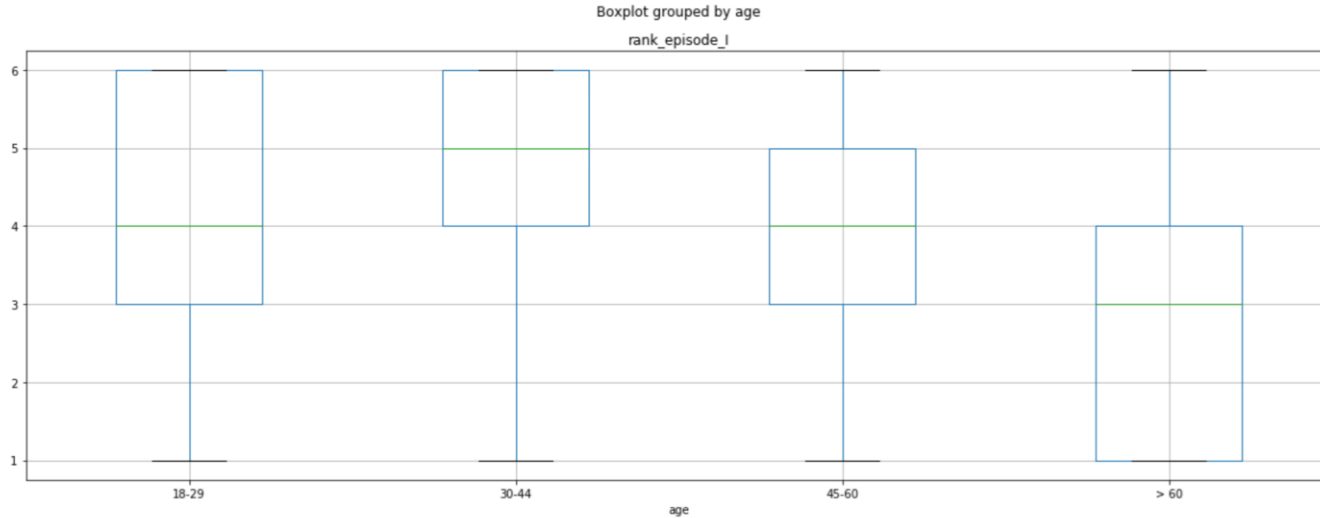
EPISODE I WAS VERY POLITICAL,
WITH MANY EXTENDED
SEQUENCES OF THE MOVIE
DEDICATED TO DIPLOMACY

HYPOTHESIS

OLDER GENERATIONS WOULD'VE UNDERSTOOD
THE POLITICAL ASPECT BETTER AND ENJOYED
THIS MOVIE BETTER THAN YOUNGER
AUDIENCES, THUS RANKING IT AS THE BEST
OVER ALL OTHERS

BOXPLOT OF THE RELATIONSHIP BETWEEN RESPONDENT'S AGE AND HOW THEY RANK

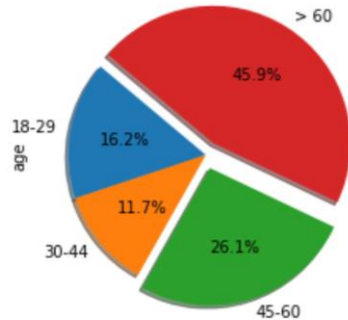
EPISODE I: THE PHANTOM MENACE



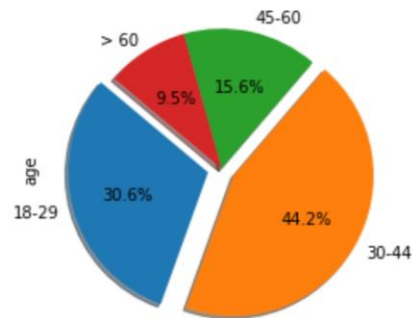
PIE CHARTS OF THE RELATIONSHIP BETWEEN RESPONDENT'S AGE AND HOW THEY RANK

EPISODE I: THE PHANTOM MENACE

Which age groups thought Episode I was the best



Which age groups thought Episode I was the worst



72%

OF RESPONDENTS WHO THOUGHT EPISODE I WAS
THE BEST ARE OVER THE AGE OF 45

74.8%

OF RESPONDENTS WHO THOUGHT EPISODE I
WAS THE WORST ARE AGED 18-44



EXPLORATION PAIR 2

RELATIONSHIP BETWEEN RESPONDENTS
WHO CONSIDER THEMSELVES FANS, AND
THEIR SUBSEQUENT ANSWER TO THE
NOTORIOUSLY DEBATED QUESTION

'WHO SHOT FIRST?'

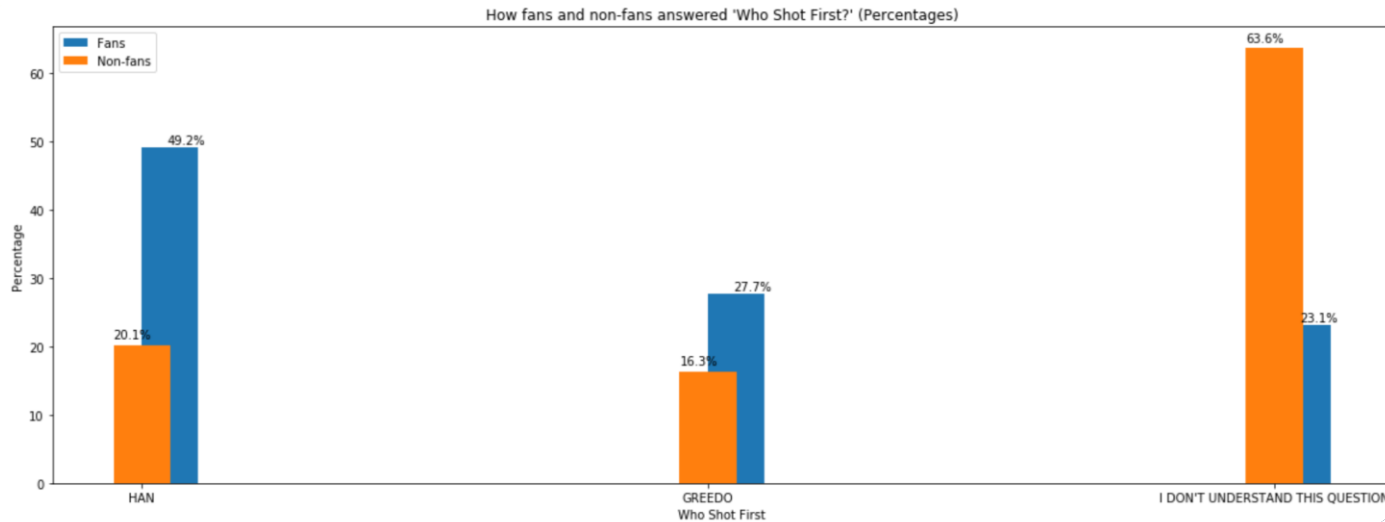
BACKGROUND

THE ORIGINAL SCENE IN EPISODE IV
HAD HAN SOLO SHOOT GREEDO FIRST,
BUT THIS WAS LATER CHANGED IN
SUBSEQUENT VERSIONS OF THE MOVIE

HYPOTHESIS

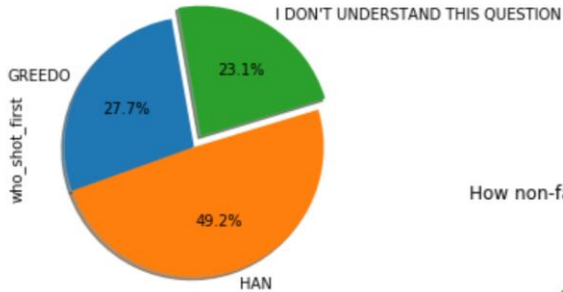
REAL FANS SHOULD KNOW THAT HAN
SHOT FIRST, OR AT THE VERY LEAST,
UNDERSTAND THE QUESTION.

Bar Graph of How respondents answered the question **'WHO SHOT FIRST?'**

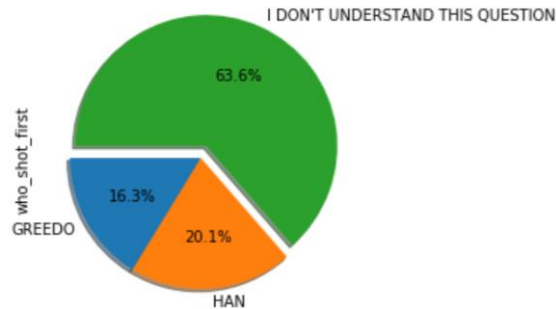


PIE CHARTS OF HOW FANS AND non-FANS answered the question 'WHO SHOT FIRST?'

How fans answered 'Who Shot First'



How non-fans answered 'Who Shot First'



49.2%
of respondents who are fans answered
'HAN'

76.9%
of respondents who are fans
understood the question

63.6%
of respondents who are not fans did not
understand the question



EXPLORATION PAIR 3

RELATIONSHIP BETWEEN RESPONDENTS
WHO ARE STAR WARS FANS AND STAR
TREK FANS

BACKGROUND

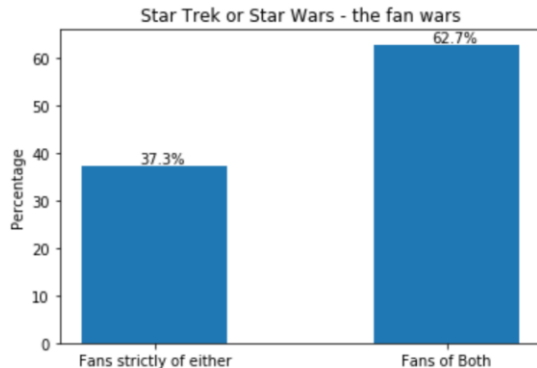
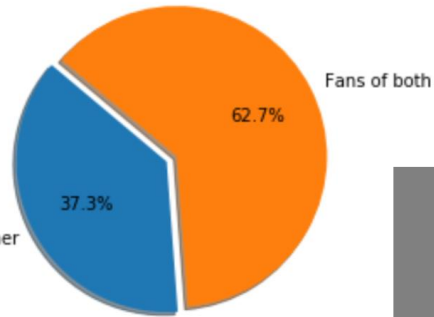
PEOPLE'S LOVE FOR BOTH THESE
PROPERTIES IS SO MUCH, THAT
THERE IS A BIT OF A RIVALRY
BETWEEN THE TWO FANDOMS.

HYPOTHESIS

DUE TO THIS INTENSE RIVALRY, IT IS
LIKELY THAT WE MAY FIND IN THE DATA
THAT PEOPLE ONLY CONSIDER
THEMSELVES FANS OF ONE OF THE
FRANCHISES, BUT NOT BOTH.

Bar Graph of How respondents answered the question **'WHO SHOT FIRST?'**

Star Trek or Star Wars - the fan wars



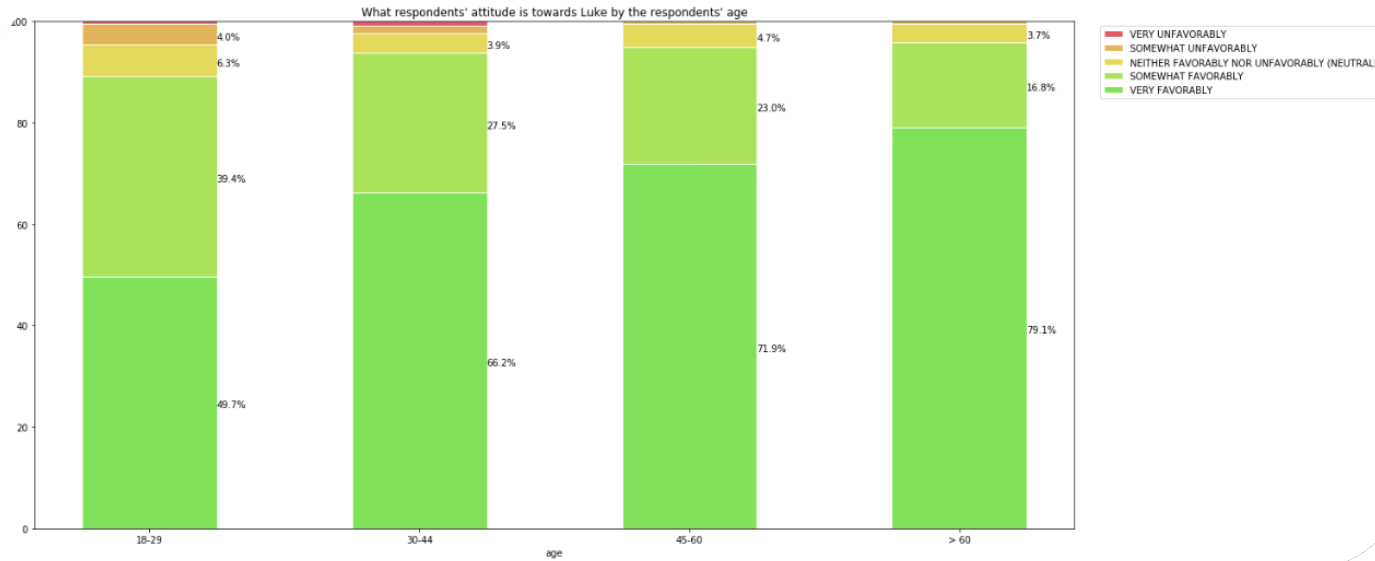
62.7%
OF RESPONDENTS ARE FANS OF BOTH
FRANCHISES

37.3%
OF RESPONDENTS ARE FANS STRICTLY OF
EITHER ONE

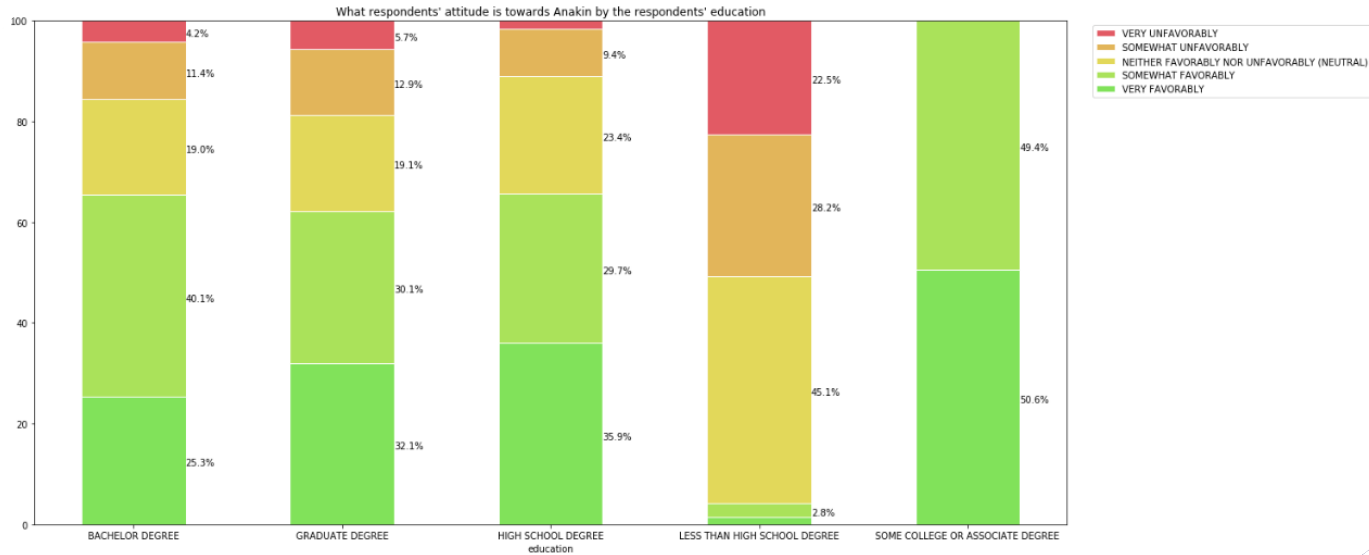


CHARACTER AFFINITIES BY DEMOGRAPHICS

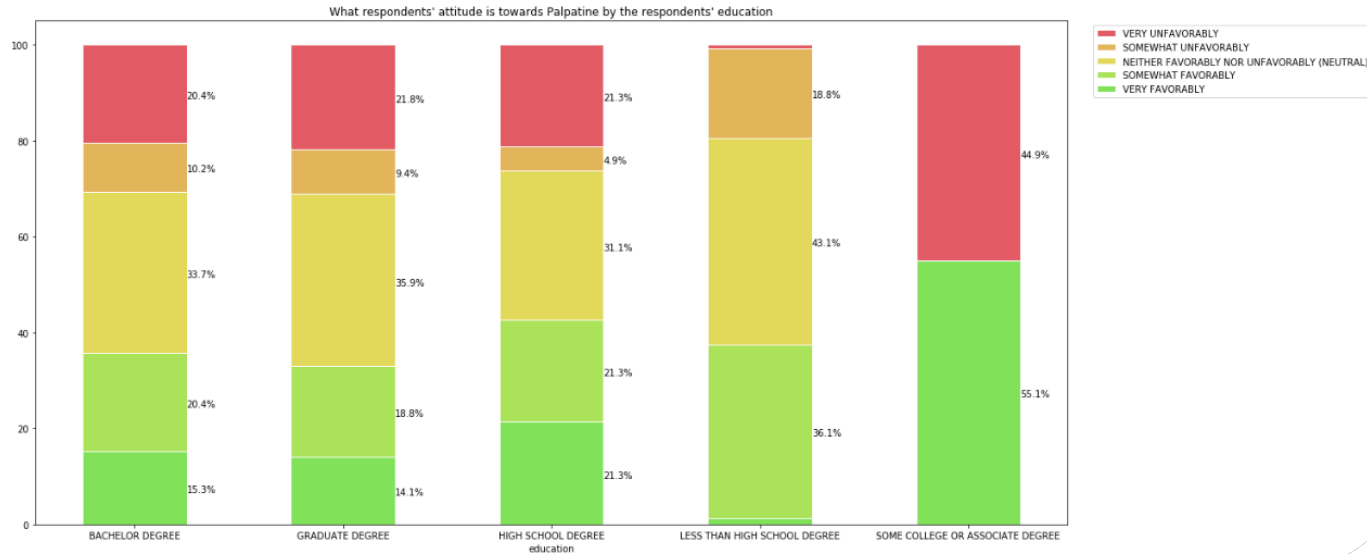
RESPONDENTS' attitude towards LUKE SKYWALKER BY age



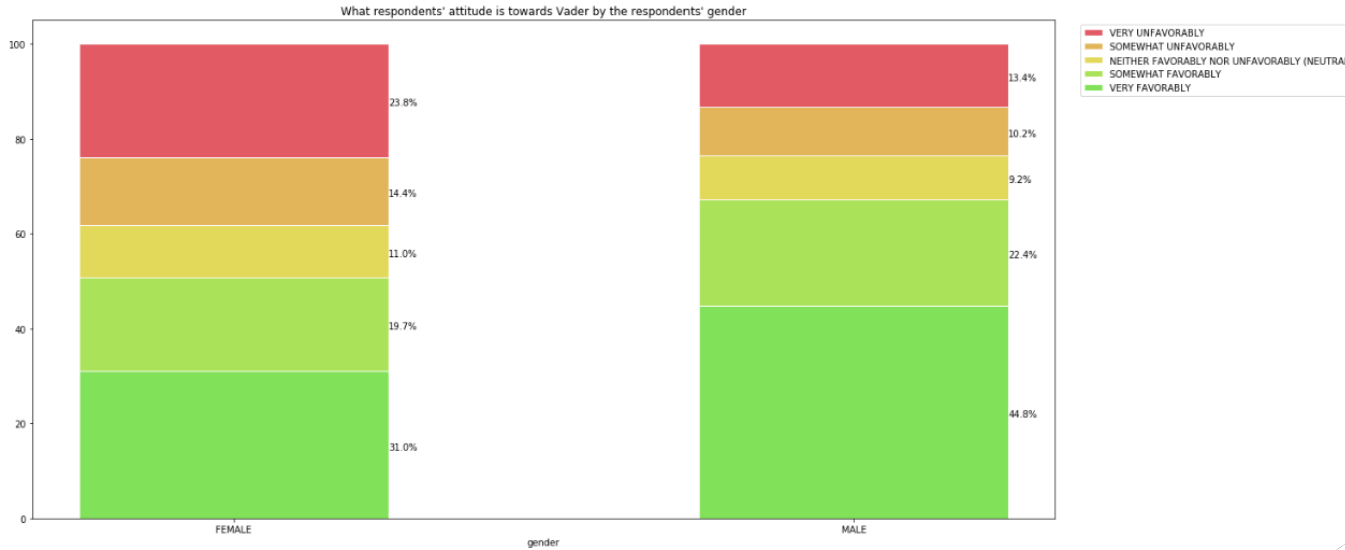
RESPONDENTS' attitude towards Anakin BY age



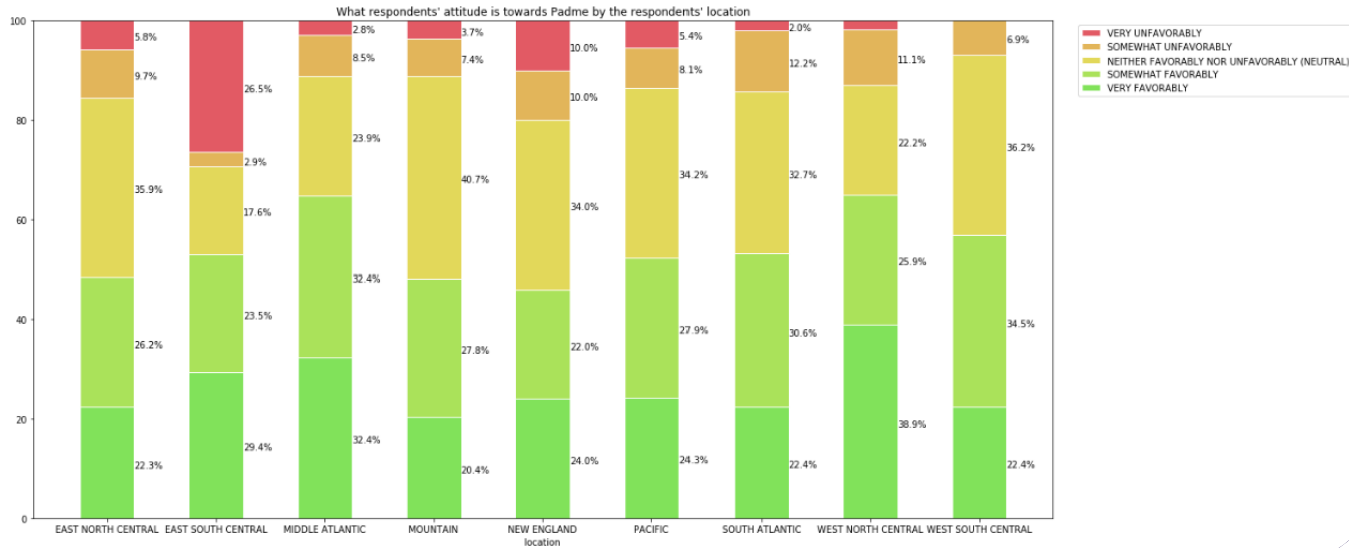
RESPONDENTS' attitude towards PALPATINE BY education



RESPONDENTS' attitude towards DARTH VADER BY GENDER



RESPONDENTS' attitude towards Padme BY LOCATION





THANKS!