

NETFLIX

Gifford Tompkins
Project 04 | DSI 09
21 October 2019

Problem Statement

As netflix expands and seeks to create more content, they are looking to speed up their evaluation process for incoming show ideas, scripts and pilots.

We will attempt to build a model that will predict if a NetFlix original show will be renewed for another season.

Data Collection

We scraped our data from the
*“List of original programs
distributed by Netflix”* page on
Wikipedia.

Title ↕	Genre ↕	Premiere ↕	Seasons ↕	Length ↕
<i>House of Cards</i>	Political drama	February 1, 2013	6 seasons, 73 episodes	42–59 min.
<i>Hemlock Grove</i>	Horror/thriller	April 19, 2013	3 seasons, 33 episodes	45–58 min.
<i>Orange Is the New Black</i>	Comedy-drama	July 11, 2013	7 seasons, 91 episodes	50–92 min.
<i>Marco Polo</i>	Historical drama	December 12, 2014	2 seasons, 20 episodes	48–65 min.
<i>Bloodline</i>	Thriller	March 20, 2015	3 seasons, 33 episodes	48–68 min.

Data Cleaning:

Data Cleaning consisted of:

- Converting 'seasons', 'length' and 'premiere' into appropriate data columns.
- Expanding 'genre' into dummy columns that included each sub-genre.
- Binarizing the 'status' column into 'Renewed', or 'Ended'
- Miscellaneous standardization of string texts.

When we finished the munging process, our data set was whittled down from 1313 to 285.

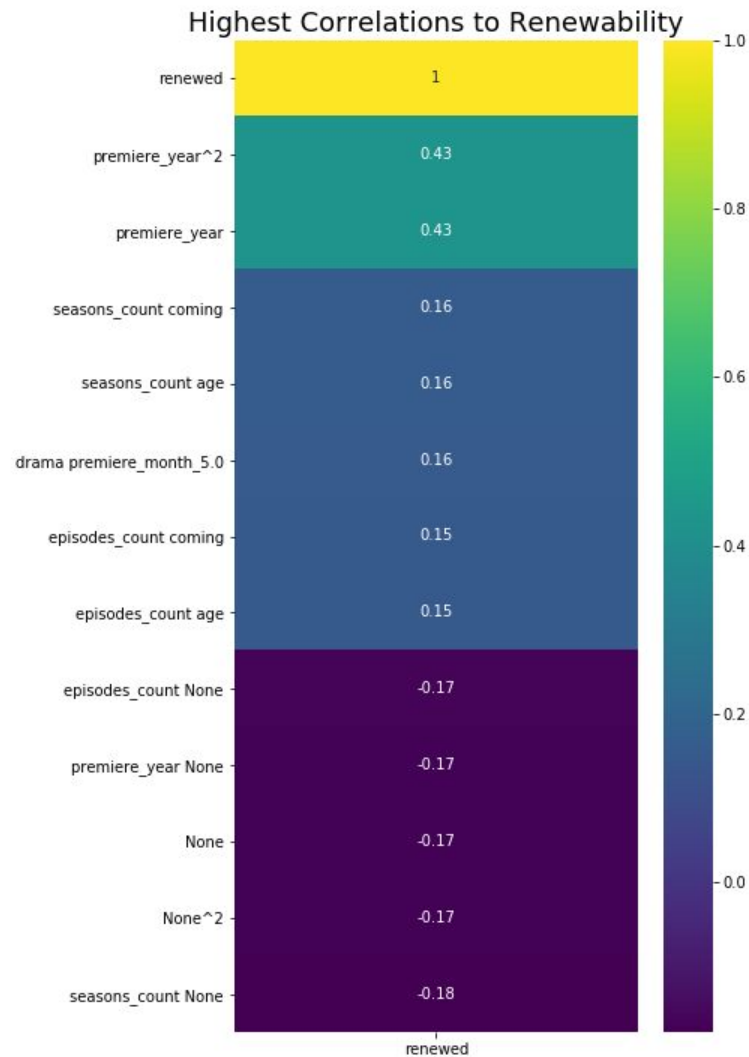
Title	Genre	Premiere	Seasons	Length
<i>House of Cards</i>	Political drama	February 1, 2013	6 seasons, 73 episodes	42–59 min.
<i>Hemlock Grove</i>	Horror/thriller	April 19, 2013	3 seasons, 33 episodes	45–58 min.
<i>Orange Is the New Black</i>	Comedy-drama	July 11, 2013	7 seasons, 91 episodes	50–92 min.
<i>Marco Polo</i>	Historical drama	December 12, 2014	2 seasons, 20 episodes	48–65 min.
<i>Bloodline</i>	Thriller	March 20, 2015	3 seasons, 33 episodes	48–68 min.

64.6 %

Baseline Score

After all of our data cleaning process, this is our goal to beat.

Correlative Columns



Decision Tree Classifier

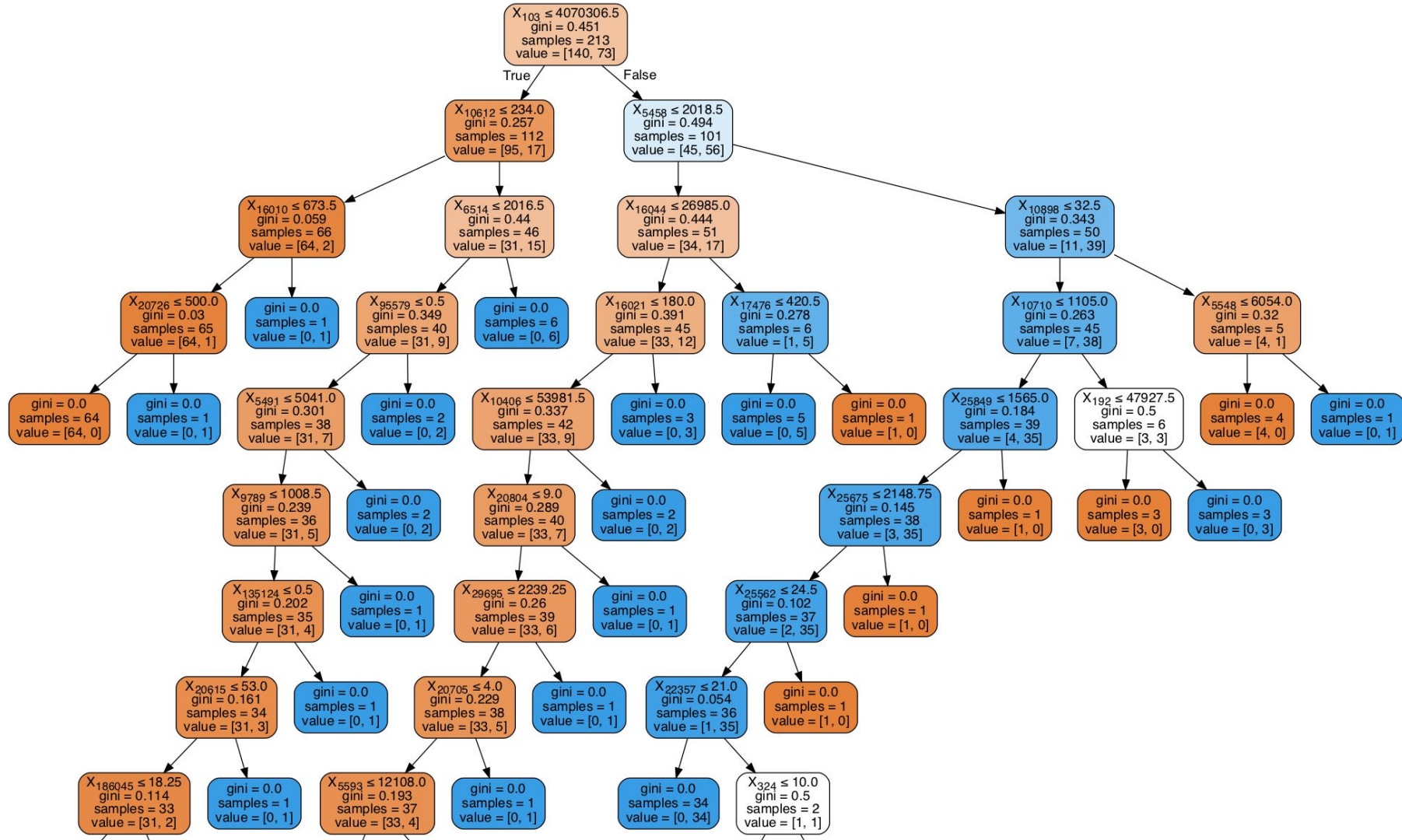
Parameters:

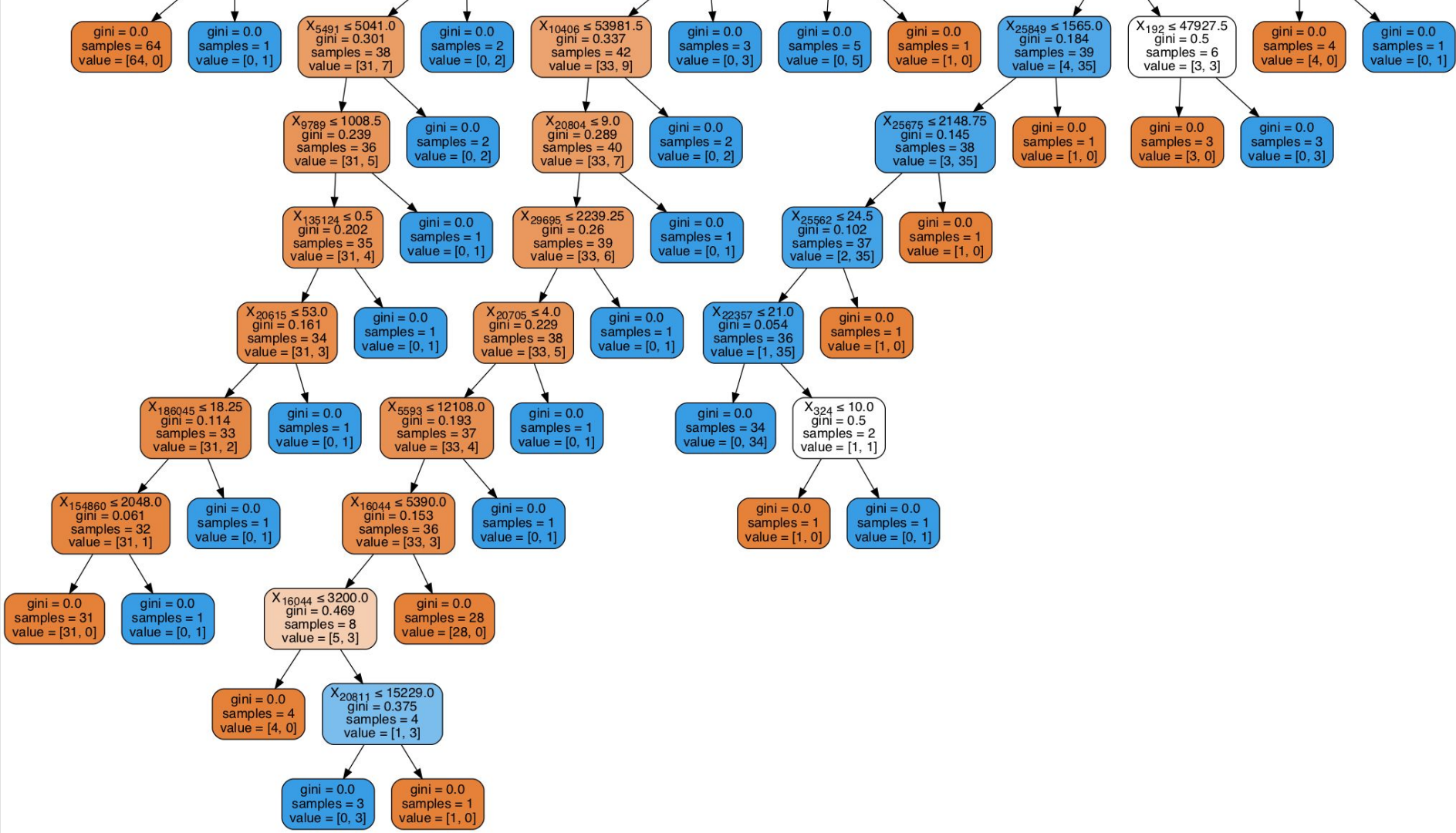
```
{'max_depth': 3,  
'min_samples_leaf': 9,  
'min_samples_split': 2,  
'random_state': 42}
```

Train score: 80%

Test Score: 77.8%

Improvement: 20.43% *from baseline*





Summary And Further Steps

There seems to be promise in being able to predict the renewability of a show in this way.

Recommendations

1. More data.
2. Play with Polynomial Features more.
3. Try a model that is a little more interpretable.