

# Do POC lead actors impact gross movie revenue?

Team 9 - Abhishek Baral, Joe Hsieh, Joao Mansur

## Abstract for Class

In the 70s Bruce Lee auditioned to be in the show Kung Fu. Set in the American Wild West, the show would feature a Shaolin monk escaping to the US to get caught in several disputes. The part was eventually awarded to David Carradine and rumors exist that Bruce Lee was passed on because “minorities can’t carry lead roles.” This report seeks to give quantitative evidence to a *casting director at a Hollywood movie studio* to provide evidence as to whether minority (non-white) leads in movies cause changes in box office revenue. Rather than try to vindicate Bruce Lee using historical data, we will use modern data to see whether a modern movie can benefit or lose from casting a minority lead. The main method used is matching minority- and white-led movies with multiple methods and showing causality through linear regression. If making a report to an actual casting director, we would exclude code but we are keeping it here for class.

## Executive Summary

Casting minority actors for lead roles in movies is beneficial for movies in the Action genre and does not cause any statistically significant loss or gain in other genres. This report proves this by scraping movie data from the IMDB and boxofficemojo.com, pruning the data by matching similar movies with POC and non-POC lead actors, then running a linear regression to find any evidence of causality. By looking at movies with similar release dates, genres, and more it was possible to isolate the effect of a lead actor’s POC status. In all cases, an actor belonging to a racial minority is not a disadvantage to a movie’s revenue and is proven to be a significant positive in Action movies.

## Scraping

Our data for this report was gathered by a process called Scraping.

Web scraping or web harvesting, is the process of using a bot or web crawler to extract content and data from a website, rather than manually copying and pasting results. Web scraping extracts the HTML code and the data stored in a database, and then can recreate the website content elsewhere. Essentially, data from websites is automatically and procedurally extracted according to specification.

For our purposes, we decided to scrape 2 websites, box office mojo and IMDB. For box office mojo, we collected four years’ worth of top domestic box office films in the United States from 2017 to 2020. The relevant data that we searched for was title, gross sales, release date, revenue, and lead actor among other variables. For IMDB, we had a similar approach, but focused on ratings and movie distributor. Lastly for the race data, we were able to find lists on IMDB of actors and their ethnicity

Once the three data sets were sourced, some level of processing was needed to clean the variables in the dataset and re-casted to an appropriate data type. Lastly, we then merged the box office mojo and imdb datasets by movie title, and then merged that newly created set with the race data set by lead actor name.

The full data results from this step is available on github at

To make sure results were accurate, we checked and cleaned the data manually. Removing any incomplete or repeating data then checking whether actor-race matches were accurate. The final dataset is then sent to analysis below.

## Initial Analysis

A linear regression can be an effective way to estimate the effect of a certain feature, like whether the lead actor is of a minority race or not. We can isolate the effect it can have on an outcome like a movie's gross revenue by using regression and accounting for other variables to make sure we isolate the effect. The way this happens is that through statistics, a linear regression can assign weights to variables to try to predict a movie's gross revenue the best it can. In these regressions, it is also possible to mathematically calculate whether we are sure an effect exists, a process summarized by the displayed p-score. If a p-score is above 0.05, we say that we aren't sure if the effect is statistically significant. If something is statistically significant, we are assuming that the resulting coefficient is not a random occurrence but can be attributed to a cause.

In other words, if we want to see if a racial minority lead actor affects gross revenue, then we can check whether the variable that accounts for POC status has a p-score of less than 0.05.

Let's begin by running a regression on the data we have scraped and cleaned. Firstly let's only use the Race variable:

Parameter	Coefficient	SE	95% CI	t(493)	p
(Intercept)	5.98e+07	4.89e+06	(5.02e+07, 6.94e+07)	12.24	< .001
Race (POC)	1.18e+07	1.06e+07	(-8.97e+06, 3.25e+07)	1.12	0.265

The print out above shows that the regression thinks that an actor being a POC, despite a negative coefficient, is not statistically significant (p-score  $[\Pr(>|t|)]$  above 0.05). The model considers a baseline (Intercept) then adds the coefficients to it; this makes the baseline a non-POC lead and the RacePOC coefficient being the calculated difference between them. Since the effect is not mathematically proven to not be random, thus having a low p-score, it seems that an actor being a POC has no significant effect if considered as the only variable in a regression.

However, we must also consider the many other variables we have, adding other variables could clarify the effect and isolate it further. Let's try a new regression with a few more variables.

Parameter	Coefficient	SE	95% CI	t(485)	p
(Intercept)	1.03e+10	9.65e+09	(-8.62e+09, 2.93e+10)	1.07	0.284
Race (POC)	9.91e+06	8.43e+06	(-6.65e+06, 2.65e+07)	1.18	0.240
year	-5.11e+06	4.78e+06	(-1.45e+07, 4.29e+06)	-1.07	0.286
genre (Adventure)	-8.39e+07	1.10e+07	(-1.05e+08, -6.23e+07)	-7.65	< .001
genre (Biography)	-9.39e+07	1.57e+07	(-1.25e+08, -6.31e+07)	-5.99	< .001
genre (Comedy)	-8.29e+07	9.89e+06	(-1.02e+08, -6.35e+07)	-8.39	< .001
genre (Documentary)	-8.11e+07	2.69e+07	(-1.34e+08, -2.83e+07)	-3.01	0.003
genre (Drama)	-8.89e+07	1.21e+07	(-1.13e+08, -6.51e+07)	-7.35	< .001
Male	1.23e+07	7.34e+06	(-2.10e+06, 2.67e+07)	1.68	0.094
opening_theaters	23745.85	3951.40	(15981.88, 31509.83)	6.01	< .001

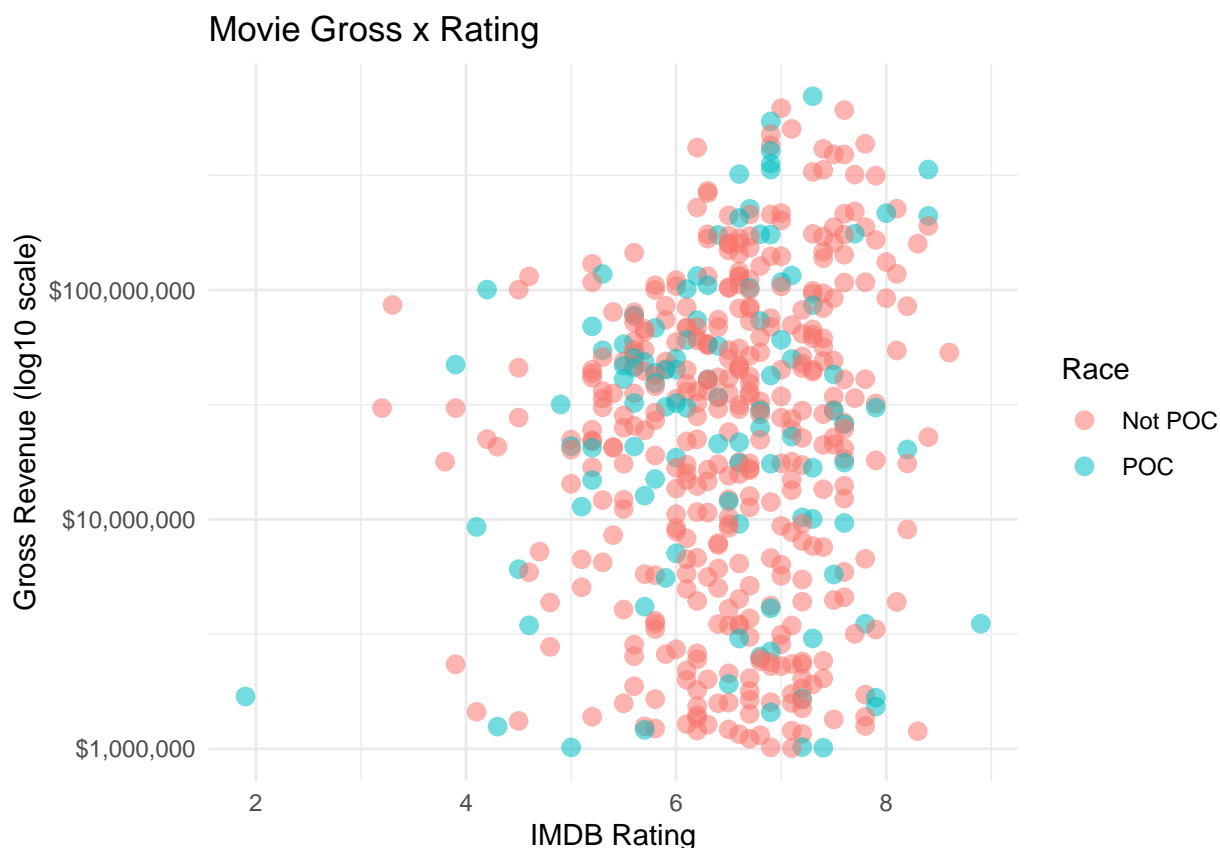
What we're looking at here is a baseline or intercept of White female lead actors in Action movies in the year 2017 with no opening theaters (an illogical but mathematically fine assumption since we correct it later). The coefficients modify the expected gross revenue of a movie depending on their characteristics. What we see here is that there is a clearly significant difference between Action movies and other genres, where Action movies tend to make more money. This is also true for male-led movies, where we see a larger gross revenue effect. We also included the number of opening theaters for the movie, which should relate to how large the movie's production is, and we see that movies launched on more theaters make more revenue.

In the end though, whether the lead is a POC or what year it was made does not seem to show any statistically significant difference.

However, the data isn't well-balanced between POC and non-POC lead actors which we can see below:

Race	Count
Not POC	389
POC	106

POC-led movies are only 106 movies or around 21.4% of the data. If we plot these movies using their IMDB Ratings and their Gross Revenue, we can see that non-POC movies occupy more space.



It seems like a lot of the data is mixed and there isn't a clear separation between POC and not POC-led movies. If we prune the data to include movies that are similar to each other but only distinguishable by the race of their lead actor, we can more closely isolate the effect of a lead actor being a minority has on gross revenue.

With the data in mind, we can begin to prepare our matching strategy.

## Matching

To perform matching, we must determine how movies are to be matched. We have a couple of obvious choices: the year, movie genre, film rating, and whether the lead is male. These variables are all fixed categories (they are not fractional and a movie can't occupy two at the same time). This means that we can match films based on whether they meet all criteria the same way except for whether their lead is a POC. This kind of matching is called exact matching, and it allows for the creation of groups of similar movies and the removal of ones that don't fit with another.

Some data is more granular and difficult to match exactly, like IMDB ratings that are fractional and the amount of theaters a movie opened in. These will require some strategy to allow exact matches to occur, a practice we call binning.

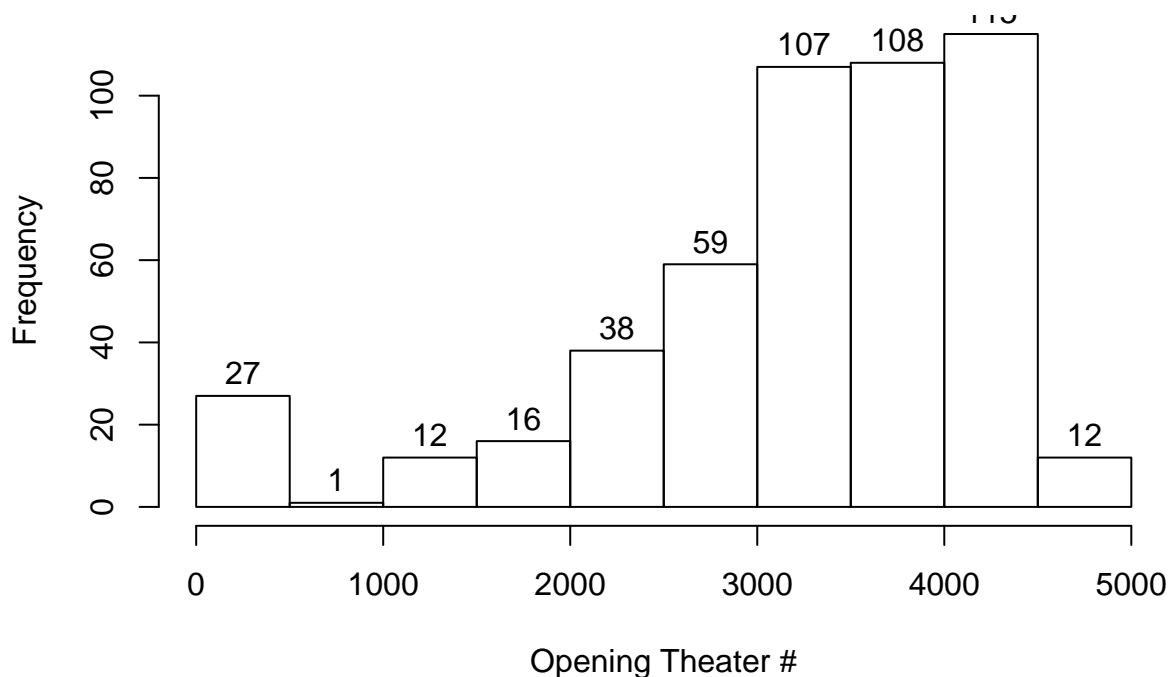
## Opening Theaters

How many theaters the movies opened in is a number with a lot of variability and finding exact matches would be difficult.

However, the number of theaters a movie opened in can be really indicative of the amount of promotion that went into the movie. More indie films are clearly released in less theaters than blockbusters.

To use this feature, which we believe to be important as mentioned, we need to bin them by separating movies into groups. A histogram can help us find these bins:

### Histogram of Movies' Opening Theaters



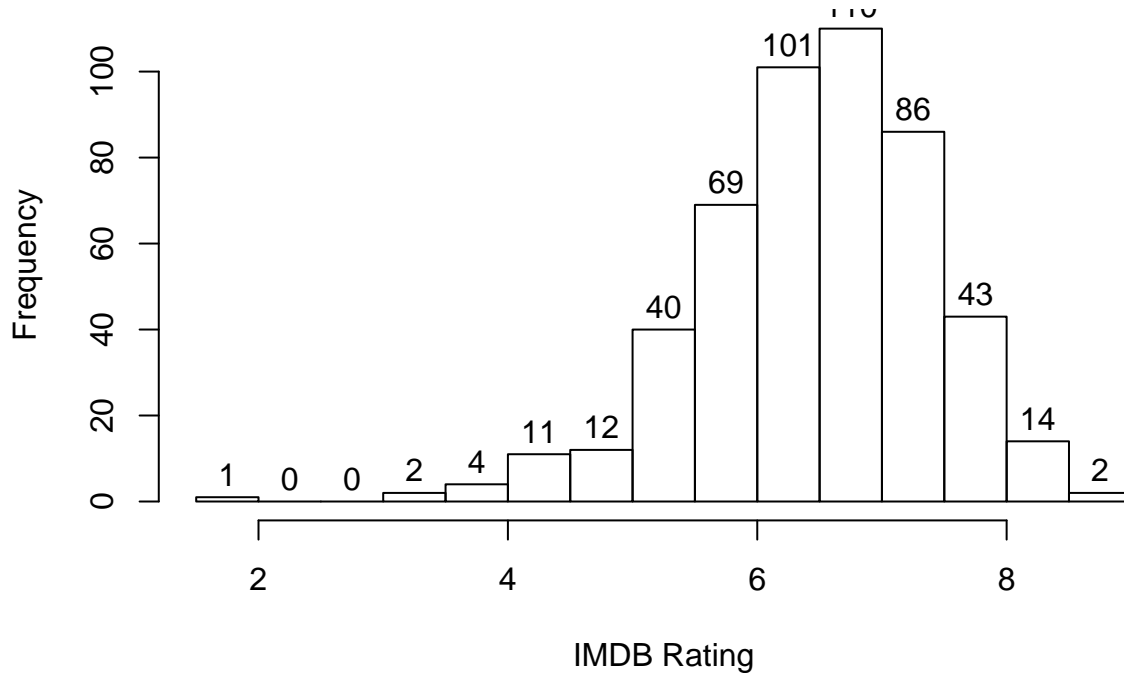
Using a histogram we can see that we can group movies by their opening theaters. Here, the choice for the size of the bin matters where there is a trade-off between how informative a bin is and how easy it will be to find matches. We decided to use movies with under 2000, 2500, 3000 then so on in 250-sized increments. These values allow for at least 40 or so movies in each bin which allows for good matching chances without losing the predictive power of the variable.

Opening Theaters	Count
< 2000	56
< 2500	38
< 3000	59
< 3250	61
< 3500	46
< 3750	62
< 4000	44
< 4250	73
>4250	56

## IMDB Ratings

We can do the same for their IMDB ratings. These ratings hopefully will account for the quality of the movie enough to be a proxy for a movie's long-term box office success.

### Histogram of Movies IMDB Ratings



IMDB ratings seem very close to a normal distribution around 6.5 or so. This means that if we round to whole numbers, most values will be 6 or 7 and the rest will be in tails of 8+ and less than 5. To make effective use of these, we can separate bins within one standard deviation from those beyond, thus achieving bins related to bad, medium, and great movies. One standard deviation in the IMDB Ratings is .94, meaning that movies within around 5.5 and 7.5 are within a standard deviation. To bin these movies, we can round their IMDB ratings to the nearest whole number, bin 6's and 7's, and then bin those above and beyond this range.

IMDB Rating	Count
5-	59
6/7	363
8+	73

Using these four simple bins should be enough to isolate the effect of extraordinarily good and bad movies from the middle of the pack. 8+ movies are one standard deviation higher in rating and 5- are one lower.

## Exact Matching

We can now finally start matching movies, where we can create subclasses of movies that are in all the same bins but with a difference in whether their lead actor is a POC.

To do this, we use an R Package called MatchIt. There are several ways to match data, but the method we have been referring to is called Exact matching. If there is even a slight difference then a match won't be found, this is what made binning the variables so important. MatchIt creates subclasses that contain movies with the same binned variables and drops any that don't have both movies with POC and movies

with non-POC lead actors. This means that we are not matching in pairs, but rather creating groups of movies and removing any outliers.

Let's run MatchIt and see how our data set changes.

```
library(MatchIt)

## Warning: package 'MatchIt' was built under R version 3.6.2

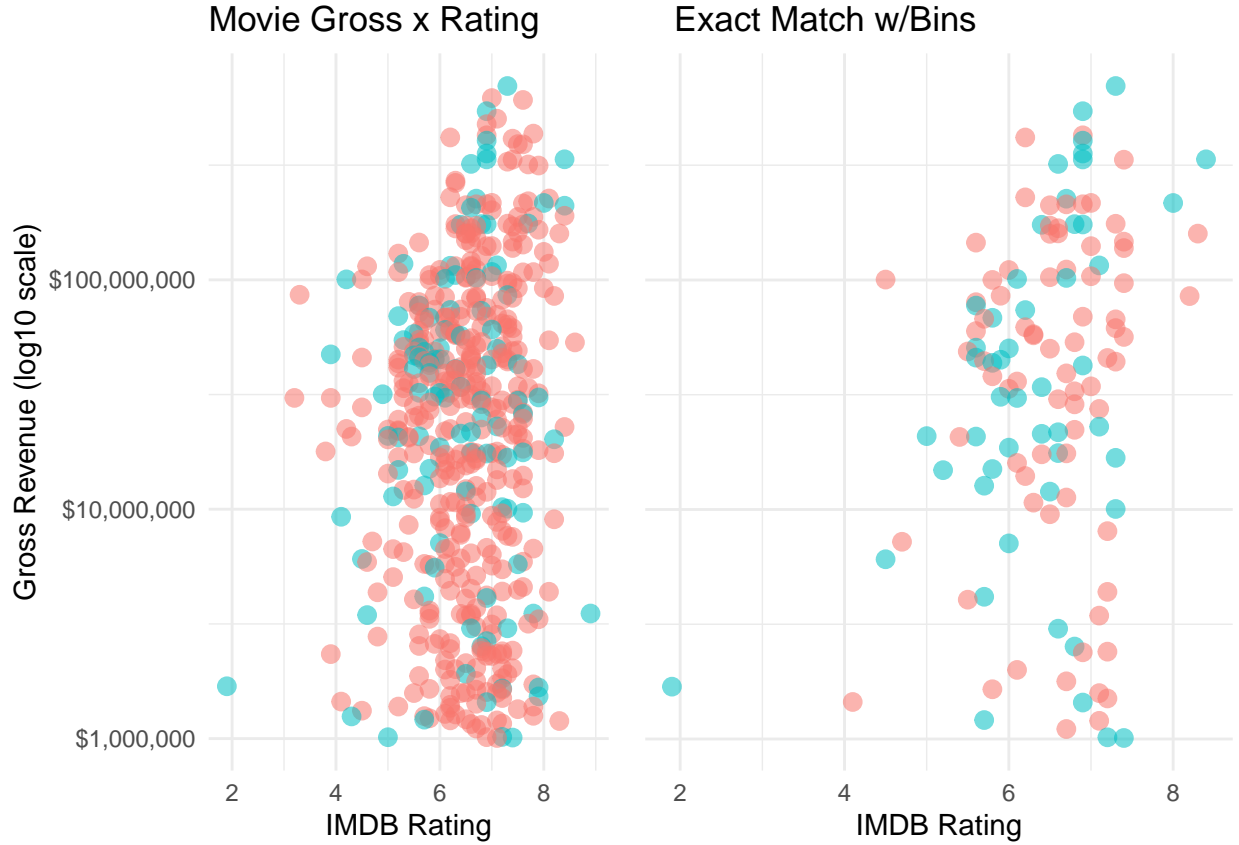
data$treat <- 0
data[data$Race == "POC", "treat"] <- 1
m.out <- matchit(treat ~ year + genre + OPT + film_rating + Male + IMDBbin, data = data,
  method = "exact", verbose = TRUE)

## Exact matching...
## Calculating matching weights... Done.
m.data <- match.data(m.out)

## [1] 495
## [1] 126
```

Race	Count
Not POC	76
POC	50

Our original data had 495 movies pre-match, after matching we have 126 movies. While POC-led movies are still around 39.6% of the data, we have eliminated movies that don't share similarities and potentially increased the visibility of the effect we are looking for. We can visualize the dataset before and after matching in the same graph:



It is clear how the data set is pruned and more balanced between the turquoise POC data and the reddish non-POC data. We obviously did not balance on Gross Revenue, which is our target variable to predict, but we have maintained the normal distribution of IMDB ratings which is evident by the same overall shape of both graphs.

We can now once again run our regression but this time with matched data. Note that while we matched using bins, we can go back to using un-binned variables in our regression. We will also utilize the weights from our matching. These weights are meant to give power to certain movie data that is representative of the overall population. This method means that although we pruned a lot of movies, certain movies can be given additional weight so that they represent information from movies that were lost.

With that, let's run the regression:

Parameter	Coefficient	SE	95% CI	t(112)	p
(Intercept)	-1.87e+10	2.49e+10	(-6.81e+10, 3.07e+10)	-0.75	0.454
Race (POC)	3.82e+07	1.60e+07	(6.56e+06, 6.99e+07)	2.39	0.018
year	9.21e+06	1.24e+07	(-1.53e+07, 3.37e+07)	0.75	0.458
genre (Adventure)	-4.60e+07	3.25e+07	(-1.10e+08, 1.84e+07)	-1.42	0.160
genre (Biography)	-6.81e+07	6.36e+07	(-1.94e+08, 5.78e+07)	-1.07	0.286
genre (Comedy)	-7.32e+07	2.40e+07	(-1.21e+08, -2.57e+07)	-3.05	0.003
genre (Documentary)	-7.33e+07	8.18e+07	(-2.35e+08, 8.88e+07)	-0.90	0.372
genre (Drama)	-6.08e+07	3.38e+07	(-1.28e+08, 6.11e+06)	-1.80	0.074
opening_theaters	30559.40	11912.65	(6956.02, 54162.78)	2.57	0.012
film_rating (PG)	1.79e+07	5.82e+07	(-9.74e+07, 1.33e+08)	0.31	0.759
film_rating (PG-13)	-1.41e+07	4.73e+07	(-1.08e+08, 7.97e+07)	-0.30	0.766
film_rating (R)	-3.78e+07	4.78e+07	(-1.32e+08, 5.69e+07)	-0.79	0.431
Male	2.72e+07	2.02e+07	(-1.28e+07, 6.72e+07)	1.35	0.180

Parameter	Coefficient	SE	95% CI	t(112)	p
imdb_rating	2.41e+07	9.64e+06	(4.98e+06, 4.32e+07)	2.50	0.014

After matching, we now see that a movie’s lead actor being a POC has not only a statistically significant effect, but a positive one. As a reminder, our baseline is the gross revenue expected from White female lead actors in Action movies in the year 2017 with no opening theaters. This means that our matched data, which keeps movies that are more standard and similar to the usual release according to our data, shows that POC-led movies tend to make more money.

However, we can still create another variable to gain more insight into movies led by POC. Since a lot of movie genres are deeply impacted by race and culture, it could be a good idea to check what happens if we allow the regression to work with the interaction between race and genre. What this means is that the regression will create a variable for movies that at the same time are led by a POC actor and belong to a particular genre. These interaction terms will allow us to see the effect of having a POC lead on movies of each genre.

That regression will come to this result:

Parameter	Coefficient	SE	95% CI	t(107)	p
(Intercept)	-1.94e+10	2.49e+10	(-6.87e+10, 2.99e+10)	-0.78	0.436
Race (POC)	7.87e+07	2.34e+07	(3.24e+07, 1.25e+08)	3.37	0.001
year	9.54e+06	1.23e+07	(-1.49e+07, 3.40e+07)	0.77	0.440
genre (Adventure)	-1.56e+07	3.81e+07	(-9.11e+07, 5.99e+07)	-0.41	0.684
genre (Biography)	-3.46e+07	7.89e+07	(-1.91e+08, 1.22e+08)	-0.44	0.661
genre (Comedy)	-4.41e+07	2.98e+07	(-1.03e+08, 1.49e+07)	-1.48	0.142
genre (Documentary)	-5.01e+07	9.12e+07	(-2.31e+08, 1.31e+08)	-0.55	0.584
genre (Drama)	-2.60e+07	3.85e+07	(-1.02e+08, 5.02e+07)	-0.68	0.500
opening_theaters	31673.82	12090.61	(7705.60, 55642.04)	2.62	0.010
film_rating (PG)	1.98e+07	5.81e+07	(-9.53e+07, 1.35e+08)	0.34	0.734
film_rating (PG-13)	-1.23e+07	4.72e+07	(-1.06e+08, 8.12e+07)	-0.26	0.794
film_rating (R)	-3.57e+07	4.77e+07	(-1.30e+08, 5.88e+07)	-0.75	0.456
Male	2.82e+07	2.01e+07	(-1.17e+07, 6.80e+07)	1.40	0.164
imdb_rating	2.19e+07	9.73e+06	(2.58e+06, 4.12e+07)	2.25	0.027
Race (POC) * genre (Adventure)	-7.28e+07	5.14e+07	(-1.75e+08, 2.90e+07)	-1.42	0.159
Race (POC) * genre (Biography)	-8.04e+07	1.15e+08	(-3.08e+08, 1.47e+08)	-0.70	0.485
Race (POC) * genre (Comedy)	-7.20e+07	4.45e+07	(-1.60e+08, 1.61e+07)	-1.62	0.108
Race (POC) * genre (Documentary)	-4.27e+07	1.16e+08	(-2.73e+08, 1.87e+08)	-0.37	0.714
Race (POC) * genre (Drama)	-8.39e+07	4.25e+07	(-1.68e+08, 4.46e+05)	-1.97	0.051

The regression printout above now isolates the effect of POC actors in each genre, with the original RacePOC being the effect on Action movies. We still see a positive effect from POC lead actors but now with an even larger estimate and better p-value (thus more statistically significant). This is unsurprising given the success of action movies like Black Panther (Rest in peace King Boesman), Joker (Joaquin Phoenix of Puerto Rican descent) and others led by Dwayne “The Rock” Johnson, Jason Momoa, and Will Smith. However, in other genres we see that the effect of POC-leads is not statistically significant.

We can safely assume that modern Action movies can benefit from minority leads while in other genres there is no proof that minority leads have an effect on box office performance. Either way, there is no evidence that minority actors can’t carry movies as lead actors; if anything, they might be a great choice for action movies right now.



## Optimal Matching

Though we decided to use Exact matching, we can test this data set by using another form of matching. The MatchIt package also offers optimal matching, which uses Mahalanobis Distances to find similar movies. These distance metrics allow for non-exact differences by mathematically finding the distance between two movies' characteristics, and matching those closest together but different by whether the lead is a POC. MatchIt's optimal matching system solves for pairs of points while minimizing the total distance between them, thus finding the optimal and most similar movie pairs. Exact matching, on the other hand, created subsets instead of pairs and finding pairs can create better comparisons for our target effect.

However, the fact that it minimizes the distance between points means that it relies on data a lot more. While in binning for exact matches we grouped movies following a logic (bad, medium, good movies and binning opening theater numbers to find approximations of promotion size), optimal matching minimizes distance between the metrics making small differences in features more important. This method makes certain matches viable, like matches at the end of one bin and the beginning of another, but it also overvalues relationships within bins where it may pair movies with extremely similar ratings though ratings are subjective and not informative past a certain threshold. The idea of using distances therefore has issues because it assumes that the distances between points is relevant, something we believed to not be true to this dataset.

By using exact matching, we utilized non-categorical data in a way that we tried to best preserve their utility with domain knowledge. Optimal matching, however, can produce similar or different results but it will unarguably provide pairs that are mathematically similar. Though our conclusion will rely mostly on exact matching, we believe optimal matching can provide an interesting take and likely the same overall result.

Let's try running optimal matching to see what happens.

```
# Note: optimal matching naturally attempts 1:1 matching but can be modified to  
# do 1:n
```

```
m.out2 <- matchit(treat ~ year + genre + opening_theaters + film_rating + Male +  
  imdb_rating, data = data, method = "optimal", verbose = TRUE)
```

```
## Optimal matching...  
## Calculating matching weights... Done.
```

```
m.data2 <- match.data(m.out2)
```

```
# original  
paste("Original Data rows:", nrow(data))
```

```
## [1] "Original Data rows: 495"  
paste("Exact-matched Data rows:", nrow(m.data))
```

```
## [1] "Exact-matched Data rows: 126"  
paste("Optimal-matched Data rows:", nrow(m.data2))
```

```
## [1] "Optimal-matched Data rows: 212"
```

Optimal matching results in more rows than exact matching, this makes sense since exact matching is stricter in the pairs it can allow. Optimal data also balances perfectly between our treated variable (POC actors):

Race	Count
Not POC	106
POC	106

Let's look at the differences visually.



Between Exact and Optimal there is an increase in amount of points while the overall shape is still maintained. Notably, non-POC movies matched seem to be more clustered at the bottom with two exceptions. Whether this clustering of lower-revenue non-POC movies will have an effect can be something we check with our regression.

Parameter	Coefficient	SE	95% CI	t(192)	p
(Intercept)	-4.24e+09	1.73e+10	(-3.84e+10, 2.99e+10)	-0.25	0.807
Race (POC)	5.50e+07	1.86e+07	(1.82e+07, 9.18e+07)	2.95	0.004
year	2.02e+06	8.58e+06	(-1.49e+07, 1.89e+07)	0.24	0.814
genre (Adventure)	-4.48e+07	2.56e+07	(-9.53e+07, 5.62e+06)	-1.75	0.081
genre (Biography)	-4.54e+07	3.41e+07	(-1.13e+08, 2.18e+07)	-1.33	0.185
genre (Comedy)	-4.90e+07	2.32e+07	(-9.48e+07, -3.22e+06)	-2.11	0.036
genre (Documentary)	-1.30e+07	6.50e+07	(-1.41e+08, 1.15e+08)	-0.20	0.842
genre (Drama)	-4.76e+07	2.54e+07	(-9.77e+07, 2.49e+06)	-1.87	0.062
opening_theaters	21770.96	6817.00	(8325.13, 35216.79)	3.19	0.002
film_rating (Not Rated)	4.96e+07	9.46e+07	(-1.37e+08, 2.36e+08)	0.52	0.601
film_rating (PG)	8.75e+07	9.28e+07	(-9.56e+07, 2.71e+08)	0.94	0.347
film_rating (PG-13)	7.59e+07	9.17e+07	(-1.05e+08, 2.57e+08)	0.83	0.409
film_rating (R)	4.57e+07	9.15e+07	(-1.35e+08, 2.26e+08)	0.50	0.618
Male	8.81e+06	1.28e+07	(-1.65e+07, 3.41e+07)	0.69	0.493
imdb_rating	1.62e+07	6.10e+06	(4.13e+06, 2.82e+07)	2.65	0.009
Race (POC) * genre (Adventure)	-5.41e+07	3.38e+07	(-1.21e+08, 1.26e+07)	-1.60	0.111
Race (POC) * genre (Biography)	-5.99e+07	5.12e+07	(-1.61e+08, 4.10e+07)	-1.17	0.243
Race (POC) * genre (Comedy)	-4.06e+07	3.22e+07	(-1.04e+08, 2.30e+07)	-1.26	0.210
Race (POC) * genre (Documentary)	-4.02e+07	7.87e+07	(-1.95e+08, 1.15e+08)	-0.51	0.610

Parameter	Coefficient	SE	95% CI	t(192)	p
Race (POC) * genre (Drama)	-4.61e+07	3.33e+07	(-1.12e+08, 1.96e+07)	-1.38	0.168

Once again, we see the same significant variables in the regression. Action movies with POC leads seem to generate more gross revenue while there is no statistically significant difference for the other genres.

## Conclusion

When evaluating casting for a movie, there is significant evidence to prove that there is no negative causal effect from casting a racial minority in the lead role. In fact, action movies may benefit from minority leads which follows the success of several blockbuster action movies led by POC. Through the use of matching, we managed to isolate the effect of POC leads on movies by only considering similar movie groups or pairs. Our models have R-squared values between .3 and .5 showing that at least half of the variation of data is not captured, thus we should not think that these are the only factors that determine a movie's box office success. However, we have mathematically and scientifically found substantial evidence that, even outside of morality, there is no reason to discriminate against casting POC to lead movie roles.

From Team 9, we hope this has been a useful look into data and that this may fuel more representative behavior in the industry.