

Met4 Home Exam

Movies: Blockbusters and Flops

Dept. of Business and Management Science

15.11.2017 09:00 - 17.11.2017 14:00

1 Introduction

In this case, we are interested in costs and revenues from movies. How are costs and revenues related? Can we say something about the expected return from movies? And what about risk - do some genres have a higher variability of returns than others?

*The Movie Database*¹ (TMDB) is a large online database of movies, containing a lot of information about each movie. Movies are classified by genres, and each movie may be tagged in several genres. Further, users on TMDB may rate movies on a scale from 1 to 10, where a higher value indicates a more positive view on the movie.

In this home exam, we will use a subset of the data from TMDB. Firstly, we will only consider movies released in 1990 or later years, and secondly only movies where the first genre is either “*Action*”, “*Adventure*”, “*Comedy*”, “*Drama*” or “*Horror*”. See the appendix for an overview of the variables in the dataset.

Although this is not the primary focus of this home exam, keep in mind the data quality may be low, and therefore that there might be some errors in the dataset. Please see TMDB if you need any additional information on the variables. The dataset is distributed in two files: `data.csv` and `data.RData`. These files are identical except for the formatting. Please use the file most suited for your statistical software.

¹<https://www.themoviedb.org/>

2 Assignments

1. Present and *briefly* discuss relevant statistics of the data set.
2. Use OLS to regress $\log(\text{revenue})$ on $\log(\text{budget})$, and answer the following questions:
 - (a) Briefly discuss your results.
 - (b) Are there any exceptional outliers, and if so - which ones? Are some genres more likely to be outliers?²
 - (c) Are your OLS-estimates sensitive to the presence of the outliers?
3. A risk averse investor is considering buying shares in a movie that has already been made, but not yet been released to the public. The return on investment for the investor is given by:

$$r = \frac{.5 * \text{revenue} - \text{budget}}{\text{budget}} \quad (1)$$

where `revenue` and `budget` are variables in your dataset. The movie in question has a budget of \$10 000 000, and in a test screening viewers gave the movie on average 7.5 points with a similar measure as the variable `vote_average`. What is the expected return, and what is the variance of the investment?³ What additional information could potentially improve your estimates?

4. The dataset contains more variables and may provide insights beyond the first three assignments. Use maximum one page to present and discuss one additional finding from the dataset.

²Note: see appendix B for some hints!

³Note: see appendix C for some hints!

A Variables in the dataset

The following table lists all the variables in the data set:

Variable name	Definition
budget	Budgeted production costs of movie in million USD
revenue	Box office revenues in million USD
original_language	Original language of movie
original_title	Title of movie in original language
runtime	Runtime of movie in minutes
title	Movie title in English
vote_average	Average vote by users on TMDB
vote_count	Number of user votes for a movie on TMDB
month	Month movie was released
year	Year movie was released
genre	Main genre of movie
prod_comp	Main production company

B Largest entries of a variable

B.1 With R

With R, if `df` is a data frame, the command:

```
large.values <- tail(order(abs(df$Variable), na.last = F), n=10)
```

creates a collection in `large.values` with the *indices* of the 10 largest, absolute values of `df$Variable`, ignoring missing values. Hence, issuing the following command

```
df[large.values,]
```

returns the rows of `df` with the ten largest absolute values of `Variable`.

B.2 With Gretl (software used in earlier versions of Met4)

The following commands can be used in Gretl, if we want to sort the dataset by the absolute values of `Variable`:

```
absVariable = abs(Variable)
```

```
dataset sortby absVariable
```

```
smpl missing(absVariable)==0 --restrict
```

The data set is now sorted by the absolute values of `Variable`, ignoring missing values.

C Variance of predictions with OLS

C.1 With R

If we estimate a model in R with OLS, e.g.

```
reg <- lm(Y ~ X, data=data.df)
```

and have a data frame `data.new` where we want to make predictions:

```
pred <- predict.lm(reg, data.new, se.fit=TRUE)
```

We can use `pred` to find the variance of the *expected value of Y*:

```
var.conf <- pred$se.fit^2
```

If we want the variance of a predicted, single value of *Y*, we can find this with:

```
var.pred <- var.conf+pred$residual_scale^2
```

C.2 With Gretl (software used in earlier versions of Met4)

Assuming you have subset the data set to the observations you want to use for running the regression,⁴ you run OLS with the command:

```
ols Y 0 X
```

We first run the following command to activate the whole dataset. This allows us to create predictions not only to the observations used for estimating the model:

```
smp1 full
```

Thereafter, you generate predictions with

```
fcast --mean-y
```

The third column of the output displays the standard error of the predicted, expected value of *Y*, and squaring it gives the variance. To obtain the variance of the predicted value of a single observation of *Y*, we use

```
fcast
```

where again we square the third column to get the variance.

⁴E.g. the command `smp1 1 10` subsets to the first 10 observations

Regulations for the home exam

The home exam in Met4 must be handed in by groups of sizes 2, 3 or 4 students. See in particular section 9 the REGULATIONS FOR EXAMINATIONS AT NHH, as well as part 2 of SUPPLEMENTARY REGULATIONS TO THE REGULATIONS FOR EXAMINATIONS (FULL-TIME PROGRAMMES). It is not permitted to discuss the exam with students not in your group after the dataset has been released.

Grading and formal requirements

The reports will be graded in accordance with the assessment rubric posted on ItsLearning as well as in table 1. You may write your report in Norwegian or English.

If there is a need for clarifications, you may send an email to Ole-Petter Moe Hansen (s9705@nhh.no) *and* Håkon Otneim (Hakon.Otneim@nhh.no). Any extra information will be announced to all groups on It'sLearning.

The report should be maximum 10 pages. Tables figures and references are included in the ten pages. If the report has a front page without any answers to the assignments, this may come in addition to the 10 pages. Prioritize what you include in the report!

The report must typeset with Times New Roman, with font size 12 and line spacing 1.15. Text in tables and figures may be set as low as font size 9.

The exam is handed in through Wiseflow. We accept a single file of type pdf for submission - MS Word files or r-scripts are not accepted.

Met4 Fall 2017: Rubric for grading home exam

Points:	3	2	1	0	Weight
Presentation of tables and figures	Tables and figures i: are self-explanatory, ii: give a meaningful contribution to the results of the report, and iii: are visually appealing. The report demonstrates independence in the choice of graphics.	Tables and figures i: are self-explanatory, ii: give meaningful contribution to the results of the report, and iii: are visually appealing.	Tables and figures only partially i: are self-explanatory, ii: give meaningful contribution to the results of the report, and iii: are visually appealing.	Tables and figures are difficult to understand, and do not contribute to the results of the report.	0.10
Choice of methods	The methods chosen are appropriate for the problem at hand. The report explains which choices and assumptions are used. The report demonstrates great degree of independence in the choice of methods.	The methods chosen are appropriate for the problem at hand. The report explains which choices and assumptions are used.	The methods may be appropriate, but the justification for choice of methods is lacking.	The methods chosen are not appropriate for the problem at hand. Assumptions and choices are not explained. The report displays lack of independence in the choice of methods.	0.30
Application of methods	The methods are skillfully implemented, without any errors.	The methods are skillfully implemented, with only minor errors .	There are some errors in the implementation, but the the results are overall correct.	There are several errors made in the application of the methods.	0.25
Discussion of findings	Demonstates an outstanding understanding of both the statistical and economical significance of results. Results are discussed in a relevant context, e.g. policy prescriptions, discussion of causality, contrast to theoretical results, as applicable.	Demonstrates an good understanding of both the statistical and economical significance of results. Results are discussed in a relevant context, e.g. policy prescriptions, discussion of causality, contrast to theoretical results, as applicable.	Demonstates an some understanding of both the statistical and economical significance of results. Results are to a limited extent discussed in a relevant context.	Lists statistical results without demonstrating an understanding of their meaning. Does not highlight the economic importance of the results.	0.30
Ethics	Discusses the ethical implications of the analysis (if applicable).			Ignores the ethical implications of the analysis (if applicable).	0.05
<p>Each divergence from the formal requirements of the report (e.g. page lenght, font..) gives -1 points from the overall score.</p>					

Table 1: We calculate the the overall score for a submission by summing together the points for each row, using the weights in the rightmost column. The cutpoints between grades will be set after the exam.

Solution Proposal

Assignment 1

Table 2 shows summary statistics for the numerical variables in the dataset. Other descriptive statistics, such as histograms or tables with frequencies over categorical variables, could also be relevant.

Table 2: Summary statistics over the numerical variables in the data set.

Statistic	N	Mean	St. Dev.	Min	Max
budget	2,488	39.792	43.517	0.008	380.000
revenue	2,230	118.037	186.196	0.002	2,787.965
runtime	3,216	106.877	20.037	0	276
vote_average	3,217	6.047	1.071	0.000	10.000
vote_count	3,217	691.994	1,265.847	0	13,752
month	3,217			1	12
year	3,217			1990	2017

Assignment 2

Table 3 shows the baseline regression result in column (1), showing that there is almost a one-for-one relationship between $\log(\text{budget})$ and $\log(\text{revenue})$. However, there is considerable heteroskedasticity, as shown in figure C.2.

The ten largest outliers from column (1) in table 3 are shown in table 4. As we see, outliers are present in all categories, but only horror movies have revenues far over budgets.

Finally, we can note that dropping these outliers do not influence the results in the initial regression all that much, as shown in column (2) in table 3.

Table 3: Assignment 2 regressions. Column (1) is on the full dataset, column (2) excludes the ten largest outliers from column (1).

	<i>Dependent variable:</i>	
	log(revenue)	
	(1)	(2)
log(Budget)	0.932*** (0.024)	0.949*** (0.022)
Constant	0.786*** (0.081)	0.748*** (0.077)
Observations	2,125	2,115
R ²	0.425	0.459
Adjusted R ²	0.424	0.458
Residual Std. Error	1.396 (df = 2123)	1.312 (df = 2113)
F Statistic	1,566.366*** (df = 1; 2123)	1,789.393*** (df = 1; 2113)

Note:

*p<0.1; **p<0.05; ***p<0.01

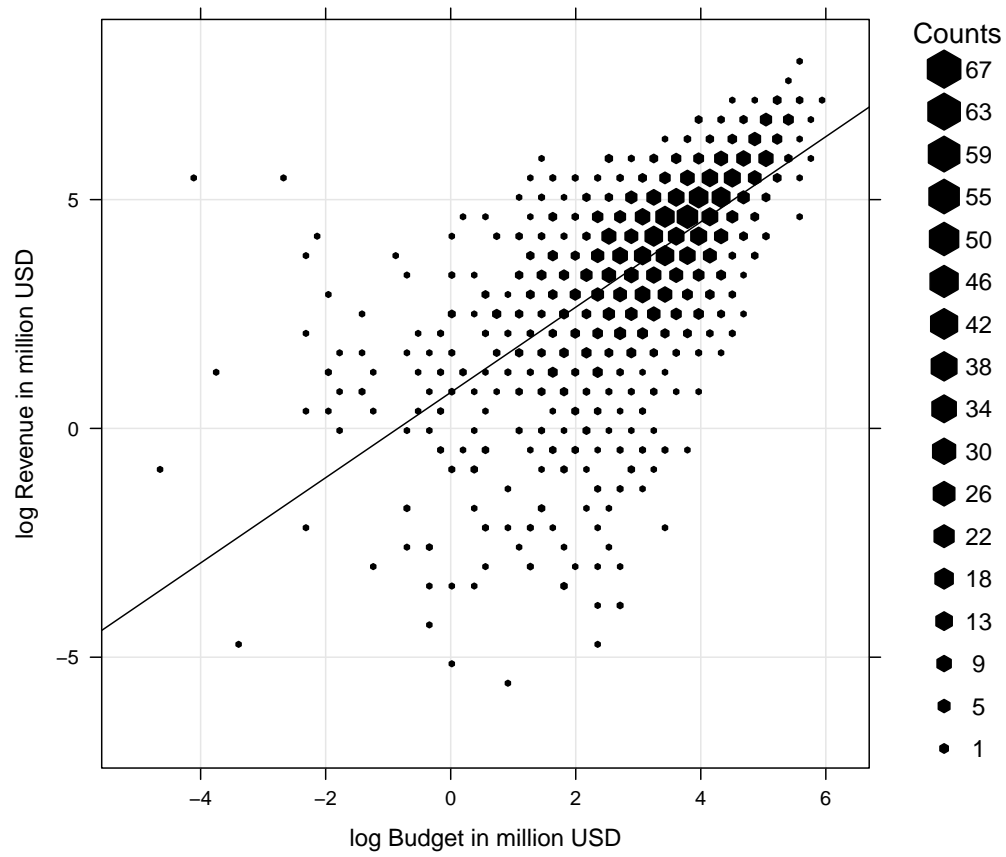


Figure 1: Scatterplot over $\log(\text{budget})$ and $\log(\text{revenue})$, with the data pooled in bins (the “counts” show how many observations there are in each bin). The straight line is the regression line, coming from column(1) in table 3.

Table 4: The ten largest outliers from column (1) in table 3

Budget, mill \$	Revenue Mill \$	Title	Genre
15	0.023	All The Queen's Men	Action
30	0.108	1911	Adventure
10	0.046	An Alan Smithee Film: Burn, Hollywood, Burn	Comedy
15	0.020	The Good Night	Comedy
10	0.009	Janky Promoters	Comedy
14	0.046	Margaret	Drama
10	0.017	Strangerland	Drama
2.100	0.003	Mi America	Drama
0.060	248	The Blair Witch Project	Horror
0.015	193.356	Paranormal Activity	Horror

Assignment 3

There are several ways to answer this assignment. One way of solving it is to estimate the relationship between r and the two variables where we have data for this movie, and use the regression to predict the expected return and variance for the movie in question.

Table 5 shows the results, both in total and by genre. We can note that horror movies are very different than the other categories, and that R^2 increase by a factor of around 3-8 when we estimate the model by genres separately. Hence, knowing the genre increases the fit of the model considerably.

Finally, table 6 shows expected returns and variance for the movie in question, where we can see that both statistics vary considerable across the categories. We don't know how the investor weighs expected returns against variance, so we cannot make strong claims about which categories it is worthwhile investing in. However, Drama movies have both a lower expected return and higher variance than Comedies, so a Comedy should be preferred over a Drama.

Horror movies have very large variance, so the investor would have to be close to risk neutral to want to invest here. However, a good answer will note that the extreme results of horror movies are driven a few movies with small budgets and very large revenues. Hence, it may be very unlikely that those extreme returns are relevant for a movie with a budget of 10 million. Dropping movies with a budget

smaller than e.g. 1 million makes the results on horror movies look a lot more like the other genres.

Table 5: Regression results on return by genre

	<i>Dependent variable:</i>					
	$r = \frac{.5 * \text{revenue} - \text{budget}}{\text{budget}}$					
	All	Action	Adventure	Comedy	Drama	Horror
	(1)	(2)	(3)	(4)	(5)	(6)
vote_average	-1.467 (3.661)	0.470*** (0.051)	0.643*** (0.112)	0.977*** (0.227)	0.468 (0.406)	2.350 (57.019)
log(budget)	-18.694*** (2.446)	0.022 (0.044)	-0.101 (0.103)	-1.365*** (0.156)	-1.994*** (0.257)	-186.343*** (28.605)
Constant	73.914*** (24.457)	-2.742*** (0.341)	-3.099*** (0.799)	-0.780 (1.514)	3.662 (2.775)	491.918 (343.292)
Observations	2,125	521	233	577	644	150
R ²	0.027	0.142	0.125	0.159	0.087	0.225

Note:

*p<0.1; **p<0.05; ***p<0.01

Table 6: Expected return and variance for the movie with 7.5 vote average and budget of 10 mill \$, total and by genre.

	All	Action	Adventure	Comedy	Drama	Horror
Expected Return	19.87	0.83	1.49	3.40	2.58	80.48
Variance	21,069.78	0.98	1.91	17.30	78.17	249,930.50

Assignment 4

This is an opportunity for students to show of their skills. However, given that there is much work to be done in the first three assignments, the requirements here are not very high.


```

44 reg <- lm(lrev~lbud , data=df)
45
46 # Next, want to identify the largest outliers. First, calculate
47 # residuals for each observation:
48 df$residuals <- df$lrev-predict.lm(reg , df)
49
50 # Next, find a list with the ten largest outliers. Use absolute
51 # deviations from the regression line:
52 outlier.list <- tail(order(abs(df$residuals) , na.last = F) ,n=10)
53
54 # Create a small dataframe with the largest outliers
55 df.outliers <- df[outlier.list , c(1,2,6,11)]
56 df.outliers <- df.outliers[order(df.outliers$genre) ,]
57
58 # Print list of the largest outliers:
59 stargazer(df.outliers , summary=F,rownames = F, type="text")
60
61 # Re-estimate the regression without the 10 outliers:
62 reg.no.outliers <- lm(lrev~lbud , data=df ,
63                       subset=-c(outlier.list))
64
65 # Compare with and without outliers:
66 stargazer(reg , reg.no.outliers , type="text")
67
68 # We can also show this with a plot:
69 pdf("plot.pdf")
70 hexbinplot(df$lrev~df$lbud ,
71            inv=function(x) x ,
72            type=c('g','r') ,
73            xlab="log Budget in million USD" ,
74            ylab="log Revenue in million USD" ,
75            style="lattice")
76 dev.off()
77
78
79 ## Assignment 3 – Expected return _____
80
81 # First, create the variable with returns:
82 df$r <- (df$revenue*.5 - df$budget)/(df$budget)
83
84 # Next, regress returns on vote average and log(budget).
85 # Do this for the entire dataset, as well as for each
86 # genre separately:
87 eq <- formula(r~vote_average+lbud)
88 reg <- lm(eq, data=df )

```

```

89 reg1 <- lm(eq, data=df, subset=genre=="Action" )
90 reg2 <- lm(eq, data=df, subset=genre=="Adventure" )
91 reg3 <- lm(eq, data=df, subset=genre=="Comedy" )
92 reg4 <- lm(eq, data=df, subset=genre=="Drama" )
93 reg5 <- lm(eq, data=df, subset=genre=="Horror" )
94
95 # Compare the results:
96 stargazer(reg ,
97           reg1 ,
98           reg2 ,
99           reg3 ,
100          reg4 ,
101          reg5 ,
102          type="text",
103          column.labels = c("All",
104                            "Action",
105                            "Adventure",
106                            "Comedy",
107                            "Drama",
108                            "Horror")
109          )
110
111 # Next, want to predict expected utility. Create a data frame
112 # with a single observation, with values equal to
113 # the movie the investor is considering:
114 predframe <- data.frame(vote_average=7.5, lbud=log(10))
115
116 # create a prediction list for each of the regressions:
117 p <- predict.lm(reg , predframe , se.fit = T)
118 p1 <- predict.lm(reg1 , predframe , se.fit = T)
119 p2 <- predict.lm(reg2 , predframe , se.fit = T)
120 p3 <- predict.lm(reg3 , predframe , se.fit = T)
121 p4 <- predict.lm(reg4 , predframe , se.fit = T)
122 p5 <- predict.lm(reg5 , predframe , se.fit = T)
123
124 # Finally, we calculate expected return and variance in
125 # each of the cases:
126 exp.ret <- matrix(NA, ncol=6,nrow=2)
127 colnames(exp.ret) <- c("All",
128                       "Action",
129                       "Adventure",
130                       "Comedy",
131                       "Drama",
132                       "Horror")
133 rownames(exp.ret) <- c("Expected Return",

```

```

134                                     "Variance")
135 exp.ret[1,1] <- p $fit
136 exp.ret[1,2] <- p1$fit
137 exp.ret[1,3] <- p2$fit
138 exp.ret[1,4] <- p3$fit
139 exp.ret[1,5] <- p4$fit
140 exp.ret[1,6] <- p5$fit
141 exp.ret[2,1] <- p $se.fit^2+p $residual.scale^2
142 exp.ret[2,2] <- p1$se.fit^2+p1$residual.scale^2
143 exp.ret[2,3] <- p2$se.fit^2+p2$residual.scale^2
144 exp.ret[2,4] <- p3$se.fit^2+p3$residual.scale^2
145 exp.ret[2,5] <- p4$se.fit^2+p4$residual.scale^2
146 exp.ret[2,6] <- p5$se.fit^2+p5$residual.scale^2
147
148 # Print out expected returns. As we see, the investor
149 # has highest expected returns from comedies. Horrors movies
150 # have an extreme variance.
151 stargazer(exp.ret, type="text", digits = 2)
152
153
154 ## The End

```

solution_proposal.R