

Data Literacy

University of Tübingen, Winter Term 2020/21

Exercise Sheet 5

© 2020 Prof. Dr. Philipp Hennig & Marius Hobbahn

This sheet is due on Tuesday 10 December 2020 at 12 noon sharp (i.e. before the start of the lecture).

This week's exercise sheet is inspired by a recently viral "Who is talking in popular films" post which you can find [here](#).

We want to investigate whether the average user rating of films on [IMDb](#) is related to how much of the dialogue is done by male vs female characters. This pop-culture data is a less serious variant of a [data analysis](#) problem that was [hotly debated](#) in [Germany](#) earlier [this year](#) in the context of the pandemic (the connection will be made clearer in lectures 6 and 7).

The goal of the exercise is to gain an intuition for how to inspect and analyse data without necessarily following a strict statistical recipe. In fact we will see that pre-packaged stats can be dangerous or silly if applied without context, and that sometimes a less formal but richer visual analysis can be more useful.

title

```
In [1]: import matplotlib.pyplot as plt
import matplotlib
import numpy as np
import pandas as pd

# Make inline plots vector graphics
%matplotlib inline
from IPython.display import set_matplotlib_formats

set_matplotlib_formats("pdf", "svg")

matplotlib.rcParams["font", **{"family": "serif", "serif": ["Computer Modern"]}]
plt.rcParams["text.usetex"] = True
plt.rcParams["text.latex.preamble"] = r"\usepackage{amsfonts} \usepackage{amsmath}"
```

Part I: Data preparation

For this exercise we do the entire data wrangling for you such that you can exclusively focus on the statistical analysis. The data was extracted from the [original tableau datafile](#) by us. Everything you need is given in `df_combined` which contains meta data about movies.

```
In [2]: # import all .csv files
df_combined = pd.read_csv("df_combined.csv")
```

```
In [33]: print("length of df_combined: ", len(df_combined))
df_combined.head()

length of df_combined: 1996
```

```
Out[33]:
```

Unnamed: 0	script_id	words	words_male	words_female	f_percentage	imdb_id	title	year	gross	lines_data	averageRating
0	0	280	6394	2631	3763	0.588520	tt0112579	The Bridges of Madison County	1995	142.0	4.332023e+72
1	1	623	9108	7584	1524	0.167325	tt0179626	15 Minutes	2001	37.0	7.777778e+127
2	2	625	4401	4246	155	0.035219	tt0062622	2001: A Space Odyssey	1968	376.0	7.777734e+52
3	3	630	10132	9059	1073	0.105902	tt0307901	25th Hour	2002	19.0	7.777775e+102
4	4	633	9029	8163	866	0.095913	tt1019452	A Serious Man	2009	10.0	1.456768e+87

Part II: Understanding the data

We begin by visualizing the data in a scatter plot. We recreate the "Who is talking in popular films" scatter plot.

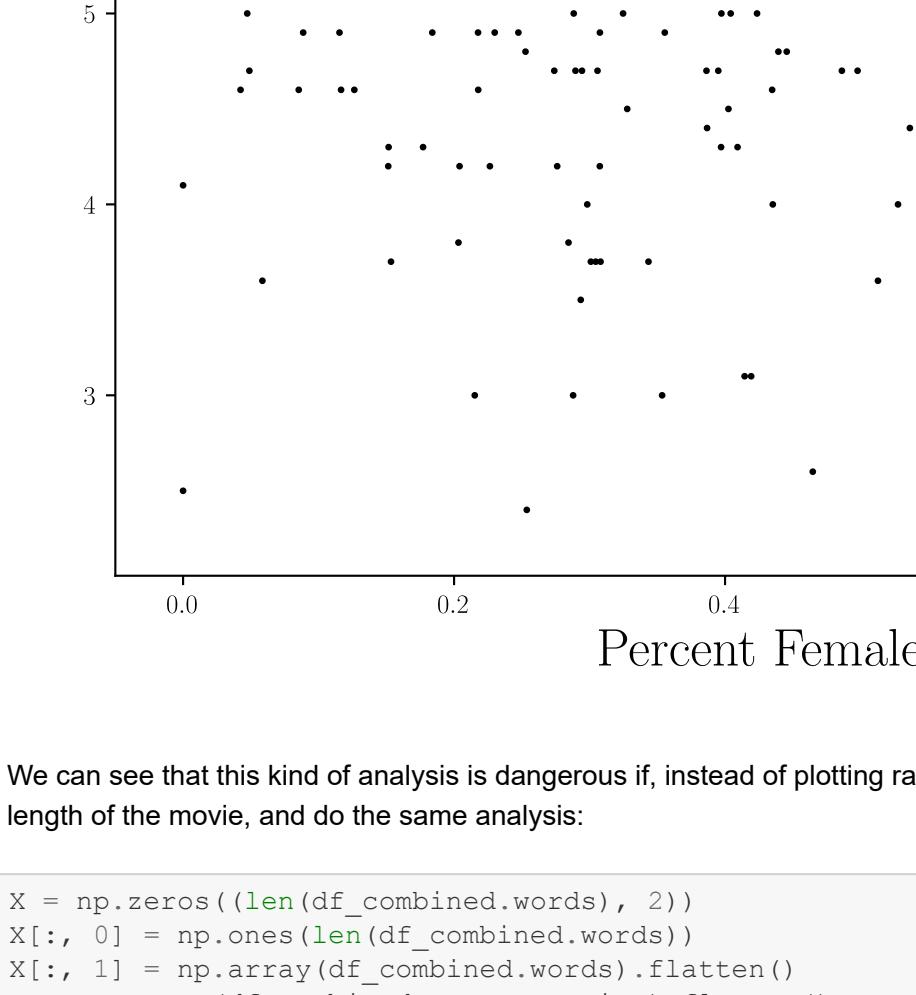
The necessary components of that plot are `percentFemaleDialogue` on the x-axis, `userRating` on the y-axis and two different color schemes for majority male or majority female films. Everything else, such as showing some titles, adapting the size or the marginal plot above is optional. You are also not supposed to filter for films above 200k ratings.

```
In [32]: # recreating the "Who is talking in popular films" scatter plot
majority_speaker = df_combined["f_percentage"] > 0.5
import matplotlib.patches as mpatches

yellow_patch = mpatches.Patch(color='yellow', label='female')
purple_patch = mpatches.Patch(color='purple', label='male')

plt.scatter(df_combined.f_percentage, df_combined.averageRating, s=1.5, c= majority_speaker.astype(int))
plt.xlabel("Mean Movie Ratings by Percentage of male and female main speakers")
plt.ylabel("Movie Rating")
plt.legend(handles=[yellow_patch, purple_patch])
plt.show()
```

Mean Movie Ratings by Percentage of male and female main speakers



Part III: Pre-packaged Statistics

The tempting next step is to do some `tests`. For example, we may want to apply binary tests for contingency tables, as introduced in lecture 04, to test whether there is a bias against films with majority female dialogue. Doing so is actually *not* a good idea here (more below), but we will try it anyway.

1. Compute the mean of the average rating of all films. This is just one number - not a `DataFrame`.
2. Create a contingency table that has larger/smaller than 50% female dialogue on the x-axis and larger/smaller than the mean average rating on the y-axis.
3. Choose either the binomial test or Fisher's exact test and justify your choice.

--> we choose the Fisher's test because we want to test a relation of two variables. The Binomial test is for testing hypothesis which we don't have at this point.

1. Compute a *p*-value for your test of choice.

```
In [37]: # recreating the "Who is talking in popular films" scatter plot
# add the average rating
mean_rating = df_combined["averageRating"].mean()
print("Mean of all movie ratings")
print(mean_rating)
# TODO: this is optional but very helpful to understand what we are doing
```

Mean of all movie ratings
6.8044589178356825

```
In [57]: # determine the totals on the edge
femalespeaker_above = df_combined[df_combined["f_percentage"] > 0.5].count().title
femalespeaker_below = df_combined[df_combined["f_percentage"] <= 0.5].count().title
averageRating_above = df_combined[df_combined["averageRating"] > mean_rating].count().title
averageRating_below = df_combined[df_combined["averageRating"] <= mean_rating].count().title
# determine all 4 boxes in the contingency table
femalespeaker_aboveRating = df_combined[(df_combined["f_percentage"] > 0.5) & (df_combined["averageRating"] > mean_rating)].count().title
femalespeaker_belowRating = df_combined[(df_combined["f_percentage"] <= 0.5) & (df_combined["averageRating"] < mean_rating)].count().title
malespeaker_aboveRating = df_combined[(df_combined["f_percentage"] <= 0.5) & (df_combined["averageRating"] > mean_rating)].count().title
malespeaker_belowRating = df_combined[(df_combined["f_percentage"] > 0.5) & (df_combined["averageRating"] <= mean_rating)].count().title

contingency_table = pd.DataFrame([["above average", "female", "male", "total"],
["above average", femalespeaker_aboveRating, malespeaker_aboveRating, averageRating_above],
["below average", femalespeaker_belowRating, malespeaker_belowRating, averageRating_below],
["total", femalespeaker_above, femalespeaker_below, df_combined.count().title]])
print(contingency_table)
```

	female	male	total
above average	113	923	1036
below average	189	771	960
total	302	1694	1996

```
In [63]: # Compute the test(s)
from scipy import stats
result_fisher, p = stats.fisher_exact([[femalespeaker_aboveRating, malespeaker_aboveRating], [femalespeaker_belowRating, malespeaker_belowRating]])
print("The p value of the Fisher test is:")
print(p)
```

The p value of the Fisher test is:
4.648303305358867e-08

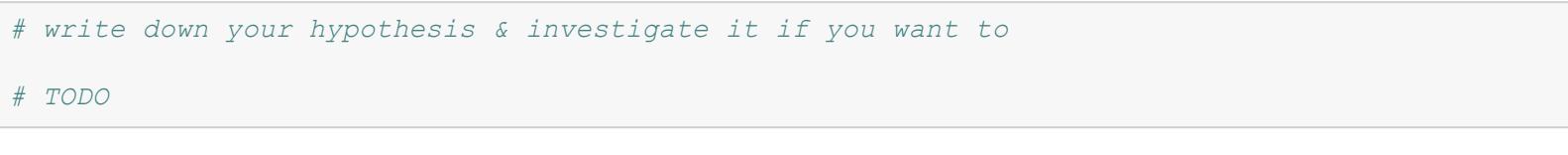
Part IV: custom analysis

The problem tests such as the one above is that they contain a lot of arbitrary choices and discretizations. Try to think for yourself why this kind of test is not a good idea.

However, methods that do not arbitrarily split the data into binary categories are better. For example, we could find the *linear least squares* best fit for a linear function $f(x) = ax + b$ if x is the percentage of female dialogue and $f(x)$ is the IMBD rating:

```
In [64]: X = np.zeros((len(df_combined.f_percentage), 2))
X[:, 0] = np.ones(len(df_combined.f_percentage))
X[:, 1] = np.array(df_combined.f_percentage).flatten()
Y = np.array(df_combined.averageRating).flatten()
w = np.linalg.solve(X.T @ X, (X.T @ Y))
plt.xlabel("IMDB Rating", size=20)
plt.ylabel("Percent Female Dialogue", size=20)
plt.title("Percent Female Dialogue vs. IMBD Rating - linear regression", size=20)
plt.show()
```

Percent Female Dialogue vs. IMBD Rating - linear regression



We can see that this kind of analysis is dangerous if, instead of plotting rating against the female participation, we plot the rating against the length of the movie, and to the same analysis.

```
In [65]: X = np.zeros((len(df_combined.words), 2))
X[:, 0] = np.ones(len(df_combined.words)).flatten()
Y = np.linalg.lstsq(X.T @ X, df_combined.averageRating).flatten()
fig, ax1 = plt.subplots(1, 1, figsize=(10, 10))
ax1.plot(df_combined.f_percentage, df_combined.averageRating, "k", ms=3)
ax1.plot([0, 1], [w[0] + w[1] * e4], "r", lw=5)
plt.xlabel("IMDB Rating", size=20)
plt.ylabel("Percent Female Dialogue", size=20)
plt.title("Percent Female Dialogue vs. IMBD Rating - linear regression", size=20)
plt.show()
```

length of movie vs. IMBD Rating - linear regression

note that the line is very different from the one above. It is not a good idea to do this kind of analysis.

Before moving on, think about how you *believe* this data may be causally related. Do people rank movies low because they contain women?

Or do unpopular movies have a higher chance of containing many female characters? What is the most striking thing about the plots above: a potential correlation between female participation and rating, or isn't it rather the fact that there are so few movies with majority female conversation?

In situations like this, it can be helpful to look at the data not as a supervised problem, but as an unsupervised problem: What is the joint distribution between two variables, such as female participation and rating?

and instead consider it as an unsupervised problem. The Binomial test is for testing hypothesis which we don't have at this point.

1. Compute a *p*-value for your test of choice.

```
In [37]: # recreating the "Who is talking in popular films" scatter plot
# add the average rating
mean_rating = df_combined["averageRating"].mean()
print("Mean of all movie ratings")
print(mean_rating)
# TODO: this is optional but very helpful to understand what we are doing
```

Mean of all movie ratings
6.8044589178356825

```
In [57]: # determine the totals on the edge
femalespeaker_above = df_combined[df_combined["f_percentage"] > 0.5].count().title
femalespeaker_below = df_combined[df_combined["f_percentage"] <= 0.5].count().title
averageRating_above = df_combined[df_combined["averageRating"] > mean_rating].count().title
averageRating_below = df_combined[df_combined["averageRating"] <= mean_rating].count().title
# determine all 4 boxes in the contingency table
femalespeaker_aboveRating = df_combined[(df_combined["f_percentage"] > 0.5) & (df_combined["averageRating"] > mean_rating)].count().title
femalespeaker_belowRating = df_combined[(df_combined["f_percentage"] <= 0.5) & (df_combined["averageRating"] < mean_rating)].count().title
malespeaker_aboveRating = df_combined[(df_combined["f_percentage"] <= 0.5) & (df_combined["averageRating"] > mean_rating)].count().title
malespeaker_belowRating = df_combined[(df_combined["f_percentage"] > 0.5) & (df_combined["averageRating"] <= mean_rating)].count().title

contingency_table = pd.DataFrame([["above average", "female", "male", "total"],
["above average", femalespeaker_aboveRating, malespeaker_aboveRating, averageRating_above],
["below average", femalespeaker_belowRating, malespeaker_belowRating, averageRating_below],
["total", femalespeaker_above, femalespeaker_below, df_combined.count().title]])
print(contingency_table)
```

	female	male	total
above average	113	923	1036
below average	189	771	960
total	302	1694	1996

```
In [63]: # Compute the test(s)
from scipy import stats
result_fisher, p = stats.fisher_exact([[femalespeaker_aboveRating, malespeaker_aboveRating], [femalespeaker_belowRating, malespeaker_belowRating]])
print("The p value of the Fisher test is:")
print(p)
```

The p value of the Fisher test is:
4.648303305358867e-08

The problem tests such as the one above is that they contain a lot of arbitrary choices and discretizations. Try to think for yourself why this kind of test is not a good idea.

However, methods that do not arbitrarily split the data into binary categories are better. For example, we could find the *linear least squares* best fit for a linear function $f(x) = ax + b$ if x is the percentage of female dialogue and $f(x)$ is the IMBD rating:

```
In [64]: X = np.zeros((len(df_combined.f_percentage), 2))
X[:, 0] = np.ones(len(df_combined.f_percentage))
X[:, 1] = np.array(df_combined.f_percentage).flatten()
Y = np.array(df_combined.averageRating).flatten()
w = np.linalg.solve(X.T @ X, (X.T @ Y))
plt.xlabel("IMDB Rating", size=20)
plt.ylabel("Percent Female Dialogue", size=20)
plt.title("Percent Female Dialogue vs. IMBD Rating - linear regression", size=20)
plt.show()
```

Percent Female Dialogue vs. IMBD Rating - linear regression

We can see that this kind of analysis is dangerous if, instead of plotting rating against the female participation, we plot the rating against the length of the movie, and to the same analysis.

```
In [65]: X = np.zeros((len(df_combined.words), 2))
X[:, 0] = np.ones(len(df_combined.words)).flatten()
Y = np.linalg.lstsq(X.T @ X, df_combined.averageRating).flatten()
fig, ax1 = plt.subplots(1, 1, figsize=(10, 10))
ax1.plot(df_combined.f_percentage, df_combined.averageRating, "k", ms=3)
ax1.plot([0, 1], [w[0] + w[1] * e4], "r", lw=5)
plt.xlabel("IMDB Rating", size=20)
plt.ylabel("Percent Female Dialogue", size=20)
plt.title("Percent Female Dialogue vs. IMBD Rating - linear regression", size=20)
plt.show()
```

length of movie vs. IMBD Rating - linear regression

note that the line is very different from the one above. It is not a good idea to do this kind of analysis.

Before moving on, think about how you *believe* this data may be causally related. Do people rank movies low because they contain women?

Or do unpopular movies have a higher chance of containing many female characters? What is the most striking thing about the plots above: a potential correlation between female participation and rating, or isn't it rather the fact that there are so few movies with majority female conversation?

In situations like this, it can be helpful to look at the data not as a supervised problem, but as an unsupervised problem: What is the joint distribution between two variables, such as female participation and rating?

and instead consider it as an unsupervised problem. The Binomial test is for testing hypothesis which we don't have at this point.

1. Compute a *p*-value for your test of choice.

```
In [37]: # recreate the "Who is talking in popular films" scatter plot
# add the average rating
mean_rating = df_combined["averageRating"].mean()
print("Mean of all movie ratings")
print(mean_rating)
# TODO: this is optional but very helpful to understand what we are doing
```

Mean of all movie ratings
6.8044589178356825

```
In [57]: # determine the totals on the edge
femalespeaker_above = df_combined[df_combined["f_percentage"] > 0.5].count().title
femalespeaker_below = df_combined[df_combined["f_percentage"] <= 0.5].count().title
averageRating_above = df_combined[df_combined["averageRating"] > mean_rating].count().title
averageRating_below = df_combined[df_combined["averageRating"] <= mean_rating].count().title
# determine all 4 boxes in the contingency table
femalespeaker_aboveRating = df_combined[(df_combined["f_percentage"] > 0.5) & (df_combined["averageRating"] > mean_rating)].count().title
femalespeaker_belowRating = df_combined[(df_combined["f_percentage"] <= 0.5) & (df_combined["averageRating"] > mean_rating)].count().title
malespeaker_aboveRating = df_combined[(df_combined["f_percentage"] <= 0.5) & (df_combined["averageRating"] < mean_rating)].count().title
malespeaker_belowRating = df_combined[(df_combined["f_percentage"] > 0.5) & (df_combined["averageRating"] <= mean_rating)].count().title

contingency_table = pd.DataFrame([["above average", "female", "male", "total"],
["above average", femalespeaker_aboveRating, malespeaker_aboveRating, averageRating_above],
["below average", femalespeaker_belowRating, malespeaker_belowRating, averageRating_below],
["total", femalespeaker_above, femalespeaker_below, df_combined.count().title]])
print(contingency_table)
```

	female	male	total
above average	113	923	1036
below average	189	771	960
total	302	1694	1996