

# Problem set 5

Dili Maduabum, Joshua Bailey

2024-02-16

## Problem 1

1

Paper read!

2

Social utility refers to the benefit or satisfaction that individuals derive from their actions within a social context, not solely based on their own direct outcomes but also influenced by the actions and outcomes of others in their social network or community. It's a concept that encapsulates the notion that individuals' utility or happiness can be affected by the choices, behaviours, and welfare of others, in addition to their own personal outcomes.

Social learning is the process through which individuals acquire new information, behaviours, or strategies by observing or interacting with others. It involves the diffusion of knowledge and practices through a social network, where individuals can learn from the experiences and decisions of their peers.

The difference between social utility and social learning lies in their focus and mechanism. Social utility is about the derived satisfaction or benefit from being part of a social setting and how others' actions and outcomes affect an individual's own utility. Social learning, in contrast, is focused on the acquisition of knowledge, behaviours, or norms through the observation of or interaction with others. It is a mechanism for information diffusion and behavioural adaptation within a social network, emphasizing the cognitive aspect of learning from the environment and others' experiences.

3

Designing an ideal experiment to differentiate between social utility and social learning involves isolating the effects of social influence on individual decision-making from the effects of learning from others' experiences. The key distinction between these two mechanisms is the role of private versus public information and the subsequent effects on individual behavior. Social utility refers to the value individuals derive from conforming to the behaviors of their social group, often driven by network externalities, where the utility of a product or action increases as more people adopt it. On the other hand, social learning involves updating one's beliefs and decisions based on observing the actions and outcomes of others, under the assumption of imperfect information about the quality of an option.

This requires a setup where individuals make decisions based not on learning about the quality of options from others but on the intrinsic value they place on aligning with the behaviors of their social group. Conversely, to isolate social learning, the experimental design should allow individuals to update their beliefs about the quality of options based on the observed choices of others, without the influence of wanting to conform to a group norm for its own sake.

## 4

A regression of the number of viewers in later weeks on the number of viewers in the first week would not have a causal interpretation due to potential unobserved confounding variables. One such confounder could be the inherent attractiveness or quality of the movie, which simultaneously impacts viewership in the first week and viewership in later weeks. For example, a highly anticipated movie might attract a large audience in its first week due to its perceived quality, marketing, or star cast, and continue to attract viewers in subsequent weeks through word of mouth or sustained interest. This inherent quality affects both initial and later viewership but is not directly measured in a simple regression, leading to confounded estimates of the effect of first-week viewership on later viewership.

## 5

The authors exploit weather variations as an instrumental variable. This approach relies on the assumption that weather variations are random and thus not directly related to the movie's quality, but can significantly affect viewership numbers due to people's preference for indoor activities, like watching movies, during bad weather. By using weather as an instrument, the analysis can isolate the effect of early viewership on later viewership, effectively controlling for the unobserved movie quality that might simultaneously influence viewership across different weeks. This methodology allows for a cleaner estimation of the causal effect of social spillovers on movie consumption, sidestepping the confounding effect of inherent movie quality.

## 6

If Gilchrist and Sands had observed the confounding variable and implemented selection on observables under the conditional independence assumption, they could have potentially estimated the causal effect of first-week viewership on later-week viewership. This approach assumes that controlling for observed confounders—variables that affect both the treatment (first-week viewership) and the outcome (later-week viewership)—would allow for a causal interpretation of the relationship between the two.

However, merely observing the confounding variable and controlling for it would not necessarily help in distinguishing between social learning and social utility. The conditional independence assumption helps identify the causal effect of first-week viewership on later-week viewership, assuming all confounders are observed and correctly modeled. Still, this method alone does not differentiate whether the mechanism behind this effect is due to social learning (where individuals update their information based on others' actions) or social utility (where individuals derive utility from consuming a product in a social context, independent of information changes).

To disentangle social learning from social utility, one would need to identify mechanisms that specifically relate to information transmission or shared consumption experiences. For example, if the effect of first-week viewership on later-week viewership was found to be stronger for movies with more ambiguous quality (suggesting viewers are updating their beliefs based on others' actions), this would indicate social learning. Conversely, if the effect was uniform across movies of different quality levels, suggesting that the social aspect of viewing is driving the effect, this could indicate social utility. Gilchrist and Sands' approach using weather as an instrumental variable aims to capture the exogenous variation in first-week viewership unrelated to movie quality, potentially offering insights into the social spillover effects. However, observing the confounding variable and controlling for it through selection on observables might not provide clear evidence to distinguish between these two mechanisms without further assumptions or additional data on the nature of social interactions related to moviegoing.

## 7

1. **Variation in Responses Based on Information Precision:** Social learning theory suggests that the impact of observing others' actions (e.g., early viewership on later viewership) would vary depend-

ing on the precision of the information available to the observers. In contexts where there's greater uncertainty about a movie's quality, the influence of early viewers on later viewers should be more pronounced because potential viewers have more to learn from the actions of early viewers. This variation in response is consistent with social learning because it underlines the role of information acquisition in decision-making. In contrast, social utility, which derives from the shared experience of consumption, does not inherently vary with the precision of information about the product's quality.

2. **Differential Impact Based on Observable Characteristics of Consumers:** Social learning models also predict that the impact of early viewership on later viewership may differ across groups with different propensities for learning. For instance, if younger viewers are more influenced by peer behavior than older viewers, one might observe stronger social spillover effects among movies targeted at younger audiences. This implication is consistent with social learning as it reflects the idea that the propensity to update beliefs based on others' actions can vary across different segments of the population. Social utility does not offer a clear reason for why such differential impacts based on consumer characteristics should exist, as the enjoyment derived from shared consumption experiences does not inherently depend on the consumers' propensity to learn from others.

## 8

Gilchrist and Sands proxy for realized movie quality using ratings by expert reviewers, specifically from IMDb's top 1000 voters. This group is characterized by IMDb as the 1000 people who have voted for the most titles in their ratings polls. Movies are then categorized into quality terciles based on these ratings, capturing a significant dimension of demand as top tercile movies, on average, sell significantly more tickets over the first 6 weeks compared to bottom tercile movies.

They also use information about the movie's production budget. The rationale behind this is that advertising budgets, which influence the amount of information available to the public before a movie's release, are generally set as a fixed percentage of production budgets. Production budgets are thus positively correlated with the precision of prior knowledge about a movie.

While these proxies are insightful and useful for empirical analysis, they may not fully capture the nuances of movie quality or the exact precision of information. Movie quality, for instance, is subjective and can be influenced by factors beyond critic ratings, such as viewer preferences and cultural impacts, which are harder to quantify. Similarly, the precision of information might be affected by factors other than advertising, such as word-of-mouth or social media influence, which are not directly tied to the production budget.

## 9

Gilchrist and Sands do not expect the conditional expectation of a movie's total viewership with respect to weather to have a causal interpretation because of the strategic behavior of film studios in choosing release dates. Film studios tend to release movies at times when they expect higher demand, both in terms of the number and quality of movies available in theaters. This strategic timing means that the observed relationship between weather and movie viewership could be confounded by the studios' anticipation of demand variations across different times of the year.

The seasonality of movie viewership is driven by both underlying demand and supply-side considerations, where the supply side (i.e., the studios) takes into account expected demand in timing releases. This includes releasing certain types of movies during specific seasons to maximize viewership i.e. blockbuster movies are often released during summer and holiday seasons when audience turnout is expected to be higher. Thus, any observed correlation between weather and viewership could be influenced by these strategic release decisions, making it difficult to interpret the relationship causally.

## 10

The regression of opening weekend viewership on indicators for day of the week of the year, week, year, and holidays is not a saturated regression because it does not include every possible combination of these indicators as separate variables. A saturated model would include a unique indicator variable for every unique combination of the levels of these categorical variables, effectively allowing for a separate intercept for every combination. Given the vast number of combinations (especially when considering all weeks of all years, along with holidays), such a regression would be extremely high-dimensional and possibly exceed the number of observations available, leading to overfitting where the model perfectly fits the training data but performs poorly on new data.

They likely do not consider a saturated regression because it would involve an enormous number of parameters, leading to issues with overfitting and interpretability. High-dimensional models can be computationally challenging to estimate and may require regularization techniques to obtain stable estimates. That said, one advantage of such a model is that the coefficients of a saturated model can be straightforwardly interpreted.

## 11

**11.a** A movie should appear in the dataset at least 18 times. Each has a record for the weekend (Friday, Saturday and Sunday) from the opening weekend to at least 6 weekends later (for the ones kept). The ones dropped were not in theaters for more than 6 weekends.

### 11.b

```
#keeping films that aren't dropped
films_used <- films |>
  filter(dropped != 1)
```

### 11.c

```
# day when 12 Rounds came in
round_12_date <- as.Date("2009-03-27")

# Define the number of days to add
days_before <- 17984 #number under 12 Rounds "date" column

# Days prior to the
reference_date <- round_12_date - days_before + 1

# Print the new date
print(reference_date)
```

```
## [1] "1960-01-01"
```

### 11.d

```
films_used_d <- films_used |>
  mutate(movie_date = as.Date(reference_date + date)) |>
  #putting the release_date in the 4th column
  select(title, production_budget, release_yr,
         movie_date, sat_date, everything())

films_used_d[, c("title", "movie_date")]
```

```
## # A tibble: 24,855 x 2
##   title          movie_date
##   <chr>          <date>
## 1 (500) Days of Summer 2009-08-08
## 2 12 Rounds          2009-03-28
## 3 127 Hours           2010-11-26
## 4 13 Going on 30      2004-04-25
## 5 1408                 2007-06-24
## 6 16 Blocks           2006-03-05
## 7 17 Again            2009-04-19
## 8 2 Fast 2 Furious    2003-06-07
## 9 2012                 2009-11-15
## 10 21                  2008-03-28
## # i 24,845 more rows
```

11.e

```
#first using sat_date to get the date for each saturday
films_used_date <- films_used_d |>
  #getting the day for saturday
  mutate(sat_day = reference_date + sat_date) |>
  mutate(sat_day_of_week = wday(sat_day, label = TRUE)) |>
  mutate(
    fri_dummy = ifelse(movie_date == sat_day - 1, 1, 0),
    sat_dummy = ifelse(movie_date == sat_day, 1, 0),
    #reasoning... there was no movie released on Sunday...
    sun_dummy = ifelse(movie_date == sat_day + 1, 1, 0)
  ) |> arrange(title)

films_used_date[, c("title", "movie_date", "sat_day", "fri_dummy", "sat_dummy", "sun_dummy")]
```

```
## # A tibble: 24,855 x 6
##   title          movie_date sat_day    fri_dummy sat_dummy sun_dummy
##   <chr>          <date>    <date>      <dbl>      <dbl>      <dbl>
## 1 (500) Days of Summer 2009-08-08 2009-08-08         0         1         0
## 2 (500) Days of Summer 2009-08-07 2009-08-08         1         0         0
## 3 (500) Days of Summer 2009-08-09 2009-08-08         0         0         1
## 4 (500) Days of Summer 2009-08-14 2009-08-15         1         0         0
## 5 (500) Days of Summer 2009-08-15 2009-08-15         0         1         0
## 6 (500) Days of Summer 2009-08-16 2009-08-15         0         0         1
## 7 (500) Days of Summer 2009-08-21 2009-08-22         1         0         0
## 8 (500) Days of Summer 2009-08-23 2009-08-22         0         0         1
## 9 (500) Days of Summer 2009-08-22 2009-08-22         0         1         0
## 10 (500) Days of Summer 2009-08-30 2009-08-29         0         0         1
## # i 24,845 more rows
```

11.f

```
#creating dummies for week using fastDummies
films_used_date <- films_used_date |>
  arrange(title, sat_day) |>
  group_by(title) |>
  # Assign numeric labels to unique elements of sat_date within each title
```

```
mutate(week = as.integer(factor(sat_date)))

#Now using fast dummies...
films_used_date <- dummy_cols(films_used_date, select_columns = 'week')
films_used_date[, c("title", "movie_date", "week_1", "week_2", "week_3")]
```

```
## # A tibble: 24,855 x 5
##   title          movie_date week_1 week_2 week_3
##   <chr>          <date>    <int> <int> <int>
## 1 (500) Days of Summer 2009-08-08      1     0     0
## 2 (500) Days of Summer 2009-08-07      1     0     0
## 3 (500) Days of Summer 2009-08-09      1     0     0
## 4 (500) Days of Summer 2009-08-14      0     1     0
## 5 (500) Days of Summer 2009-08-15      0     1     0
## 6 (500) Days of Summer 2009-08-16      0     1     0
## 7 (500) Days of Summer 2009-08-21      0     0     1
## 8 (500) Days of Summer 2009-08-23      0     0     1
## 9 (500) Days of Summer 2009-08-22      0     0     1
## 10 (500) Days of Summer 2009-08-30      0     0     0
## # i 24,845 more rows
```

11.g

```
#using the "Fast Dummies" library... to automatically create dummies for year
film <- dummy_cols(films_used_date, select_columns = 'release_yr')

film[, c("title", "release_yr", "release_yr_2009", "release_yr_2010")]
```

```
## # A tibble: 24,855 x 4
##   title          release_yr release_yr_2009 release_yr_2010
##   <chr>          <dbl>        <int>        <int>
## 1 (500) Days of Summer    2009            1            0
## 2 (500) Days of Summer    2009            1            0
## 3 (500) Days of Summer    2009            1            0
## 4 (500) Days of Summer    2009            1            0
## 5 (500) Days of Summer    2009            1            0
## 6 (500) Days of Summer    2009            1            0
## 7 (500) Days of Summer    2009            1            0
## 8 (500) Days of Summer    2009            1            0
## 9 (500) Days of Summer    2009            1            0
## 10 (500) Days of Summer    2009            1            0
## # i 24,845 more rows
```

11.h

```
#combine the weekends
temp <- film |>
mutate(weekend = case_when(
  sat_dummy == 1 ~ "Saturday",
  fri_dummy == 1 ~ "Friday",
  sun_dummy == 1 ~ "Sunday",
```

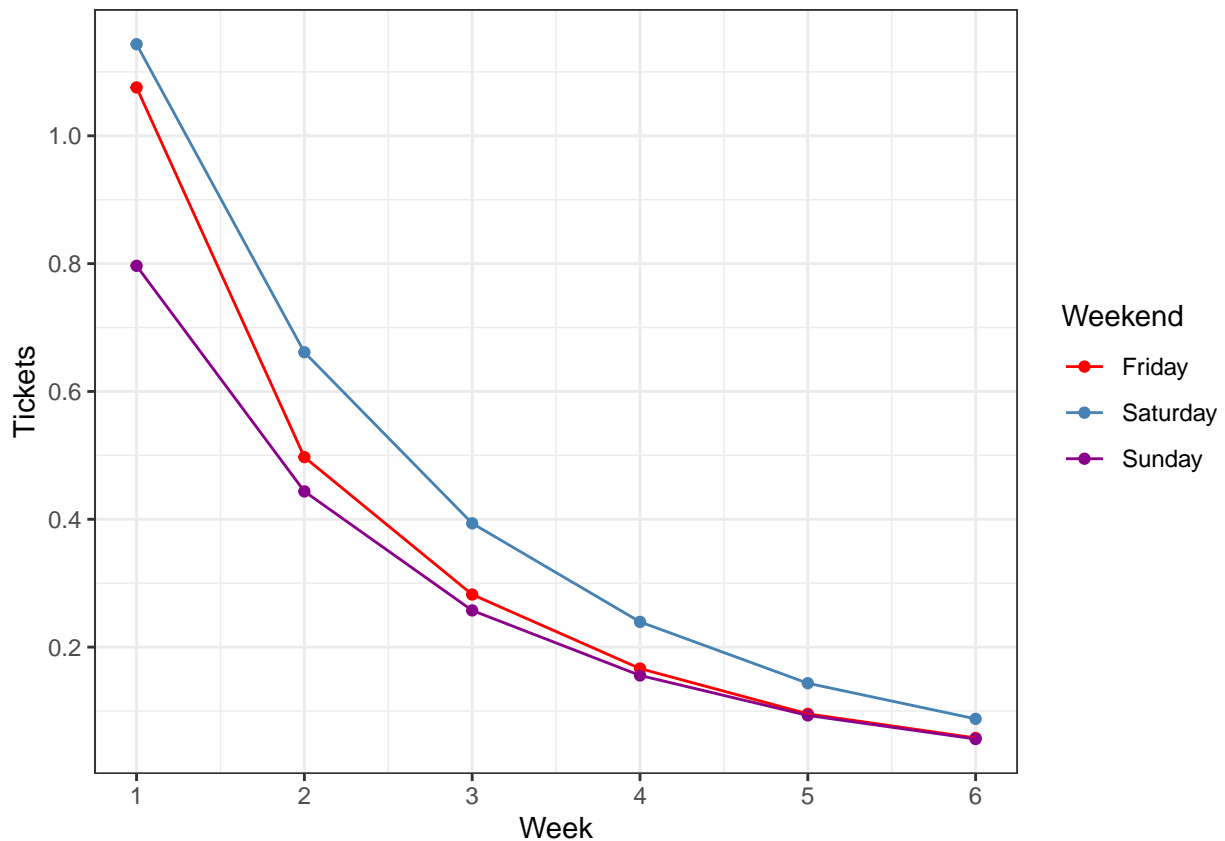
```

)) |>
  group_by(week, weekend) |>
  summarize(mean = mean(tickets, na.rm = TRUE))

temp |>
  ggplot(aes(x = week, y = mean, color = as.factor(weekend))) +
  geom_point() +
  geom_line() +
  scale_color_manual(values = c("Saturday" = "#4682B4",
                                "Friday" = "red",
                                "Sunday" = "#8B008B")) +

  labs(color = "Weekend",
        y = "Tickets",
        x = "Week") +
  scale_x_continuous(breaks = scales::pretty_breaks(n = 6)) + # Set x-axis ticks
  scale_y_continuous(breaks = scales::pretty_breaks(n = 6)) + # Set y-axis ticks
  theme_bw()

```



12

NOT NEEDED

13

```
#subset colnames that have the hh in them
holiday <- str_subset(colnames(film), "hh")

#make the things in holiday "add"
holiday_dummy <- str_c(holiday, collapse = " + ")

#day of the week dummies
weekend_dummy <- str_c(str_subset(colnames(film), "dummy"), collapse = " + ")

#week of the year dummies
week_dummy <- str_c(str_subset(colnames(film), "week_"), collapse = " + ")

#year of the week dummy
year_dummy <- str_c(str_subset(colnames(film), "release_yr_"), collapse = " + ")

#combine
mod1 <- glue("tickets ~ {weekend_dummy} + {week_dummy} + {year_dummy} + {holiday_dummy}")

#fit a regression model
reg_mod1 <- lm(as.formula(mod1), data = film)

film <- film |>
  mutate(pred_tickets = predict(reg_mod1, film)) |>
  mutate(abnormal_viewership = tickets - pred_tickets)

film[, c("tickets", "pred_tickets", "abnormal_viewership", "sat_day")]
```

```
## # A tibble: 24,855 x 4
##   tickets pred_tickets abnormal_viewership sat_day
##   <dbl>      <dbl>          <dbl> <date>
## 1  0.185        1.07          -0.890 2009-08-08
## 2  0.159        0.991          -0.833 2009-08-08
## 3  0.155        0.933          -0.777 2009-08-08
## 4  0.126        0.518          -0.393 2009-08-15
## 5  0.153        0.602          -0.449 2009-08-15
## 6  0.117        0.460          -0.343 2009-08-15
## 7  0.0981       0.296          -0.198 2009-08-22
## 8  0.0808       0.237          -0.156 2009-08-22
## 9  0.125        0.379          -0.254 2009-08-22
## 10 0.0660       0.114          -0.0478 2009-08-29
## # i 24,845 more rows
```

14

```
weather <- read_dta("data/weather_collapsed_day.dta")

#adding www to the column names
```



```
original_cols <- colnames(weather)

# adding prefix using the paste
colnames(weather) <- paste("www", original_cols, sep = "_")

weather

## # A tibble: 4,932 x 39
##   www_date   www_sat_date www_snow www_rain www_mat5_10 www_mat5_15 www_mat5_20
##   <date>     <date>         <dbl>   <dbl>     <dbl>     <dbl>     <dbl>
## 1 1982-01-03 1982-01-02         0.165  0.533         0.0391    0.0276    0.0318
## 2 1982-01-01 1982-01-02         0.227  0.446         0.0476    0.0405    0.0472
## 3 1982-01-02 1982-01-02         0.233  0.408         0.0150    0.0157    0.0508
## 4 1982-01-10 1982-01-09         0.422  0.0916        0.464     0.0742    0.0529
## 5 1982-01-08 1982-01-09         0.222  0.0314        0.0709    0.0606    0.119
## 6 1982-01-09 1982-01-09         0.435  0.00428       0.156     0.0510    0.0857
## 7 1982-01-17 1982-01-16         0.255  0.0538        0.386     0.0455    0.0352
## 8 1982-01-15 1982-01-16         0.420  0.0533        0.0916    0.0940    0.102
## 9 1982-01-16 1982-01-16         0.494  0.135         0.174     0.0322    0.0455
## 10 1982-01-24 1982-01-23         0.216  0.0675        0.158     0.0589    0.0388
## # i 4,922 more rows
## # i 32 more variables: www_mat5_25 <dbl>, www_mat5_30 <dbl>, www_mat5_35 <dbl>,
## #   www_mat5_40 <dbl>, www_mat5_45 <dbl>, www_mat5_50 <dbl>, www_mat5_55 <dbl>,
## #   www_mat5_60 <dbl>, www_mat5_65 <dbl>, www_mat5_70 <dbl>, www_mat5_75 <dbl>,
## #   www_mat5_80 <dbl>, www_mat5_85 <dbl>, www_mat5_90 <dbl>, www_mat5_95 <dbl>,
## #   www_prec_0 <dbl>, www_prec_1 <dbl>, www_prec_2 <dbl>, www_prec_3 <dbl>,
## #   www_prec_4 <dbl>, www_prec_5 <dbl>, www_cloud_0 <dbl>, ...
```

```
weather_film <- film |>
  left_join(weather |> #all variables except the saturday variable
    select(-c("www_sat_date")),
    #combine on dates, automatically filters out dates that don't match
    by = c("movie_date" = "www_date"))

#test <- weather_film |>
# arrange(title) |>
# select(title, movie_date, sat_day, contains("www"))
```

15

```
# Select columns with names containing "www_"
www_columns <- str_subset(colnames(weather_film), "www_")

# Create a copy of the original dataframe
df <- weather_film

# Define regression formula with dummy variables
regressors <- glue("~ {weekend_dummy} + {week_dummy} + {year_dummy} + {holiday_dummy}")

# Iterate over columns with names containing "www_"
for (columns in www_columns) {
```

```

# Construct regression formula
model <- paste(columns, regressors)

# Generate names for predicted values and residuals
pred_name <- paste("pred", columns, sep = "_")
resid_name <- paste("abnormal", columns, sep = "_")

# Add predicted values and residuals to the dataframe
df <- df |>
  mutate(!pred_name := predict(lm(as.formula(model), data = df), df)) |>
  #residuals = column - predicted_value_for_column
  mutate(!resid_name := eval(parse(text = columns)) - eval(parse(text = pred_name)))
}

#remove the predicted and original values, keeping only the residuals
new_weather <- df |>
  select(-c(contains("pred_www"), starts_with("www")))

```

16

```

#combine
#fit a regression model
week_2_data <- new_weather |>
  filter(week_2 == 1)

#using the same regression
reg_mod2 <- lm(as.formula(mod1), data = week_2_data)

new_weather_film_wk2 <- week_2_data |>
  mutate(pred_tickets_wk_2 = predict(reg_mod2, week_2_data)) |>
  mutate(abnormal_viewership_wk_2 = tickets - pred_tickets_wk_2)

new_weather_film_wk2[, c("tickets", "pred_tickets_wk_2", "week_2", "abnormal_viewership_wk_2")]

```

```

## # A tibble: 4,143 x 4
##   tickets pred_tickets_wk_2 week_2 abnormal_viewership_wk_2
##   <dbl>         <dbl> <int>         <dbl>
## 1  0.126         0.455     1         -0.330
## 2  0.153         0.615     1         -0.462
## 3  0.117         0.394     1         -0.277
## 4  0.689         0.412     1          0.277
## 5  0.671         0.350     1          0.321
## 6  0.976         0.572     1          0.404
## 7  0.0668        0.394     1         -0.327
## 8  0.116         0.455     1         -0.340
## 9  0.119         0.615     1         -0.497
## 10 0.0554        0.394     1         -0.339
## # i 4,133 more rows

```

17

```
#Make
#subsetting the data to just be week 1
week_1_data <- new_weather |>
  filter(week_1 == 1)

#creating the "abnormal viewerships in week 1"-----
mod1 <- glue("tickets ~ {weekend_dummy} + {week_dummy} + {year_dummy} + {holiday_dummy}")

#fit a regression model
reg_mod1 <- lm(as.formula(mod1), data = week_1_data)

new_weather_film_wk1 <- week_1_data |>
  mutate(pred_tickets_wk1 = predict(reg_mod1, week_1_data)) |>
  mutate(abnormal_viewership_wk1 = tickets - pred_tickets_wk1)
```

17.a OLS;

```
abnormal_weather_wk1_names <-
  str_subset(colnames(new_weather_film_wk1), "abnormal_www")

abnormal_weather_wk1 <-
  str_c(abnormal_weather_wk1_names, collapse = "+")

ols_glue <- glue("abnormal_viewership_wk1 ~ {abnormal_weather_wk1}")
ols_mod <- lm(as.formula(ols_glue),
  new_weather_film_wk1)

#modelssummary(list(ols_mod), output = "gt")
```

17.b

```
#subset the data to include the variables of interest
leaps_data <- new_weather_film_wk1 |>
  select(c(abnormal_viewership_wk1, all_of(abnormal_weather_wk1_names)))

forward <- regsubsets(abnormal_viewership_wk1 ~ .,
  data = leaps_data, method = "forward")
```

## Reordering variables and trying again:

```
# Get summary of the models
summary_forward <- summary(forward)

# Find the index of the model with the highest R-squared Adjusted
best_model_index_fwd <- which.max(summary_forward$adjr2) #9th model has the highest

# Get the names of predictors (coef) in the best model (9), without the intercept([-1])
best_adjr_predictors <- names(coef(forward, id = best_model_index_fwd)[-1])

# Print the selected predictors and the corresponding R-squared Adjusted value
best_adjr_predictors
```

```
## [1] "abnormal_www_rain"      "abnormal_www_mat5_60" "abnormal_www_mat5_85"
## [4] "abnormal_www_mat5_90"    "abnormal_www_prec_1"  "abnormal_www_cloud_0"
## [7] "abnormal_www_cloud_4"    "abnormal_www_cloud_5" "abnormal_www_cloud_8"
```

```
#running regressions based on the model from forward (adj R^2)
regs_fwd <- str_c(best_adjr_predictors, collapse = " + ")

fwd_glue <- glue("abnormal_viewership_wk1 ~ {regs_fwd}")

fwd_mod <- lm(as.formula(fwd_glue), data = new_weather_film_wk1)
```

17.c

```
#only show the last steps (trace = 0)
backward <- step(ols_mod, direction = "backward", trace=0)
best_bkwd_predictors <- names(coefficients(backward)[-1])

best_bkwd_predictors
```

```
## [1] "abnormal_www_rain"      "abnormal_www_mat5_45" "abnormal_www_mat5_55"
## [4] "abnormal_www_mat5_70"    "abnormal_www_mat5_75" "abnormal_www_prec_0"
## [7] "abnormal_www_cloud_3"    "abnormal_www_cloud_4" "abnormal_www_mat_la"
```

```
#running regressions based on the model from backward
regs_bkwd <- str_c(best_bkwd_predictors, collapse = " + ")

bkwd_glue <- glue("abnormal_viewership_wk1 ~ {regs_bkwd}")

bkwd_mod <- lm(as.formula(bkwd_glue), data = new_weather_film_wk1)
```

17.d i

```
ridge_mod <- cv.glmnet(
  x = as.matrix(new_weather_film_wk1 |>
    select(all_of(abnormal_weather_wk1_names))),
  y = new_weather_film_wk1 |>
    pull(abnormal_viewership_wk1), #pull gets the numeric values
  alpha = 0, # Ridge penalty
  nfolds = 5 # 5 fold cross validation
)

#lambda that minimizes the cross validation error
nfold_ridge <- ridge_mod$lambda.min

coefs_nfold_ridge <- coef(ridge_mod, s = nfold_ridge)
```

ii

```
# Find lambda within 1 standard deviation of the minimum lambda (storind estimates for now. Will predic
nfold_ridge_1se <- ridge_mod$lambda.1se

coefs_nfold_ridge_1se <- coef(ridge_mod, s = nfold_ridge_1se)
```

17.e *i*

```
lasso_mod <- cv.glmnet(  
  x = as.matrix(new_weather_film_wk1 |>  
    select(all_of(abnormal_weather_wk1_names))),  
  y = new_weather_film_wk1 |>  
    pull(abnormal_viewership_wk1), #pull gets the numeric values  
  alpha = 1, # Lasso penalty  
  nfolds = 5 # 5 fold cross validation  
)  
  
#lambda that minimizes the cross validation error  
nfold_lasso = lasso_mod$lambda.min
```

*ii*

```
# Find lambda within 1 standard deviation of the minimum lambda  
nfold_lasso_1se <- lasso_mod$lambda.1se
```

*iii*

## 20

The absence of Lord of the Rings: The Return of the King, Lord of the Rings: The Two Towers, Toy Story 3, X2: X-Men United, and Finding Nemo from the upper decile made Kuehn and Lampe suspicious.

Kuehn and Lampe collect underlying user ratings data and reconstruct the quality categories for movies, addressing anomalies they found in the categorization of movies into high and low-quality categories by Gilchrist and Sands. While maintaining the use of Gilchrist and Sands' empirical framework, Kuehn and Lampe tested additional implications of social learning models, such as the impact of sequels versus original movies on sales momentum and the effect of audience age demographics (e.g., movies with a larger share of teenage viewers) on momentum.

Kuehn and Lampe find strong evidence of social learning, contradicting Gilchrist and Sands' findings which suggested a limited role for learning and emphasized network (or demand) externalities. Kuehn and Lampe's results indicate that movies in the top quality tercile exhibit significant momentum, which is not observed for movies in the bottom tercile. This pattern is consistent with social learning but does not support the presence of network externalities. Kuehn and Lampe's tests on the flexibility of prior beliefs (comparing sequels to originals) and the size of social networks (comparing movies with larger shares of teenage viewers to those with smaller shares) further support the presence of social learning. They find momentum to be significant for originals and movies with larger shares of teenage viewers, especially within the PG-13 category, reinforcing the idea that social learning mechanisms are at play.

## 21

```
#movies <- read_csv("data/movie_lens_20m/movie.csv")  
#ratings <- read_csv("data/movie_lens_20m/rating.csv")
```

The paper is an economics paper because it examines the impact of social interactions on consumer behavior, specifically how early viewership influences later movie attendance through social spillovers. This topic falls within the realm of economics as it deals with understanding how individuals make decisions based on information and experiences shared within their social networks, and how these decisions aggregate to broader market outcomes. This is of fundamental interest to the field of economics.

## Problem 2

### 1

Given that  $Y | X, \beta \sim \mathcal{N}(X\beta, \sigma^2 I_n)$ , the density of  $Y$  conditional on  $X$  and  $\beta$  is given by the probability density function of the multivariate normal distribution. The pdf for a multivariate normal distribution with mean  $\mu$  and covariance matrix  $\Sigma$  is given by:

$$f(y; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (y - \mu)^\top \Sigma^{-1} (y - \mu) \right)$$

For our specific case,  $\mu = X\beta$  and  $\Sigma = \sigma^2 I_n$ , thus the conditional density of  $Y$  given  $X$  and  $\beta$  is:

$$f(Y | X, \beta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left( -\frac{1}{2\sigma^2} (Y - X\beta)^\top (Y - X\beta) \right)$$

This expression gives the likelihood of observing  $Y$  given the predictors  $X$  and the regression coefficients  $\beta$ , under the assumption of normally distributed errors with variance  $\sigma^2$ .

### 2

Given that  $\beta \sim \mathcal{N}(0, \tau^2 I_p)$ , the density of  $\beta$  is given by the probability density function of the multivariate normal distribution, where  $\mu = 0$  (a  $p$ -dimensional zero vector) and  $\Sigma = \tau^2 I_p$  (a  $p \times p$  covariance matrix where the variance of each component of  $\beta$  is  $\tau^2$  and all covariances are 0). The pdf for a multivariate normal distribution with mean  $\mu$  and covariance matrix  $\Sigma$  is:

$$f(x; \mu, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

In our case,  $k = p$ ,  $\mu = 0$ , and  $\Sigma = \tau^2 I_p$ , thus the density of  $\beta$  is:

$$f(\beta) = \frac{1}{(2\pi\tau^2)^{p/2}} \exp \left( -\frac{1}{2\tau^2} \beta^\top \beta \right)$$

This expression gives the prior belief about the distribution of the regression coefficients  $\beta$ , assuming they are drawn from a multivariate normal distribution centered at zero with covariance  $\tau^2 I_p$ , indicating a belief in the shrinkage towards zero of the coefficients (consistent with the principles of Ridge Regression).

### 3

To show that the posterior likelihood of  $\beta$  given  $Y$  and  $X$  satisfies the given expression, we'll use the densities for  $f(Y | X, \beta)$  and  $f(\beta)$  that we have previously derived and apply Bayes' theorem.

Given:

1. The conditional density of  $Y$  given  $X$  and  $\beta$ :

$$f(Y | X, \beta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(Y - X\beta)^\top(Y - X\beta)\right)$$

2. The prior density of  $\beta$ :

$$f(\beta) = \frac{1}{(2\pi\tau^2)^{p/2}} \exp\left(-\frac{1}{2\tau^2}\beta^\top\beta\right)$$

3. Bayes' theorem for the posterior likelihood:

$$f(\beta | Y, X) = \frac{f(Y | X, \beta)f(\beta)}{f(Y)}$$

where  $f(Y)$  does not depend on  $\beta$  and serves as a normalizing constant.

Substituting the expressions for  $f(Y | X, \beta)$  and  $f(\beta)$  into Bayes' theorem:

$$f(\beta | Y, X) = \frac{\frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(Y - X\beta)^\top(Y - X\beta)\right) \times \frac{1}{(2\pi\tau^2)^{p/2}} \exp\left(-\frac{1}{2\tau^2}\beta^\top\beta\right)}{f(Y)}$$

Since  $f(Y)$  does not depend on  $\beta$ , and noting that the constants involving  $(2\pi\sigma^2)^{n/2}$  and  $(2\pi\tau^2)^{p/2}$  do not depend on  $\beta$ , we can denote all the terms that do not depend on  $\beta$  by a constant  $C'$ . Thus, the expression simplifies to:

$$f(\beta | Y, X) = C' \exp\left(-\frac{1}{2\sigma^2}(Y - X\beta)^\top(Y - X\beta) - \frac{1}{2\tau^2}\beta^\top\beta\right)$$

Here,  $C'$  is a constant that may involve factors from the normalizing constant  $f(Y)$  and the constants pulled out of the exponential terms but crucially does not depend on  $\beta$ .

Given that  $C$  is defined as a constant that does not depend on  $\beta$ , we can express  $C'$  simply as  $C$  for clarity. Therefore, we show that:

$$f(\beta | Y, X) = C \exp\left(-\frac{1}{2\sigma^2}(Y - X\beta)^\top(Y - X\beta) - \frac{1}{2\tau^2}\beta^\top\beta\right)$$

## 4

To conclude that the  $\beta$  which maximizes the posterior likelihood (the mode of the posterior distribution) is given by the solution to the specified optimization problem, we start with the expression for the posterior likelihood derived in the previous question:

$$f(\beta | Y, X) = C \exp \left( -\frac{1}{2\sigma^2} (Y - X\beta)^t (Y - X\beta) - \frac{1}{2\tau^2} \beta^t \beta \right)$$

Maximizing the posterior likelihood  $f(\beta | Y, X)$  with respect to  $\beta$  is equivalent to minimizing the exponent's argument because the exponential function is monotonically increasing. Thus, we focus on minimizing the negative exponent in the expression, which involves two terms:

1. The term  $-\frac{1}{2\sigma^2} (Y - X\beta)^t (Y - X\beta)$  corresponds to the squared error between the observed and predicted values, scaled by the variance  $\sigma^2$ .
2. The term  $-\frac{1}{2\tau^2} \beta^t \beta$  corresponds to the penalty on the size of  $\beta$ , scaled by the variance  $\tau^2$ .

Dropping the constants and the negative sign, the problem of maximizing the posterior likelihood can be restated as minimizing the expression:

$$\frac{1}{2\sigma^2} (Y - X\beta)^t (Y - X\beta) + \frac{1}{2\tau^2} \beta^t \beta$$

Multiplying through by  $2\sigma^2$  (which does not change the location of the minimum) to simplify, we get:

$$(Y - X\beta)^t (Y - X\beta) + \frac{\sigma^2}{\tau^2} \beta^t \beta$$

Recognizing that  $\beta^t \beta = \|\beta\|_2^2$ , the optimization problem becomes:

$$\min_{\beta} (Y - X\beta)^t (Y - X\beta) + \frac{\sigma^2}{\tau^2} \|\beta\|_2^2$$

This is precisely the ridge regression optimization problem, where  $(Y - X\beta)^t (Y - X\beta)$  represents the least squares error term, and  $\frac{\sigma^2}{\tau^2} \|\beta\|_2^2$  is the ridge penalty on the size of  $\beta$ , with the penalty strength controlled by the ratio  $\frac{\sigma^2}{\tau^2}$ . Thus, the  $\beta$  that maximizes the posterior likelihood under the Bayesian interpretation is indeed given by the solution to this optimization problem, highlighting the connection between Bayesian inference and ridge regression.

## 5

In Bayesian inference, we begin with a prior belief about the parameters (in this case,  $\beta$ ) of our linear model. This belief is quantified by assuming that  $\beta$  follows a normal distribution centered around zero, with its dispersion controlled by  $\tau^2$ . This assumption reflects our expectation that the true values of the regression coefficients are likely to be small or close to zero, enforcing a form of regularization even before observing the data.

When we observe the data (represented by  $Y$  and  $X$ ), we update our beliefs about  $\beta$  by calculating the posterior distribution, which combines our prior belief and the likelihood of observing the data given  $\beta$ . The likelihood of the data given  $\beta$  is also modeled using a normal distribution, where the variance  $\sigma^2$  captures the uncertainty or noise in the observations.



Maximizing the posterior likelihood of  $\beta$ , or equivalently, finding the mode of the posterior distribution, involves adjusting our estimate of  $\beta$  to best fit the observed data while also adhering to our initial regularization imposed by the prior. This balance is captured mathematically by an optimization problem that seeks to minimize the sum of two terms: one representing the fit of the model to the data (the least squares error term) and the other representing the size of the coefficients (the ridge penalty term), weighted by the ratio of the variance of the observation noise to the variance of the prior on  $\beta$  ( $\frac{\sigma^2}{\tau^2}$ ).

This optimization problem is precisely that of ridge regression, which seeks to find regression coefficients that not only fit the data well but are also small in magnitude, to prevent overfitting and improve the generalizability of the model. The strength of the regularization is controlled by the ratio  $\frac{\sigma^2}{\tau^2}$ , linking the Bayesian interpretation directly to the regularization parameter in ridge regression. This parameter balance reflects the trade-off between fitting the data closely (minimizing prediction error) and keeping the model simple and robust to new data (minimizing the size of the coefficients).