

Customer Analytics Assignment 4

Beiming Zhang

BU.450.760.K1.SP24

*Light-colored text boxes contain code,
and dark-colored text boxes contain
output results.*

Completely independent assignment.



Task 1 (a)

- Import and process data, converting some non-numeric variables into factors and constructing a treatment indicator(`treat_ind`).
- Build a linear regression model using the treatment indicator and shares.

```
1 # Preparation of operating environment
2 library(tableone)
3 library(MatchIt)
4 library(lattice)
5 options(width = 120)
6 setwd(paste0(
7   "/Users/velen/Documents/",
8   "文稿-iCloud/Learning/JHU/Spring I/Customer Analytics/Class 5"
9 ))
10 rm(list = ls())
11 ds <- read.csv("D5.2 Mashable.csv")
12
13 # Data cleaning
14 summary(ds)
15 table(ds$disclosed)
16 ds$category <- as.factor(ds$category)
17 ds$weekday <- as.factor(ds$weekday)
18
19 # Creat the treatment indicator
20 ds$treat_ind <- ifelse(ds$num_videos > 0, 1, 0)
21
22 # Task 1a
23 # Build the linear regression
24 model_t1 <- lm(shares ~ treat_ind, data = ds)
25 summary(model_t1)
```

Task 1 (a)

- In the linear regression results, the coefficient corresponding to the treatment indicator is significantly greater than 0, and the p-value is well below 0.01.
- Therefore, we can conclude that there is a positive correlation between the treatment indicator and shares, indicating that the treatment is associated with a typically larger number of shares.

```
1 > summary(model_t1)
2
3 Call:
4 lm(formula = shares ~ treat_ind, data = ds)
5
6 Residuals:
7     Min       1Q   Median       3Q      Max
8  -4309   -2310   -1691    -491   838990
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept)  2891.47      73.58   39.30  <2e-16 ***
13 treat_ind    1418.71     123.76   11.46  <2e-16 ***
14 ---
15 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16
17 Residual standard error: 11640 on 38710 degrees of freedom
18 Multiple R-squared:  0.003383, Adjusted R-squared:  0.003358
19 F-statistic: 131.4 on 1 and 38710 DF, p-value: < 2.2e-16
```

Task 2 (a)

- Remove irrelevant variables, retaining only appropriate covariates. Here, the “url”, “timedelta” and variables related to the independent and dependent variables were deleted.
- Generating a summary table that describes the baseline characteristics of the origin dataset and providing a statistical measure of how different the groups are from each other based on those variables.
- Calculate the propensity score using logistic regression and use a histogram to determine the overlap of the propensity score.

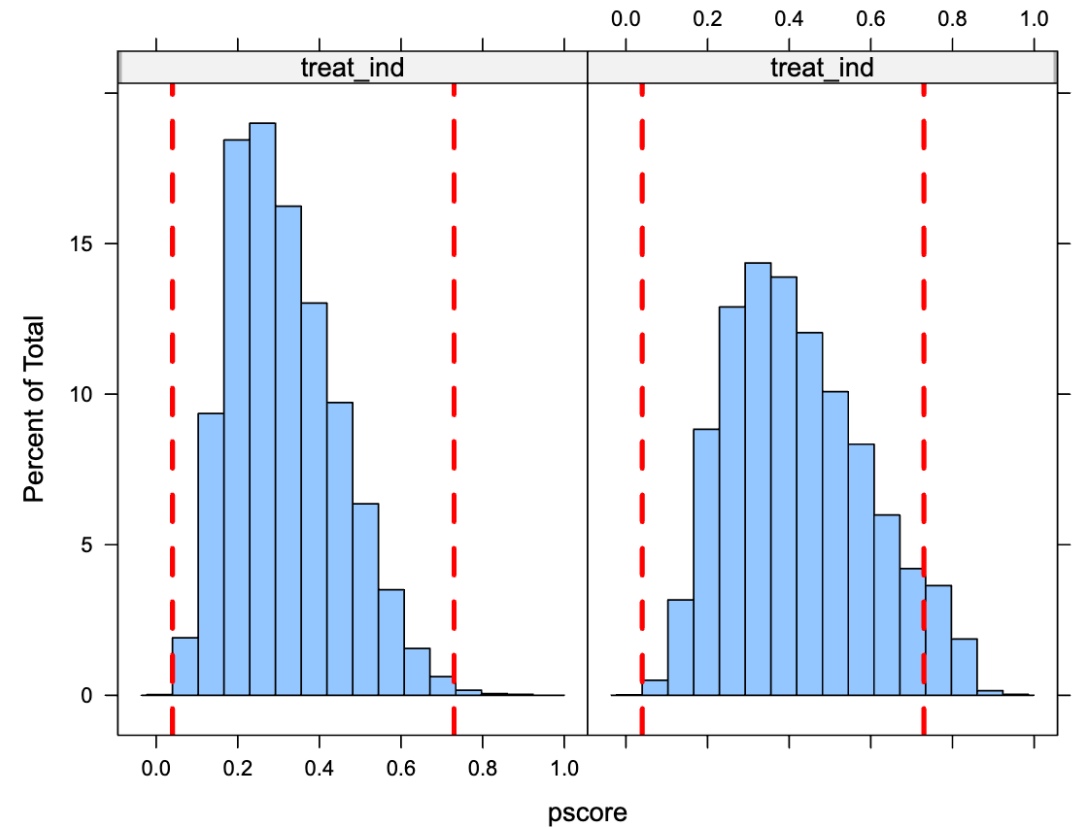
```
1 # Task 2a
2 # Remove irrelevant covariates
3 drop_list <- c("url", "timedelta", "shares", "treat_ind", "num_videos")
4 xvars <- setdiff(names(ds), drop_list)
5
6 # Pre-match assessment of balance
7 table_unmatched <- CreateTableOne(
8   vars = xvars, data = ds,
9   strata = "treat_ind"
10 )
11 print(table_unmatched, smd = TRUE)
12
13 # Calculate the pscore
14 y_x <- as.formula(paste("treat_ind ~", paste(xvars, collapse = " + ")))
15 model_t2 <- glm(y_x, family = binomial, data = ds)
16
17 # Assess pscore overlap
18 ds$pscore <- predict(model_t2, type = "response")
19 histogram(~ pscore | treat_ind, data = ds)
20
21 # Draw reference lines
22 histogram(~ pscore | treat_ind,
23   data = ds,
24   panel = function(x, ...) {
25     panel.histogram(x, ...)
26     panel.abline(v = c(.04, .73), col = "red", lty = 2, lwd = 3)
27   }
28 )
```

Task 2 (a)

```

1 > print(table_unmatched, smd = TRUE)
2
3           Stratified by treat_ind
4           0           1           p           test SMD
5 n           25026           13686
6 n_tokens_title (mean (SD))           10.21 (2.10)           10.69 (2.10)           <0.001           0.226
7 n_tokens_content (mean (SD))           554.98 (449.36)           526.87 (501.64)           <0.001           0.059
8 n_unique_tokens (mean (SD))           0.56 (4.43)           0.53 (0.17)           0.429           0.010
9 n_non_stop_unique_tokens (mean (SD))           0.70 (4.11)           0.67 (0.19)           0.286           0.013
10 num_hrefs (mean (SD))           11.04 (11.31)           10.61 (11.34)           <0.001           0.038
11 num_self_hrefs (mean (SD))           3.25 (4.06)           3.34 (3.52)           0.028           0.024
12 num_imgs (mean (SD))           4.68 (8.12)           4.20 (8.42)           <0.001           0.059
13 average_token_length (mean (SD))           4.64 (0.63)           4.39 (1.11)           <0.001           0.283
14 num_keywords (mean (SD))           7.19 (1.94)           7.30 (1.84)           <0.001           0.057
15 kwshares_worst (mean (SD))           318.09 (510.46)           302.34 (768.52)           0.016           0.024
16 kwshares_best (mean (SD))           236112.80 (121726.01)           271367.86 (133997.25)           <0.001           0.275
17 kwshares_avg (mean (SD))           2969.13 (1164.23)           3259.51 (1453.11)           <0.001           0.221
18 self_reference_min_shares (mean (SD))           4027.48 (21038.58)           4022.10 (17739.38)           0.980           <0.001
19 self_reference_max_shares (mean (SD))           8906.44 (35549.18)           12945.53 (49804.74)           <0.001           0.093
20 self_reference_avg_shares (mean (SD))           5873.04 (23925.79)           7433.21 (25319.44)           <0.001           0.063
21 global_subjectivity (mean (SD))           0.44 (0.10)           0.44 (0.14)           0.018           0.024
22 global_sentiment_polarity (mean (SD))           0.12 (0.09)           0.11 (0.10)           <0.001           0.109
23 global_rate_positive_words (mean (SD))           0.04 (0.02)           0.04 (0.02)           <0.001           0.090
24 global_rate_negative_words (mean (SD))           0.02 (0.01)           0.02 (0.01)           <0.001           0.088
25 rate_positive_words (mean (SD))           0.70 (0.17)           0.65 (0.22)           <0.001           0.229
26 rate_negative_words (mean (SD))           0.29 (0.15)           0.29 (0.17)           0.097           0.017
27 avg_positive_polarity (mean (SD))           0.36 (0.09)           0.35 (0.12)           <0.001           0.049
28 min_positive_polarity (mean (SD))           0.09 (0.07)           0.10 (0.08)           <0.001           0.043
29 max_positive_polarity (mean (SD))           0.76 (0.23)           0.74 (0.28)           <0.001           0.071
30 avg_negative_polarity (mean (SD))           -0.25 (0.12)           -0.27 (0.14)           <0.001           0.095
31 min_negative_polarity (mean (SD))           -0.51 (0.28)           -0.53 (0.30)           <0.001           0.062
32 max_negative_polarity (mean (SD))           -0.11 (0.09)           -0.11 (0.11)           0.003           0.031
33 title_subjectivity (mean (SD))           0.27 (0.32)           0.31 (0.34)           <0.001           0.125
34 title_sentiment_polarity (mean (SD))           0.07 (0.26)           0.07 (0.28)           0.760           0.003
35 abs_title_subjectivity (mean (SD))           0.34 (0.19)           0.34 (0.19)           0.068           0.019
36 abs_title_sentiment_polarity (mean (SD))           0.15 (0.22)           0.17 (0.24)           <0.001           0.094
37 category (%)
38   business           4893 (19.6)           1365 (10.0)
39   entertainment           3152 (12.6)           2973 (21.7)
40   lifestyle           1631 ( 6.5)           468 ( 3.4)
41   socialmedia           1642 ( 6.6)           681 ( 5.0)
42   tech           5275 (21.1)           2071 (15.1)
43   world           8433 (33.7)           6128 (44.8)
44 weekday (%)
45   friday           3561 (14.2)           2008 (14.7)
46   monday           4164 (16.6)           2313 (16.9)
47   saturday           1677 ( 6.7)           724 ( 5.3)
48   sunday           1791 ( 7.2)           853 ( 6.2)
49   thursday           4619 (18.5)           2505 (18.3)
50   tuesday           4541 (18.1)           2683 (19.6)
51   wednesday           4673 (18.7)           2600 (19.0)

```



Task 2 (a)

- From the table, we can see that the gap between the covariates of the two sample groups is not large, with the p-values generally being less than 0.001.
- From the overlap diagram, the common overlap between the two is concentrated in the range of $0.04 \leq n \leq 0.73$, indicating a relatively large overlap range.

Task 2 (b)

- Using the nearest neighbor matching method, perform the matching process based on the propensity score, and filter the matched dataset to ensure it only includes samples within the propensity score overlap region ($0.04 \leq p \leq 0.73$).
- In the result, there are 13,686 values corresponding to 0 and 13,686 values corresponding to 1, which is balanced.

```
1 # Task 2b
2 # Perform matching by pscore
3 ds_matched <- ds[ds$pscore >= .04 & ds$pscore <= .73, ]
4 matched <- matchit(y_x, method = "nearest", data = ds)
5 ds_matched <- match.data(matched)
6 table(ds_matched$treat_ind)
```

Task 2 (c)

- Construct a table to assess the sample in terms of covariate balance, comparing samples from two different datasets.
- After comparing the two (see next page), I believe the propensity score matching process was successful because it effectively reduced the SMD, improved the balance between the two groups, and the corresponding p-value also increased. A reduction in SMD means that the differences in these covariates between the two groups have been effectively minimized, enhancing the accuracy of the estimation.

```
1 # Task 2c
2 # Match assessment of balance
3 table_matched <- CreateTableOne(
4   vars = xvars, data = ds_matched,
5   strata = "treat_ind"
6 )
7 print(table_unmatched, smd = TRUE)
8 print(table_matched, smd = TRUE)
```


Task 2 (c)

```
1 > print(table_unmatched, smd = TRUE)
2
3           Stratified by treat_ind
4           0           1           p           test SMD
5 n
6 n_tokens_title (mean (SD))      10.21 (2.10)      10.69 (2.10)      <0.001      0.226
7 n_tokens_content (mean (SD))    554.98 (449.36)    526.87 (501.64)      <0.001      0.059
8 n_unique_tokens (mean (SD))     0.56 (4.43)      0.53 (0.17)      0.429      0.010
9 n_non_stop_unique_tokens (mean (SD)) 0.70 (4.11)      0.67 (0.19)      0.286      0.013
10 num_hrefs (mean (SD))          11.04 (11.31)    10.61 (11.34)      <0.001      0.038
11 num_self_hrefs (mean (SD))      3.25 (4.06)      3.34 (3.52)      0.028      0.024
12 num_imgs (mean (SD))            4.68 (8.12)      4.20 (8.42)      <0.001      0.059
13 average_token_length (mean (SD)) 4.64 (0.63)      4.39 (1.11)      <0.001      0.283
14 num_keywords (mean (SD))        7.19 (1.94)      7.30 (1.84)      <0.001      0.057
15 kwshares_worst (mean (SD))      318.09 (510.46)  302.34 (768.52)      0.016      0.024
16 kwshares_best (mean (SD))      236112.80 (121726.01) 271367.86 (133997.25) <0.001      0.275
17 kwshares_avg (mean (SD))        2969.13 (1164.23)  3259.51 (1453.11) <0.001      0.221
18 self_reference_min_shares (mean (SD)) 4027.48 (21038.58) 4022.10 (17739.38) 0.980      <0.001
19 self_reference_max_shares (mean (SD)) 8906.44 (35549.18) 12945.53 (49804.74) <0.001      0.093
20 self_reference_avg_sharess (mean (SD)) 5873.04 (23925.79) 7433.21 (25319.44) <0.001      0.063
21 global_subjectivity (mean (SD)) 0.44 (0.10)      0.44 (0.14)      0.018      0.024
22 global_sentiment_polarity (mean (SD)) 0.12 (0.09)      0.11 (0.10)      <0.001      0.109
23 global_rate_positive_words (mean (SD)) 0.04 (0.02)      0.04 (0.02)      <0.001      0.090
24 global_rate_negative_words (mean (SD)) 0.02 (0.01)      0.02 (0.01)      <0.001      0.088
25 rate_positive_words (mean (SD)) 0.70 (0.17)      0.65 (0.22)      <0.001      0.229
26 rate_negative_words (mean (SD)) 0.29 (0.15)      0.29 (0.17)      0.097      0.017
27 avg_positive_polarity (mean (SD)) 0.36 (0.09)      0.35 (0.12)      <0.001      0.049
28 min_positive_polarity (mean (SD)) 0.09 (0.07)      0.10 (0.08)      <0.001      0.043
29 max_positive_polarity (mean (SD)) 0.76 (0.23)      0.74 (0.28)      <0.001      0.071
30 avg_negative_polarity (mean (SD)) -0.25 (0.12)      -0.27 (0.14)      <0.001      0.095
31 min_negative_polarity (mean (SD)) -0.51 (0.28)      -0.53 (0.30)      <0.001      0.062
32 max_negative_polarity (mean (SD)) -0.11 (0.09)      -0.11 (0.11)      0.003      0.031
33 title_subjectivity (mean (SD)) 0.27 (0.32)      0.31 (0.34)      <0.001      0.125
34 title_sentiment_polarity (mean (SD)) 0.07 (0.26)      0.07 (0.28)      0.760      0.003
35 abs_title_subjectivity (mean (SD)) 0.34 (0.19)      0.34 (0.19)      0.068      0.019
36 abs_title_sentiment_polarity (mean (SD)) 0.15 (0.22)      0.17 (0.24)      <0.001      0.094
37 category (%)
38   business      4893 (19.6)      1365 (10.0)
39   entertainment 3152 (12.6)      2973 (21.7)
40   lifestyle     1631 ( 6.5)      468 ( 3.4)
41   socialmedia   1642 ( 6.6)      681 ( 5.0)
42   tech          5275 (21.1)      2071 (15.1)
43   world         8433 (33.7)      6128 (44.8)
44 weekday (%)
45   friday        3561 (14.2)      2008 (14.7)
46   monday        4164 (16.6)      2313 (16.9)
47   saturday      1677 ( 6.7)      724 ( 5.3)
48   sunday        1791 ( 7.2)      853 ( 6.2)
49   thursday      4619 (18.5)      2505 (18.3)
50   tuesday       4541 (18.1)      2683 (19.6)
51   wednesday     4673 (18.7)      2600 (19.0)
```

```
1 > print(table_matched, smd = TRUE)
2
3           Stratified by treat_ind
4           0           1           p           test SMD
5 n
6 n_tokens_title (mean (SD))      10.55 (2.11)      10.69 (2.10)      <0.001      0.064
7 n_tokens_content (mean (SD))    553.93 (457.12)    526.87 (501.64)      <0.001      0.056
8 n_unique_tokens (mean (SD))     0.58 (5.99)      0.53 (0.17)      0.356      0.011
9 n_non_stop_unique_tokens (mean (SD)) 0.72 (5.55)      0.67 (0.19)      0.288      0.013
10 num_hrefs (mean (SD))          11.02 (11.40)    10.61 (11.34)      0.003      0.036
11 num_self_hrefs (mean (SD))      3.40 (4.45)      3.34 (3.52)      0.263      0.014
12 num_imgs (mean (SD))            4.57 (7.19)      4.20 (8.42)      <0.001      0.047
13 average_token_length (mean (SD)) 4.55 (0.79)      4.39 (1.11)      <0.001      0.169
14 num_keywords (mean (SD))        7.28 (1.96)      7.30 (1.84)      0.602      0.006
15 kwshares_worst (mean (SD))      300.11 (591.01)    302.34 (768.52)      0.787      0.003
16 kwshares_best (mean (SD))      254273.47 (117490.41) 271367.86 (133997.25) <0.001      0.136
17 kwshares_avg (mean (SD))        3107.92 (1276.15)  3259.51 (1453.11) <0.001      0.111
18 self_reference_min_shares (mean (SD)) 4214.66 (21212.23) 4022.10 (17739.38) 0.415      0.010
19 self_reference_max_shares (mean (SD)) 10396.47 (43135.68) 12945.53 (49804.74) <0.001      0.055
20 self_reference_avg_sharess (mean (SD)) 6546.76 (25996.14) 7433.21 (25319.44) 0.004      0.035
21 global_subjectivity (mean (SD)) 0.45 (0.11)      0.44 (0.14)      <0.001      0.070
22 global_sentiment_polarity (mean (SD)) 0.12 (0.10)      0.11 (0.10)      0.002      0.038
23 global_rate_positive_words (mean (SD)) 0.04 (0.02)      0.04 (0.02)      <0.001      0.045
24 global_rate_negative_words (mean (SD)) 0.02 (0.01)      0.02 (0.01)      0.968      <0.001
25 rate_positive_words (mean (SD)) 0.68 (0.18)      0.65 (0.22)      <0.001      0.113
26 rate_negative_words (mean (SD)) 0.30 (0.15)      0.29 (0.17)      <0.001      0.049
27 avg_positive_polarity (mean (SD)) 0.36 (0.10)      0.35 (0.12)      <0.001      0.063
28 min_positive_polarity (mean (SD)) 0.10 (0.07)      0.10 (0.08)      0.632      0.006
29 max_positive_polarity (mean (SD)) 0.76 (0.24)      0.74 (0.28)      <0.001      0.068
30 avg_negative_polarity (mean (SD)) -0.27 (0.13)      -0.27 (0.14)      0.124      0.019
31 min_negative_polarity (mean (SD)) -0.55 (0.29)      -0.53 (0.30)      <0.001      0.043
32 max_negative_polarity (mean (SD)) -0.11 (0.10)      -0.11 (0.11)      0.998      <0.001
33 title_subjectivity (mean (SD)) 0.30 (0.33)      0.31 (0.34)      0.004      0.035
34 title_sentiment_polarity (mean (SD)) 0.07 (0.27)      0.07 (0.28)      0.296      0.013
35 abs_title_subjectivity (mean (SD)) 0.34 (0.19)      0.34 (0.19)      0.631      0.006
36 abs_title_sentiment_polarity (mean (SD)) 0.16 (0.23)      0.17 (0.24)      0.004      0.034
37 category (%)
38   business      1403 (10.3)      1365 (10.0)
39   entertainment 2751 (20.1)      2973 (21.7)
40   lifestyle     515 ( 3.8)      468 ( 3.4)
41   socialmedia   754 ( 5.5)      681 ( 5.0)
42   tech          2280 (16.7)      2071 (15.1)
43   world         5983 (43.7)      6128 (44.8)
44 weekday (%)
45   friday        1996 (14.6)      2008 (14.7)
46   monday        2258 (16.5)      2313 (16.9)
47   saturday      796 ( 5.8)      724 ( 5.3)
48   sunday        933 ( 6.8)      853 ( 6.2)
49   thursday      2512 (18.4)      2505 (18.3)
50   tuesday       2603 (19.0)      2683 (19.6)
51   wednesday     2588 (18.9)      2600 (19.0)
```

Task 3 (a)

- Reconstruct the linear regression model using the matched dataset, where the coefficient of `treat_ind` represents the estimated ATE of videos on the number of shares.

```
1 # Task 3a
2 # Build the linear regression based on ds_matched
3 model_t3 <- lm(shares ~ treat_ind, data = ds_matched)
4 summary(model_t3)
```

Task 3 (a)

- Based on the regression results, with the coefficient of `treat_ind` being positive, we can conclude that videos indeed increase the number of shares, with an estimated increase of about 1320.8 shares.

```
1 > summary(model_t3)
2
3 Call:
4 lm(formula = shares ~ treat_ind, data = ds_matched)
5
6 Residuals:
7     Min       1Q   Median       3Q      Max
8  -4309   -2810   -1989    -689   838990
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept)   2989.3      114.3   26.15 < 2e-16 ***
13 treat_ind     1320.8      161.7    8.17 3.2e-16 ***
14 ---
15 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16
17 Residual standard error: 13370 on 27370 degrees of freedom
18 Multiple R-squared:  0.002433, Adjusted R-squared:  0.002397
19 F-statistic: 66.76 on 1 and 27370 DF, p-value: 3.205e-16
```

Task 3 (b)

- The data in 1.a was unprocessed, and after reprocessing with the propensity score, the samples associated with `treat_ind` became more balanced, reducing selection bias. This leads to a smoother and more objective overall sample distribution.
- Matching makes the groups more similar in terms of important covariates, thereby reducing the bias these covariates might introduce.
- At the same time, we removed some outliers, which will not be included in the model in 3.a. This also leads to a smaller coefficient in 3.a compared to 1.a.

Task 3 (c)

- The "fudge factor" in this case may include the following:
 1. Due to copyright and other reasons, some articles cannot legally cite videos. Some video materials may have strict copyright restrictions; therefore, these articles can only display images or textual descriptions, which results in the treatment indicator being 0.
 2. Some articles do not require the use of videos because text and images are sufficient to outline the content for specific readers, which also leads to the absence of videos in the articles.
- All the above are random factors that could affect the treatment indicator but are not included in the dataset. Moreover, they are not directly related to shares.

Thank you

Beiming Zhang

BU.450.760.K1.SP24

