

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

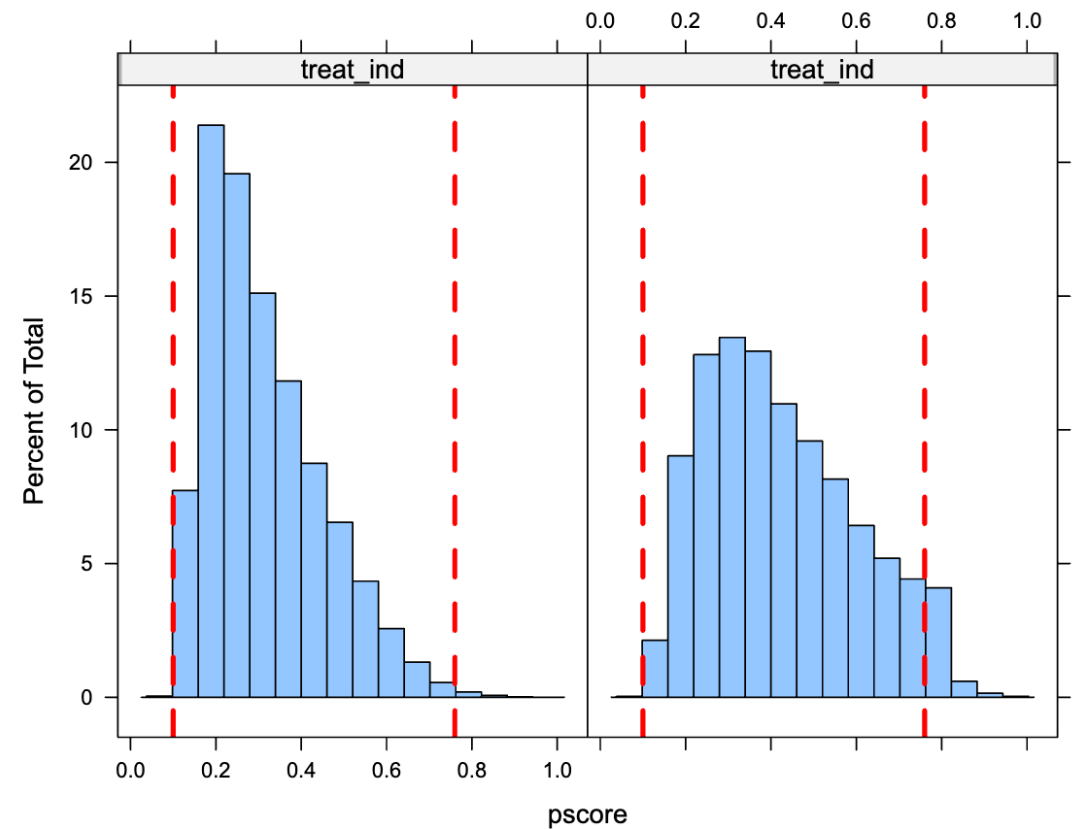
## Task 2 (a)

- 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", "num_videos", "treat_ind")
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", smd = TRUE
10 )
11 print(table_unmatched)
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(0.1, 0.76), col = "red", lty = 2, lwd = 3)
27   }
28 )
```

## Task 2 (a)

- From the overlap diagram, the common overlap between the two is concentrated in the range of  $0.1 \leq n \leq 0.76$ , 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 ( $0.1 \leq p \leq 0.76$ ).
- In the result, there are 13462 values corresponding to 0 and 12518 values corresponding to 1, which is basically balanced.

```
1 # Task 2b
2 # Perform matching
3 matched <- matchit(y_x, method = "nearest", data = ds)
4
5 # Filter based on the results of pscore overlap
6 ds_matched <- match.data(matched)
7 ds_matched <- ds_matched
8   [ds_matched$pscore ≥ .1 & ds_matched$pscore ≤ .76, ]
9 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 the following 2 pages), I believe the matching procedure has been successful, as it effectively increased the differences in covariances corresponding to the samples. Although there are still some variables with small differences, this is due to the inherently small differences among the datasets themselves.

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



# Task 2 (c) (table\_unmatched)

```

1 > print(table_unmatched)
2
3
4 n 25026
5 n_tokens_title (mean (SD)) 10.21 (2.10)
6 n_tokens_content (mean (SD)) 554.98 (449.36)
7 n_unique_tokens (mean (SD)) 0.56 (4.43)
8 n_non_stop_unique_tokens (mean (SD)) 0.70 (4.11)
9 num_hrefs (mean (SD)) 11.04 (11.31)
10 num_self_hrefs (mean (SD)) 3.25 (4.06)
11 num_imgs (mean (SD)) 4.68 (8.12)
12 average_token_length (mean (SD)) 4.64 (0.63)
13 num_keywords (mean (SD)) 7.19 (1.94)
14 kwshares_worst (mean (SD)) 318.09 (510.46)
15 kwshares_best (mean (SD)) 236112.80 (121726.01)
16 kwshares_avg (mean (SD)) 2969.13 (1164.23)
17 self_reference_min_shares (mean (SD)) 4027.48 (21038.58)
18 self_reference_max_shares (mean (SD)) 8906.44 (35549.18)
19 self_reference_avg_sharess (mean (SD)) 5873.04 (23925.79)
20 global_subjectivity (mean (SD)) 0.44 (0.10)
21 global_sentiment_polarity (mean (SD)) 0.12 (0.09)
22 global_rate_positive_words (mean (SD)) 0.04 (0.02)
23 global_rate_negative_words (mean (SD)) 0.02 (0.01)
24 rate_positive_words (mean (SD)) 0.70 (0.17)
25 rate_negative_words (mean (SD)) 0.29 (0.15)
26 avg_positive_polarity (mean (SD)) 0.36 (0.09)
27 min_positive_polarity (mean (SD)) 0.09 (0.07)
28 max_positive_polarity (mean (SD)) 0.76 (0.23)
29 avg_negative_polarity (mean (SD)) -0.25 (0.12)
30 min_negative_polarity (mean (SD)) -0.51 (0.28)
31 max_negative_polarity (mean (SD)) -0.11 (0.09)
32 title_subjectivity (mean (SD)) 0.27 (0.32)
33 title_sentiment_polarity (mean (SD)) 0.07 (0.26)
34 abs_title_subjectivity (mean (SD)) 0.34 (0.19)
35 abs_title_sentiment_polarity (mean (SD)) 0.15 (0.22)
36 shares (mean (SD)) 2891.47 (6524.88)
37 category (%)
38 business 4893 (19.6)
39 entertainment 3152 (12.6)
40 lifestyle 1631 ( 6.5)
41 socialmedia 1642 ( 6.6)
42 tech 5275 (21.1)
43 world 8433 (33.7)
44 weekday (%)
45 friday 3561 (14.2)
46 monday 4164 (16.6)
47 saturday 1677 ( 6.7)
48 sunday 1791 ( 7.2)
49 thursday 4619 (18.5)
50 tuesday 4541 (18.1)
51 wednesday 4673 (18.7)
52 pscore (mean (SD)) 0.31 (0.13)

```

```

1
2 Stratified by treat_ind
3 1 p test
4 n 13686
5 n_tokens_title (mean (SD)) 10.69 (2.10) <0.001
6 n_tokens_content (mean (SD)) 526.87 (501.64) <0.001
7 n_unique_tokens (mean (SD)) 0.53 (0.17) 0.429
8 n_non_stop_unique_tokens (mean (SD)) 0.67 (0.19) 0.286
9 num_hrefs (mean (SD)) 10.61 (11.34) <0.001
10 num_self_hrefs (mean (SD)) 3.34 (3.52) 0.028
11 num_imgs (mean (SD)) 4.20 (8.42) <0.001
12 average_token_length (mean (SD)) 4.39 (1.11) <0.001
13 num_keywords (mean (SD)) 7.30 (1.84) <0.001
14 kwshares_worst (mean (SD)) 302.34 (768.52) 0.016
15 kwshares_best (mean (SD)) 271367.86 (133997.25) <0.001
16 kwshares_avg (mean (SD)) 3259.51 (1453.11) <0.001
17 self_reference_min_shares (mean (SD)) 4022.10 (17739.38) 0.980
18 self_reference_max_shares (mean (SD)) 12945.53 (49804.74) <0.001
19 self_reference_avg_sharess (mean (SD)) 7433.21 (25319.44) <0.001
20 global_subjectivity (mean (SD)) 0.44 (0.14) 0.018
21 global_sentiment_polarity (mean (SD)) 0.11 (0.10) <0.001
22 global_rate_positive_words (mean (SD)) 0.04 (0.02) <0.001
23 global_rate_negative_words (mean (SD)) 0.02 (0.01) <0.001
24 rate_positive_words (mean (SD)) 0.65 (0.22) <0.001
25 rate_negative_words (mean (SD)) 0.29 (0.17) 0.097
26 avg_positive_polarity (mean (SD)) 0.35 (0.12) <0.001
27 min_positive_polarity (mean (SD)) 0.10 (0.08) <0.001
28 max_positive_polarity (mean (SD)) 0.74 (0.28) <0.001
29 avg_negative_polarity (mean (SD)) -0.27 (0.14) <0.001
30 min_negative_polarity (mean (SD)) -0.53 (0.30) <0.001
31 max_negative_polarity (mean (SD)) -0.11 (0.11) 0.003
32 title_subjectivity (mean (SD)) 0.31 (0.34) <0.001
33 title_sentiment_polarity (mean (SD)) 0.07 (0.28) 0.760
34 abs_title_subjectivity (mean (SD)) 0.34 (0.19) 0.068
35 abs_title_sentiment_polarity (mean (SD)) 0.17 (0.24) <0.001
36 shares (mean (SD)) 4310.18 (17476.82) <0.001
37 category (%) <0.001
38 business 1365 (10.0)
39 entertainment 2973 (21.7)
40 lifestyle 468 ( 3.4)
41 socialmedia 681 ( 5.0)
42 tech 2071 (15.1)
43 world 6128 (44.8)
44 weekday (%) <0.001
45 friday 2008 (14.7)
46 monday 2313 (16.9)
47 saturday 724 ( 5.3)
48 sunday 853 ( 6.2)
49 thursday 2505 (18.3)
50 tuesday 2683 (19.6)
51 wednesday 2600 (19.0)
52 pscore (mean (SD)) 0.43 (0.18) <0.001

```



# Task 2 (c) (table\_matched)

```
1 > print(table_matched)
2
3
4 n 13462
5 n_tokens_title (mean (SD)) 10.56 (2.12)
6 n_tokens_content (mean (SD)) 554.22 (452.40)
7 n_unique_tokens (mean (SD)) 0.53 (0.13)
8 n_non_stop_unique_tokens (mean (SD)) 0.67 (0.14)
9 num_hrefs (mean (SD)) 10.99 (11.17)
10 num_self_hrefs (mean (SD)) 3.37 (4.28)
11 num_imgs (mean (SD)) 4.60 (7.08)
12 average_token_length (mean (SD)) 4.55 (0.78)
13 num_keywords (mean (SD)) 7.29 (1.95)
14 kwshares_worst (mean (SD)) 296.39 (551.56)
15 kwshares_best (mean (SD)) 254773.98 (114739.45)
16 kwshares_avg (mean (SD)) 3104.75 (1220.46)
17 self_reference_min_shares (mean (SD)) 4098.66 (18995.66)
18 self_reference_max_shares (mean (SD)) 9217.70 (31513.51)
19 self_reference_avg_sharess (mean (SD)) 5996.02 (20855.77)
20 global_subjectivity (mean (SD)) 0.45 (0.11)
21 global_sentiment_polarity (mean (SD)) 0.12 (0.10)
22 global_rate_positive_words (mean (SD)) 0.04 (0.02)
23 global_rate_negative_words (mean (SD)) 0.02 (0.01)
24 rate_positive_words (mean (SD)) 0.68 (0.18)
25 rate_negative_words (mean (SD)) 0.30 (0.15)
26 avg_positive_polarity (mean (SD)) 0.36 (0.10)
27 min_positive_polarity (mean (SD)) 0.10 (0.07)
28 max_positive_polarity (mean (SD)) 0.76 (0.24)
29 avg_negative_polarity (mean (SD)) -0.27 (0.13)
30 min_negative_polarity (mean (SD)) -0.54 (0.29)
31 max_negative_polarity (mean (SD)) -0.11 (0.10)
32 title_subjectivity (mean (SD)) 0.29 (0.33)
33 title_sentiment_polarity (mean (SD)) 0.07 (0.27)
34 abs_title_subjectivity (mean (SD)) 0.34 (0.19)
35 abs_title_sentiment_polarity (mean (SD)) 0.16 (0.23)
36 shares (mean (SD)) 3039.33 (6158.58)
37 category (%)
38 business 1312 ( 9.7)
39 entertainment 2676 (19.9)
40 lifestyle 484 ( 3.6)
41 socialmedia 730 ( 5.4)
42 tech 2275 (16.9)
43 world 5985 (44.5)
44 weekday (%)
45 friday 1948 (14.5)
46 monday 2257 (16.8)
47 saturday 765 ( 5.7)
48 sunday 880 ( 6.5)
49 thursday 2480 (18.4)
50 tuesday 2593 (19.3)
51 wednesday 2539 (18.9)
52 pscore (mean (SD)) 0.39 (0.12)
```

```
1
2 Stratified by treat_ind
3 1 p test
4 n 12518
5 n_tokens_title (mean (SD)) 10.62 (2.09) 0.011
6 n_tokens_content (mean (SD)) 554.50 (487.85) 0.962
7 n_unique_tokens (mean (SD)) 0.55 (0.12) <0.001
8 n_non_stop_unique_tokens (mean (SD)) 0.70 (0.13) <0.001
9 num_hrefs (mean (SD)) 11.14 (11.17) 0.271
10 num_self_hrefs (mean (SD)) 3.48 (3.39) 0.022
11 num_imgs (mean (SD)) 4.36 (8.05) 0.011
12 average_token_length (mean (SD)) 4.60 (0.56) <0.001
13 num_keywords (mean (SD)) 7.29 (1.87) 0.954
14 kwshares_worst (mean (SD)) 301.08 (684.95) 0.543
15 kwshares_best (mean (SD)) 257748.39 (125972.64) 0.046
16 kwshares_avg (mean (SD)) 3132.25 (1266.04) 0.075
17 self_reference_min_shares (mean (SD)) 4136.02 (18164.96) 0.871
18 self_reference_max_shares (mean (SD)) 10691.08 (31289.92) <0.001
19 self_reference_avg_sharess (mean (SD)) 6657.88 (19952.44) 0.009
20 global_subjectivity (mean (SD)) 0.46 (0.10) <0.001
21 global_sentiment_polarity (mean (SD)) 0.12 (0.10) 0.040
22 global_rate_positive_words (mean (SD)) 0.04 (0.02) <0.001
23 global_rate_negative_words (mean (SD)) 0.02 (0.01) <0.001
24 rate_positive_words (mean (SD)) 0.69 (0.17) <0.001
25 rate_negative_words (mean (SD)) 0.30 (0.16) 0.016
26 avg_positive_polarity (mean (SD)) 0.37 (0.10) <0.001
27 min_positive_polarity (mean (SD)) 0.10 (0.07) 0.001
28 max_positive_polarity (mean (SD)) 0.78 (0.23) <0.001
29 avg_negative_polarity (mean (SD)) -0.28 (0.13) <0.001
30 min_negative_polarity (mean (SD)) -0.56 (0.29) 0.001
31 max_negative_polarity (mean (SD)) -0.11 (0.10) 0.005
32 title_subjectivity (mean (SD)) 0.30 (0.33) 0.076
33 title_sentiment_polarity (mean (SD)) 0.07 (0.28) 0.334
34 abs_title_subjectivity (mean (SD)) 0.34 (0.19) 0.925
35 abs_title_sentiment_polarity (mean (SD)) 0.17 (0.23) 0.039
36 shares (mean (SD)) 3353.21 (6213.86) <0.001
37 category (%) 0.002
38 business 1305 (10.4)
39 entertainment 2704 (21.6)
40 lifestyle 462 ( 3.7)
41 socialmedia 674 ( 5.4)
42 tech 2056 (16.4)
43 weekday (%) 0.995
44 friday 1831 (14.6)
45 monday 2118 (16.9)
46 saturday 697 ( 5.6)
47 sunday 799 ( 6.4)
48 thursday 2296 (18.3)
49 tuesday 2424 (19.4)
50 wednesday 2353 (18.8)
51 pscore (mean (SD)) 0.40 (0.14) <0.001
```

## 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 313.84 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  -3347  -2248  -1734   -348  141366
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept)  3034.49      53.09  57.154 < 2e-16 ***
13 treat_ind     313.84      76.48   4.104 4.08e-05 ***
14 ---
15 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16
17 Residual standard error: 6174 on 26096 degrees of freedom
18 Multiple R-squared:  0.0006449, Adjusted R-squared:  0.0006066
19 F-statistic: 16.84 on 1 and 26096 DF,  p-value: 4.081e-05
```

## 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. The relevance between video content and article content. If video content is highly relevant to the article and complements it, then the likelihood of shares may increase, as it reflects the overall quality of the article.
  2. The total length of the video. The length of a video generally represents the amount of information in the article; if a video is longer, it typically contains more information and is more likely to be considered "valuable" and shared.
  3. The platform where the video is located. Videos on X are usually shorter and often hastily shot, with lower video quality; whereas videos on YouTube might have higher production quality and therefore more opportunities for shares.

# Thank you

Beiming Zhang

BU.450.760.K1.SP24

