



## TZ Gaming: Optimal Targeting of Mobile Ads

Prof. Hema Yoganarasimhan, Foster School of Business, University of Washington  
Prof. Vincent Nijs, Rady School of Management, UCSD

Winter 2020

As a developer of games for mobile devices TZ gaming has achieved strong growth of its customer base. A prominent source of new customers has come from ads displayed through the Vneta ad-network. A mobile-ad network is a technology platform that serves as a broker between app developers (or publishers) looking to sell ad space and a group of advertisers.

App developers sell “impressions”, i.e., a space where an ad can be shown, through the Vneta network to companies such as TZ gaming looking to advertise to app users.

TZ uses the ads to appeal to prospective customers for their games. They generally use short (15 sec) video ads that help to emphasize the dynamic nature of the games. In the past, TZ has been able to, approximately, break-even on ad-spend with Vneta when calculating the benefits that can be directly attributed to ad click-through. TZ, however, believes there are additional, longer-term, benefits from these ads such as brand awareness, etc. that are harder to quantify.

Currently, Vneta provides only very limited targeting of ads to app users but is planning to start offering behavioral targeting to advertisers for a fee. Specifically, two options are under consideration: (a) Provide access to data that advertisers can use to determine which impressions they want to bid on or (b) Advertisers pay Vneta a data science consultancy fee to conduct data-driven targeting on the advertisers behalf.

As Vneta is developing this new business model, it has decided to work with TZ as a partner and has shared behavioral information linked to 115,488 recent impressions used to show TZ ads. Vneta have also provided a set of predictions based on their own (proprietary) algorithm that they intend to use as part of their data science consulting service.

Matt Huateng, the CEO of TZ games, is intrigued by the potential for data science to enhance the efficiency of targeted advertising on mobile devices. However, he is not convinced that the consulting services offered by Vneta will be worth the money. He has asked you to do some initial analyses on the provided data and compare the generated predictions to Vneta’s recommendations. The following three options need to be evaluated to determine the best path forward.

Options:

1. No targeting (i.e., continue with the current approach)
2. Use predictions from a logistic regression model for ad targeting
3. Use predictions generated by Vneta for ad targeting

The assumptions used for the analysis are as follows:

- Targeting of impressions to consumers covered by the Vneta ad-network to date has been random
- Cost per 1,000 video impressions (CPM) is \$10
- Conversion to sign-up as a TZ game player after clicking on an ad is 5%
- The expected CLV of customers that sign-up with TZ after clicking on an ad is approximately \$25
- The price charged for the data by Vneta is \$50K
- The price charged for the data science consulting services by Vneta is \$150K

Approach:

- Use the 87,535 rows in the data with “training == ‘train’” to estimate a model. Then generate predictions for all 115,488 rows in the dataset
- Options 1-3 should be evaluated *only* on the predictions generated for the 27,953 rows in the data with “training == ‘test’”. These are the observations that were *not* used to estimate your model
- Extrapolate the cost and benefits for options 1-3 above for an upcoming advertising campaign where TZ will purchase 20-million impressions from Vneta

Although TZ gaming has used RFM for targeting existing customers this approach is not appropriate for prospective customers. Instead, you have decided to use logistic regression. This is a powerful and widely used tool to model consumer response. It is similar to linear regression but the key difference is that the response variable is binary (e.g., click or no-click) rather than continuous. For each impression, the logistic regression model will predict the probability of click-through, which can be used for ad targeting. Like linear regression, you can include both continuous and categorical predictors in your model as explanatory variables.

Matt is eager to assess the value of logistic regression as a method to predict ad click-through and target prospects and has asked you to complete the following analyses.

## Part I: Logistic Regression (10 points)

- a. Estimate a logistic regression model using `click` as the response variable and the following as explanatory variables:

`impup, clup, ctrup, impua, clua, ctrua, imput, clut, ctrut, imppat, clpat, ctrpat`

The model should predict the probability of `click == "yes"`. See <https://radiant-rstats.github.io/docs/model/logistic.html> and <http://radiant-rstats.github.io/radiant.model/reference/logistic.html> for details. Create a new variable called `click_logit` with the predicted click-through probability linked to each impression. You can, of course, also use the `glm` function in R but you will have to do more work to get all the required output. In python, please use `statsmodels` estimate the model (see structure below).

```
import statsmodels.formula.api as smf
tz_gaming["click_yes"] = (tz_gaming["click"] == "yes").astype(int)
lr = smf.logit(
    formula="click_yes ~ impup + clup + ctrup + ...",
    data=tz_gaming.query("training == 'train'")
).fit(maxiter=1000, method='bfgs')
lr.summary()
```

- b. Summarize and interpret the logistic regression results. Which explanatory variables are statistically significant? Which variables seem to be most “important”? Make sure your model evaluation includes (1) an interpretation of the odds-ratios estimated for each of the explanatory variables and (2) an evaluation of the model as a whole.

## Part II: Decile Analysis of Logistic Regression Results (10 points)

Note: For the following questions, use only the “test” sample of impressions (i.e., 27,953 rows where “training == ‘test’”)

- Assign each impression to a decile based on the predicted probability of click through. Create a new variable `dec_logit` that captures this information. Note: The first decile should have the highest average click-through rate. If not, make sure to “reverse” the decile numbers (i.e., 10 becomes 1, 9 becomes 2, etc.). Please use the `xtile` function to create the deciles.
- Create a bar chart of click-through rates per decile (i.e., use `dec_logit` as the x-variable and `click == "yes"` as the y-variable). Note that the “click through rate” is not the same as the “predicted probability of click.” The click through rate captures the proportion of impressions in a given group (e.g., in a decile) that actually resulted in a click.
- Report the number of impressions, the number of clicks, and the click-through rate for the TZ ad per decile and save this information to a dataframe. Use the name `dec_df_logit` for the new data frame.
- Estimate a logistic regression model with `click` as the response variable and `imppat`, `clpat`, and `ctrpat` as the only explanatory variable. Make sure to “standardize” the explanatory variables before estimation. What is the interpretation of the standardized odds-ratios for the explanatory variables?
- Some of the variables in the dataset are highly correlated with each other. In particular, `imppat` and `clpat` have a positive correlation of 0.97. Discuss the implications of this (very) high level of collinearity and also different approaches to deal with it. What are the implications for the model and the interpretation of the estimated (standardized) coefficients? As part of your answer, discuss the change in the estimated (standardized) odd-ratio for `imppat` when you remove `clpat` from the model.
- Estimate another logistic regression model with `click` as the response variable and `time_fct`, `app`, `imppat`, `clpat`, and `ctrpat` as the explanatory variable. Why are the odds ratios for `imppat`, `clpat`, and `ctrpat` different in the two models? Please be specific and investigate beyond simply stating the statistical problem.

## Part III: Lift and Gains (5 points)

Note: For the following questions, use only the “test” sample of impressions (i.e., 27,953 rows where “training == ‘test’”)

- Use the data frame you created in II.c above to generate a table with lift and cumulative lift numbers for each decile
- Create a ggplot (or matplotlib or altair) chart showing the cumulative lift per decile
- Use the data frame you created in II.c above to generate a table with gains and cumulative gains numbers for each decile
- Create a ggplot (or matplotlib or altair) chart showing the cumulative gains per decile along with a reference line for the “no model”

Note: Do not use any specialized packages to construct the lift and gains tables and charts

## Part IV: Confusion matrix (5 points)

- Create a “confusion matrix” based on the predictions from the logistic regression model you estimated above for I.a. Again, use **only** data from the test set here (i.e., “training == ‘test’”). Use the financial assumptions mentioned above, and repeated in section V below, to determine an appropriate cut-off (i.e., break-even). Calculate “accuracy” based on the confusion matrix you created (see [http://lab.rady.ucsd.edu/sawtooth/RBusinessAnalytics/logit\\_models.html](http://lab.rady.ucsd.edu/sawtooth/RBusinessAnalytics/logit_models.html) for an example using R).

Note: Do not use any specialized packages to construct the confusion matrix

- b. Calculate a confusion matrix based on predictions from a logistic regression with `click` as the response variable and `rnd` as the **only** explanatory variable. As before, the model should be estimated on training sample (i.e., “training == ‘train’”). Generate predictions for all rows in the data and create the confusion matrix based only on the test set (i.e., “training == ‘test’”). Calculate “accuracy” based on the confusion matrix you created.
- c. Discuss the similarities and differences between the two confusion matrices. Which model is best based on the confusion matrix? Provide support for your conclusions.
- d. Recalculate the confusion matrices from IV.a and IV.b using 0.5 as the cutoff. Based on these new matrices, discuss again the similarities and differences. Which model is best based on the confusion matrix? Provide support for your conclusions.

## Part V: Profitability Analysis (5 points)

Use the following cost information to assess the profitability of using the logistic regression model from I.a for targeting purposes during the upcoming advertising campaign where TZ will purchase 20-million impressions from Vneta:

- Cost per 1,000 video impressions (CPM) is \$10
  - Conversion to sign-up as a TZ game player after clicking on an ad is 5%
  - The expected CLV of customers that sign-up with TZ after clicking on an ad is approximately \$25
  - The price charged for the data by Vneta is \$50K
  - The price charged for the data science consulting services by Vneta is \$150K
- a. Create a new variable `target_logit` that is TRUE if the predicted click-through probability is greater than the break-even response rate you calculated in IV.a and FALSE otherwise
  - b. For the test set (i.e., “training == ‘test’”), what is the expected profit (in dollars) and the expected return on marketing expenditures (ROME) if TZ used (1) no targeting, (2) purchased the data from Vneta and used the logistic regression from I.a for targeting, or (3) used Vneta’s data science consulting services? You can use the `click_vneta` variable to create a `target_vneta` variable and calculate the expected profit and the expected return on marketing expenditures

Note: To estimate the performance implications of “no targeting” approach use the predictions from the model you estimated in IV.b

- c. Predict the profit and ROME implications for each of the 3 options if TZ purchases 20-million impression for the upcoming ad campaign? Use the results from (b) above to project the performance implications

Note: The currently available data (+ the `click_vneta` prediction) are free as part of the partnership between Vneta and TZ-gaming. You should assume that the cost of the data (50K) and the consulting fee (150K) would apply for the 20M impression campaign.

## Part VI: Model comparison (10 points)

- a. The calculations in V.a through V.c above assume that the predicted probabilities are estimated without error. Calculate the confidence interval for the predictions from the logistic regression model in I.a. Now redo the calculations from V.a through V.c, for only the logistic regression model in I.a, adjusting for these errors. How do the results change? Example code using the `logistic` function from the `radiant.model` package are shown below:

```
result <- logistic(tz_gaming, ...)
pred <- predict(result, pred_data = tz_gaming, conf_lev = 0.9, se = TRUE)
```

You can add the columns you need from the “pred” data frame to your data set or use the `store` function in Radiant to add them:

```
tz_gaming <- store(
  tz_gaming, pred,
  name = c("click_logit", "click_logit_lb", "click_logit_ub")
)
```

In python, use the formula api for `statsmodels` and `predict_conf_int` from the `pyrsm` package (see <https://github.com/vnijs/pyrsm/blob/master/pyrsm/logit.py>). To install the latest version of the `pyrsm` package use `pip3 install --user pyrsm` from a terminal in the docker container. You can check the version number of the `pyrsm` package by using `import pyrsm` followed by `pyrsm.__version__`. Your version should be at least 0.1.9. The basic structure of your code to estimate a logistic regression in python should be as follows:

```
import statsmodels.formula.api as smf
from pyrsm import predict_conf_int
tz_gaming["click_yes"] = (tz_gaming["click"] == "yes").astype(int)
result = smf.logit(formula="... ~ ...", ...).fit(maxiter=1000, method='bfgs')
result.summary()
pred = predict_conf_int(result, df = tz_gaming, alpha = 0.1)
```

Also create a variable `target_logit_lb` that is TRUE if the predicted click-through probability is greater than the break-even response rate and FALSE otherwise.

- b. The calculations in V.b through V.d above are based on a model that did not include all available variables. Not all variables may be relevant however. To at least give each variable available in the dataset a chance of being included in the model, estimate a (“backward”) stepwise logistic regression model, starting with the following variables:

```
time_fct, app, impup, clup, ctrup, impua, clua, ctrua, imput, clut, ctrut, imppat, clpat,
ctrpat
```

Create a variable (`click_logit_stepwise`) with predicted click-through probabilities from this model. Also create a variable (`target_logit_stepwise`) that is TRUE if the predicted click-through probability is greater than the break-even response rate and FALSE otherwise.

Note: Python does not have any tools for stepwise regression so please use the `click_logit_stepwise_pre` variable in `tz_gaming.pkl` instead. Note that this variable is also available in `tz_gaming.rds`

- c. You have now estimated 4 different models and also have the predictions from Vneta (see prediction labels below). Compare the models using (1) profit calculations as in V.a through V.c and (2) a gains chart. Discuss which of these 5 models you would recommend to put into production and why.

Prediction labels to use: `click_logit`, `click_rnd`, `click_logit_step`, `click_logit_lb`, `click_vneta`

Note: For efficiency, you should adapt the `perf_calc` function you created for the Tuango case so you can use it to do the relevant performance calculations for the different models (i.e., profit, click-through rate, ROME, etc.).

## Data description

Information about the sample of 115,488 impressions is in the R dataset `tz_gaming.rds` (or python dataset `tz_gaming.pkl`) in the `data/R` (`data/python`) directory in the GitLab repo. Each row in the dataset represents an impression. For each row (impression), we have data on 22 variables. All explanatory variables are created by Vneta based on one month tracking history of users, apps, and ads. The available variables are described below.

- *training* – Dummy variable that splits the dataset into a training (“train”) and a test (“test”) set
- *inum* – Impression number
- *click* – Click indicator for the TZ ad served in the impression. Equals “yes” if the ad was clicked and “no” otherwise
- *time* – The hour of the day in which the impression occurred (1-24). For example, “2” indicates the impression occurred between 1 am and 2 am
- *time\_fct* – Same as *time* but the is coded as categorical
- *app* – The app in which the impression was shown. Ranges from 1 to 49
- *id* – Anonymized user ID
- *impup* – Number of past impressions the user has seen in the app
- *clup* – Number of past impressions the user has clicked on in the app
- *ctrup* – Past CTR (Click-Through Rate) (x 100) for the user in the app
- *impua* – Number of past impressions of the TZ ad that the user has seen across all apps
- *clua* – Number of past impressions of the TZ ad that the user has clicked on across all apps
- *ctrua* – Past CTR (x 100) of the TZ ad by the user across all apps
- *imput* – Number of past impressions the user has seen within in the hour
- *clut* – Number of past impressions the user has clicked on in the hour
- *ctrut* – Past CTR (x 100) of the user in the hour
- *imppat* – Number of past impressions that showed the TZ ad in the app in the hour
- *clpat* – Number of past clicks the TZ ad has received in the app in the hour
- *ctrpat* – Past CTR (x 100) of the TZ ad in the app in the hour
- *rnd* – Simulated data from a normal distribution with mean 0 and a standard deviation of 1
- *click\_vneta* – Predicted probability of click per impressions generated by Vneta’s proprietary machine learning algorithm
- *click\_logit\_step\_pre* – Predicted probability of click per impressions generated by a stepwise logistic regression

The last three letters of a feature name indicate the sources of variation:

- u — denotes user
- t — denotes time
- p — denotes app
- a — denotes ad

Note that there is a clear relationship between the impressions, clicks, and ctr variables within a strata. Specifically:  $ctrup = \frac{clup}{impup}$ ,  $ctr_u = \frac{cl_u}{imp_u}$ ,  $ctrut = \frac{clut}{imput}$ , and  $ctrpat = \frac{clpat}{impat}$ .

---

Professor Hema Yoganarasimhan (Foster School of Business, University of Washington) and Professor Vincent Nijs, Rady School of Management, UCSD prepared this case to provide material for class discussion rather than to illustrate either effective or ineffective handling of a business situation. Names and data have been disguised to assure confidentiality. Copyright (c) 2020 by Hema Yoganarasimhan and Vincent Nijs