



TZ Gaming: Optimal Targeting of Mobile Ads

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

Winter 2024

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 (1) app developers (or publishers) looking to sell ad space and (2) 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. Vneta acts as a broker for 50-60 millions impressions/ads per day.

TZ gaming uses 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. Many senior executives at TZ believe that there are additional, longer-term, benefits from these ads such as brand awareness, etc. that are harder to quantify.

Currently, TZ has access to very limited data from Vneta. Matt Huateng, the CEO of TZ gaming, is intrigued by the potential for data science to enhance the efficiency of targeted advertising on mobile devices. Specifically, two options are under consideration: (1) Buy access to additional data from Vneta and use TZ’s analytics team to build targeting models or (2) Subscribe to Vneta’s analytics consultancy service, which provides impression-level click-through rate predictions based on Vneta’s proprietary data and algorithms.

Vneta has shared behavioral information linked to 115,488 recent impressions used to show TZ ads and has also provided a set of predictions based on their own (proprietary) algorithm. Matt is not convinced that the consulting services offered by Vneta will be worth the money for future ad campaigns and has asked you to do some initial analyses on the provided data and compare the generated predictions to Vneta’s recommendations. The following options will be evaluated to determine the best path forward.

Options:

1. Spam all prospects
2. Continue with the current targeting approach
3. Use predictions from a logistic regression model for ad targeting
4. 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 (approximately) 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 different models. Then generate predictions for all 115,488 rows in the dataset
- Options 1-4 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 1-4 above for an upcoming advertising campaign where TZ will commit to purchase 20-million impressions from Vneta

TZ gaming has decided to use logistic regression for targeting. 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 (target) 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 (features).

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)

Note: For the following questions, use only the “training” sample of impressions (i.e., 87,535 rows where “training == ‘train’ ”).

- Estimate a logistic regression model using `click` as the response variable (target) and the following as explanatory variables (features). The model should predict the probability that `click` is equal to “yes” (2 points):

```
time_fct, app, mobile_os, impua, clua, ctrua
```

```
lr =
```

- Summarize and interpret the logistic regression results. Which of these explanatory variables are statistically significant? Which variables seem to be most “important”? Make sure your model evaluation includes (1) an interpretation of the Permutation importance and Prediction plots for the explanatory variables `mobile_os`, `impua`, `clua`, and `ctrua` and (2) an evaluation of the model as a whole using Pseudo R-squared and the Chi-square test (5 points).
- Predict the probability of a click (2 points)

The estimated logistic regression model can predict the probability of a click. Create a new variable `pred_logit` with the predicted click probabilities for each impression. Make sure to generate predictions for all rows in both the training and test data.

```
tz_gaming["pred_logit"] =
```

- d. Estimate a logistic regression with `click` as the response variable and `rnd` as the **only** explanatory variable. As before, the model should be estimated on the training sample (i.e., “training == ‘train’”). Create a new variable `pred_rnd` with the predicted click-through probabilities (1 point).

```
lr_rnd =  
tz_gaming["pred_rnd"] =
```

Part II: Understanding Multicollinearity and Omitted Variable Bias (10 points)

- a. Estimate a logistic regression model with `click` as the response variable and `imppat`, `clpat`, and `ctrpat` as the only explanatory variables. What is the interpretation of the Prediction plots for the explanatory variables? (2 points)

```
lr_mc1 =
```

Note: Make sure to watch the “Video: TZ gaming preview (12 min)” on Canvas before answering this questions so you fully understand what the variables represent

- b. Some of the variables in the dataset are highly correlated with each other. In particular, `imppat` and `clpat` have a very high 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 Prediction plots? As part of your answer, discuss the change in the Prediction plot for `imppat` when you remove `clpat` from the model you estimated for II.a (4 points).

Note: Assign your new model without `clpat` to a new object `lr_mc2`. Calculate VIF statistics for each explanatory variable in the model

```
lr_mc2 =
```

- c. 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 Prediction plots for `imppat`, `clpat`, and `ctrpat` different compared to the plots from the model you estimated in II.a? Please be specific and investigate beyond simply stating the statistical problem (4 points).

Note: You may want to test if a (set of) coefficients are equal to 0 (or Odds-ratios are equal to 1)

```
lr_mc3 =
```

Part III: Decile Analysis of Logistic Regression Results (5 points)

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

- a. Assign each impression to a decile based on the predicted probability of a click (`pred_logit`) based on the model estimated in I.a. Create a new variable `pred_logit_dec` that captures this information. Note: The first decile should have the highest average click rate. If not, make sure to “reverse” the decile numbers (i.e., 10 becomes 1, 9 becomes 2, etc.). Use the `xtile` function from the `pyrsm` package to create the deciles (2 points)

```
tz_gaming["pred_logit"] =
```

- b. Report the number of impressions (rows), the number of clicks (`click`), and the click through rate (`ctr`) (i.e., sum of clicks divided by number of impressions) for the TZ ad per decile and save this information to a new dataframe called `dec_tab` (2 points)
- c. Create a bar chart of click-through rates per decile (i.e., use `pred_logit_dec` as the x-variable and `ctr` 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 (1 point)

Part IV: Gains Curves (15 points)

Use the `dec_tab` DataFrame you created in Part III for the following calculations.

- a. Write python code to generate a table with the cumulative proportion of impressions and the cumulative gains for each decile (8 points)

Note: Do NOT use any specialized python packages to construct the gains table. Write the python code from scratch. Feel free use ChatGPT or CoPilot, but make sure that it does not use any specialized packages to construct the gains table. Be prepared to discuss the code you submit for this question in class if called upon

- b. Use `seaborn`, `matplotlib`, or `pandas` to create a chart showing the cumulative gains per decile along with a (diagonal) reference line to represent the “no model” scenario. Put cumulative gains on the Y-axis and cumulative proportion of impressions on the X-axis (7 points)

Note: Do NOT use any specialized packages to construct the gains chart. Write the python code from scratch. Feel free use ChatGPT or CoPilot, but make sure that it does not use any specialized packages to construct the gains table. Be prepared to discuss the code you submit for this question in class if called upon

Part V: Confusion matrix (10 points)

- a. Create a “confusion matrix” based on the predictions from the logistic regression model you estimated in Part I.a (i.e., the model used to generate `pred_logit`). Again, use **only** data from the test set here (i.e., “training == ‘test’”). Use the financial assumptions mentioned above, and repeated in section VI below, to determine an appropriate cut-off (i.e., breakeven). Calculate “accuracy” based on the confusion matrix you created (2 points)

Note: Do NOT use any specialized packages to construct the confusion matrix. Code the matrix from scratch. Feel free use ChatGPT or CoPilot, but make sure that it does not use any specialized packages to construct the gains table. Be prepared to discuss the code you submit for this question in class if called upon

- b. Calculate a confusion matrix based on `pred_rnd` created in Part I.e and calculate “accuracy” based on the confusion matrix you created (2 points)
- c. Discuss the similarities and differences between the two confusion matrices. Which prediction (model) is best, based on the confusion matrix? Provide support for your conclusions (3 points)
- d. Recalculate the confusion matrices from V.a and V.b using 0.5 as the cutoff. Based on these new matrices discuss the similarities and differences. Which model is best based on these new confusion matrices? Provide support for your conclusions (3 points)

Part VI: Model comparison (12 points)

Use the following cost information to assess the profitability each of these models 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 total cost of the data from Vneta is \$50K
- The total cost charged for the data science consulting services by Vneta is \$150K

Use `pred_logit`, `pred_rnd`, and the predictions from Vneta based on their proprietary model `pred_vneta` to compare model performance.

Note: The currently available data (+ the `pred_vneta` prediction) are free as part of the partnership between Vneta and TZ-gaming

- a. Create a new variable `target_logit` that is `True` if the predicted click-through (`pred_logit`) probability is greater than the break-even response rate and `False` otherwise (1 point)
- b. Create a new variable `target_rnd` that is `True` if the predicted click-through (`pred_rnd`) probability is greater than the break-even response rate and `False` otherwise (1 point)
- c. Create a new variable `target_vneta` that is `True` if the predicted click-through (`pred_vneta`) probability is greater than the break-even response rate and `False` otherwise (1 point)
- d. Based on performance in the test set (i.e, `training == 'test'`), calculate the projected expected profit (in dollars) and the expected return on marketing expenditures (ROME) for the upcoming 20M impression campaign if TZ (1) “spams” everyone, (2) continues to target using their current approach (`pred_rnd`), (3) uses the data from Vneta to build the logistic regression from I (`pred_logit`) for targeting, or (4) used Vneta’s data science consulting services (`pred_vneta`) to select the best prospects out of 20M impressions. (3 points)

Note: Calculate total profits under the assumption that options (3) and (4) are free of charge. Then compare the profit numbers to determine if these options would be worth the expense going forward.

Note: For efficiency, consider adapting the `perf_calc_actual` function you created for the Tuango case to do the relevant performance calculations for the different models.

- e. Based on the results from VI.d discuss which of these 4 approaches you would recommend and why (2 points)
- f. Calculate the profit and ROME implications for each of the 4 options mentioned in VI.d if TZ purchases exactly 20-million impressions for the upcoming ad campaign out of the +500M prospects that Vneta has access to (2 points)

Note: Calculate total profits under the assumption that options (3) and (4) are free of charge. Then compare the profit numbers to determine if these options would be worth the expense going forward.

Note: For efficiency, consider adapting the `perf_calc` function you created for the Tuango case to do the relevant performance calculations for the different models.

- g. Based on the results from VI.f, discuss which of the 4 approaches you would recommend putting into production and why. Is your recommendation different from VI.e? Why (not)? (2 points)

Note: Calculate total profits under the assumption that options (3) and (4) are free of charge. Then compare the profit numbers to determine if these options would be worth the expense going forward.

Note: For efficiency, consider adapting the `perf_calc` function you created for the Tuango case to do the relevant performance calculations for the different models.

Part VII: Generative AI (5 points)

Please describe how you used Generative AI-tools like ChatGPT to support your work on this assignment. Provide pdfs and/or screenshots of your “discussions” with these tools and comment on what things did and did not go well. Also add any questions you may have about the assignment and the support you received from GenAI so we can discuss these topics in class.

Note: No matter how you used Generative AI-tools, you will be expected to fully understand all elements of the assignment. You may be called on in class to walk us through your thought process and how different parts of your code work.

Data description

Data on 115,488 impressions is contained in the `tz_gaming.pkl` file in the `data/` directory of the GitLab repo that will be forked to your account. Each row in the dataset represents an impression that showed a TZ ad. 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 coded as categorical
- *app* – The app in which the impression was shown. Ranges from 1 to 49
- *mobile_os* – Customer’s mobile OS
- *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
- *input* – 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
- *pred_vneta* – Predicted probability of click per impressions generated by Vneta’s proprietary machine learning algorithm

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: $ctrup = \frac{clup}{impup}$, $ctru = \frac{clu}{impu}$, $ctrut = \frac{clut}{imput}$, and $ctrpat = \frac{clpat}{impat}$.

Professor Vincent Nijs, Rady School of Management, UCSD and Professor Hema Yoganarasimhan (Foster School of Business, University of Washington) 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) 2024 by Vincent Nijs and Hema Yoganarasimhan