

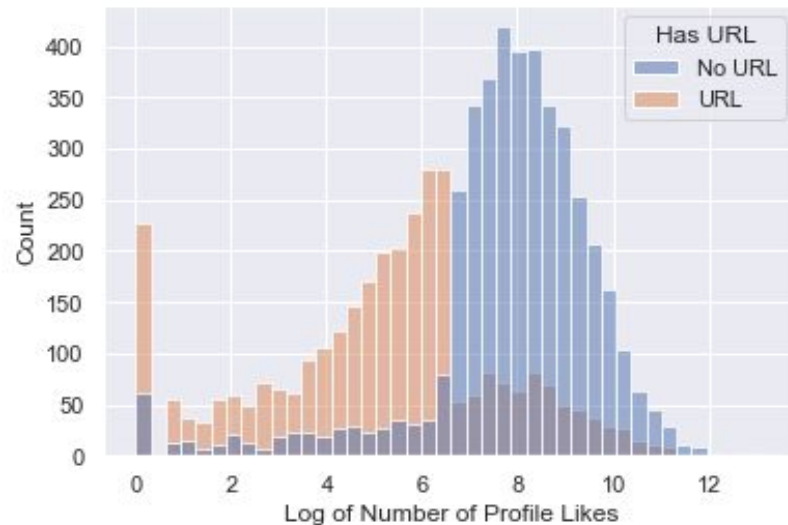
SOCIAL MEDIA PREDICTION KAGGLE COMPETITION - L'ÉQUIPE

Martin Dallaire - 923013
Benjamin Sigman - 1060990
François Paugam - 20169017
François David - 20171906

IFT-6758 Data Science

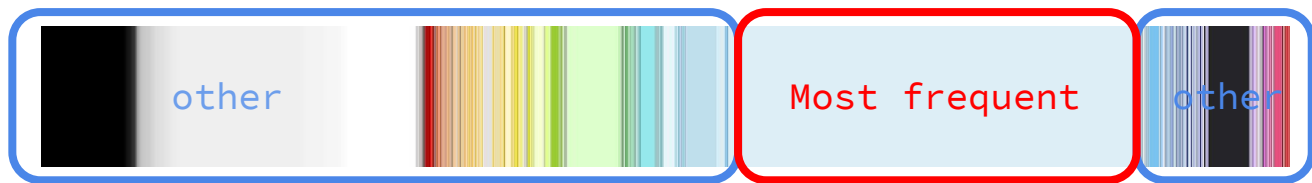
DATA CLEANING/FEATURE ENGINEERING

- Uniformisation (lower case, identify Nan)
- Change URL to “Has Url” column
- Median imputation for Nan



DATA CLEANING/FEATURE ENGINEERING

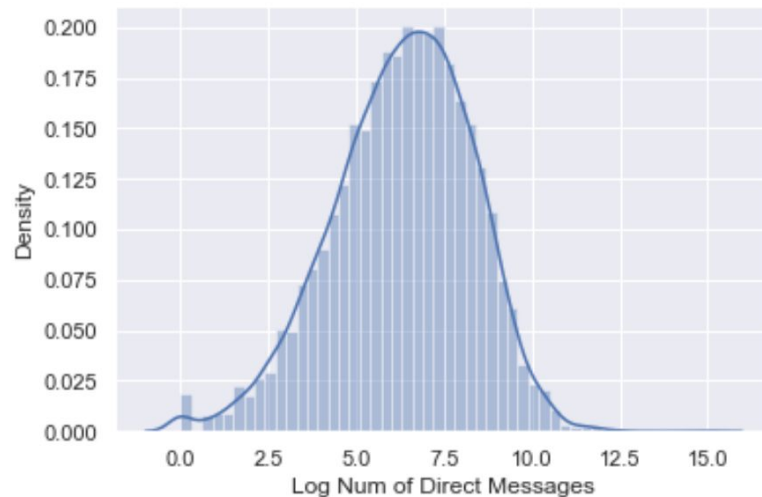
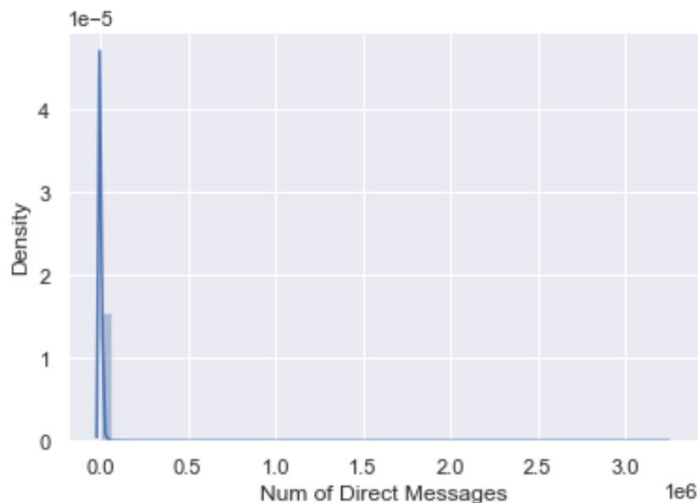
- Keep most frequent categories (that were correlated) and put remaining in “other”



Profile Page color

DATA CLEANING/FEATURE ENGINEERING

- Log transform of numerical features



DATA CLEANING/FEATURE ENGINEERING

- Estimated clicks since creation
- Bag-of-words for Location
- Target mean for categorical data
- Failed creative methods: Sin/Cos on UTC

Keep the features that are correlated to the target

MODEL SELECTION (FIRST STEPS)

- Some(log of)features decently linearly correlated with **log of target values** (**$r > 0.3$**)
 - Linear models (Linear regression, SVM)
- Many categorical features: Try decision tree-based models.
- Ensemble models involving decision trees: Random Forests

MODEL SELECTION (ENSEMBLING)

- Try more tree-based ensembling models:

- LightGBM, great right out of the box. Scores 1.71 by itself on average.

- XGBoost, needed fine tuning. Performed poorly out of the box. After hyperparameter tuning, it performed almost as well as LightGBM.

- Many models performing between 1.71 and 1.85.

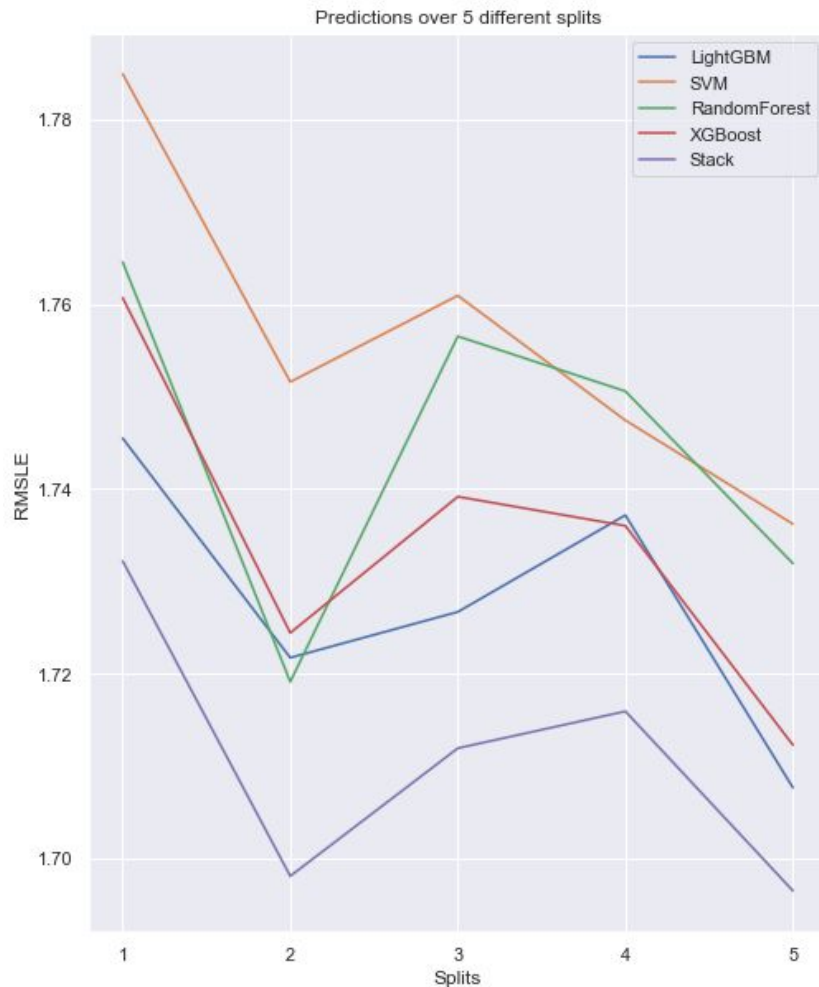
Conclusion: Stacking regressor!

MODEL PERFORMANCE

Similar performances.

LightGBM performs better of all solo models (on average)

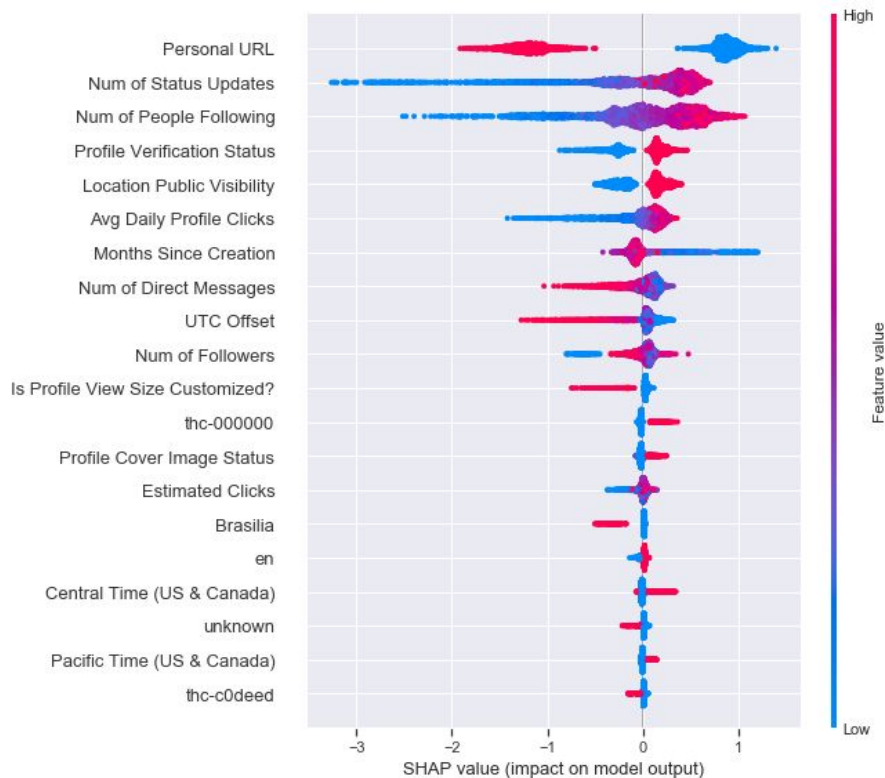
Stack performs best almost every time.



VALIDATION

- 50-fold cross validation
- Hyperparameters tuning
(gridSearchCv)
- Reduce the bias because more data points are considered (compared to val-split)
- Reduces Variances as more data points are used in the validation set.
- Reduces overfitting

SHAPLEY VALUES - INTERPRETABILITY/FEATURE SELECTION



- On the log transformation of the Number of Profile Likes.
- Red means upper range of feature value, blue means the lower range.
- Points on the right contributes to predicting higher values and the ones on the left contributes to predicting lower values.

The image shows a section of a library with three visible shelves. The top shelf is filled with books, mostly with blue and green spines, some with gold lettering. The middle shelf contains a mix of book colors including blue, white, and red, with some books having gold lettering. The bottom shelf also has a variety of book colors, including blue, white, and red, with some books having gold lettering. The books are arranged in a way that shows their spines, with some books standing upright and others lying flat. The background is a plain wall.