

Telecom Customer Churn Analysis

Akshay Abhyankar
University of Colorado Boulder
Boulder, Colorado, USA
Akshay.Abhyankar@colorado.edu

Gaurav Roy
University of Colorado Boulder
Boulder, Colorado, USA
Gaurav.Roy@colorado.edu

Aditi Ramesh Athreya
University of Colorado Boulder
Boulder, Colorado, USA
Aditi.RameshAthreya@colorado.edu

Varun Bhaskara
University of Colorado Boulder
Boulder, Colorado, USA
Varun.Bhaskara@colorado.edu

ACM Reference Format:

Akshay Abhyankar, Aditi Ramesh Athreya, Gaurav Roy, and Varun Bhaskara. 2023. Telecom Customer Churn Analysis.

1 MOTIVATION

The technique of analyzing and forecasting customer attrition in the telecom sector is known as telecom churn analysis. When users of telecom services, such as internet or mobile phone plans, cancel their subscriptions or migrate to a competitor, this is known as customer churn. Churn analysis seeks to comprehend the causes of customer turnover, identify the customers most likely to leave in the future, and create plans to lower churn rates.

For telecom providers, churn analysis holds immense significance as it contributes to safeguarding revenue, cutting costs, enhancing the overall customer experience, fine-tuning customer segmentation, enabling targeted marketing efforts, and facilitating product development. In the eyes of telecom companies, retaining existing customers is always preferable due to the fact that acquiring new customers typically incurs higher costs.

The telecom operational system is designed to cater to a specific average number of customers. Falling below this calculated threshold is considered detrimental to the company. Even a modest effort to retain an existing customer can result in a substantial increase in both revenues and profits. Hence it is essential to construct robust machine learning models that can identify the causes of churn and propose necessary enhancements for customer retention.

2 LITERATURE SURVEY

This presents a short summary of churn prediction in telecom industry as well as related work proposed by renowned researchers

Amal et al [1], does an extensive survey on the customer churn prediction by using 5 different datasets and by applying all the pre-processing, data mining techniques that have earlier been employed and evaluating those bases on parameters which have earlier been used for evaluation. The survey explores other research papers and gathers the performance of several different models (more than 8) and compares them. This paper shows all the comparison results in a tabular format so that it gets easier to compare based on various features.

Preeti et al [2], forecast's customer churn through the utilization of statistical survival analysis methods like logistic regression and decision trees. It marks the significance of meticulously choosing attributes and features to enhance the precision of churn prediction. This system holds considerable potential for telecommunications companies as it empowers them to pinpoint and retain customers at risk of churning, consequently safeguarding revenues typically allocated to customer acquisition and retention efforts. Rather than giving too much importance to models, it majorly gives more importance to feature selection as that can result to different results in case of different models.

Prabhadevi et al [3], utilized the Kaggle dataset, which comprises 7044 instances featuring 21 distinct attributes. The study began with data collection and preprocessing, which involved tasks such as one-hot encoding and label encoding to convert categorical labels into numerical format. Multiple classification algorithms were employed to categorize customers into churn and non-churn categories. These algorithms included Stochastic Gradient Booster, Random Forest, KNN, and Logistic Regression. The performance of these algorithms was assessed based on various metrics, including Accuracy, Precision, Recall, F1-score, support, ROC (Receiver Operating Characteristic), and AUC (Area Under the Curve). Notably, the Stochastic Gradient Booster algorithm achieved the highest accuracy among the evaluated models.

Ullah et al [4], uses two datasets. The first dataset is obtained from South Asia GSM telecom service provider for studying customer churn prediction problem. It has 64,107 instances with 29 features. The second dataset is a publicly available churn-bigml dataset <http://bigml.com/user/francisco/gallery/dataset/5163ad540c0b5e5b22000383>. The dataset contains 3333 instances and 16 features. The data pre-processing consists of noise removal and feature selection using Information Gain and Correlation Attributes Ranking Filter techniques. Various classification algorithms are used to classify the customers as churn and non-churn customers. Algorithms such as Random Forest, Decision Stump, J48, Random Tree, AdaboostM1 + Decision Stump and Bagging + Random Tree, Naive Bayes, Multilayer Perceptron (MLP), Logistic Regression (LR), IBK and LWL. These algorithms are evaluated using Accuracy, FP Rate, TP Rate, Precision, Recall and F-measure. The Random Forest algorithm has the highest accuracy. After classification, the churn customers are categorized into low, medium and risky customers using the

K-Means clustering algorithm.

Verbeke et al [5], explores the data mining techniques and machine learning models to predict the customer churn using real world telecom data as in many countries retention of telecommunication customers seems to be more profitable than gathering new customers. The preprocessing techniques mentioned in this paper include variable selection, which is selecting the best feature that would act as an event for the churn prediction, followed by oversampling as the datasets might have imbalanced data and then with data cleaning by removing null and unnecessary values. The paper then applies rule-based classifiers, decision trees, neural networks, ensemble methods, statistical methods, and support vector machines to predict the churns. The models are evaluated based on metrics like the area under the receiver operating curve (AUC) and top decile lift. The paper also introduces a "profit-centric" approach that evaluates models based on their ability to maximize profits by targeting the right customers for retention campaigns.

3 PROPOSED WORK

The dataset at hand, sourced from IBM Sample Data Sets, encapsulates a rich tapestry of customer information, offering insights into the factors that influence customer churn and retention. Each row within the dataset represents a unique customer, while the columns unfold a myriad of customer attributes.

In this research, we aim to predict customer churn in the telecom sector and develop targeted customer retention programs. The proposed work encompasses several key components. First, we will collect and preprocess a comprehensive dataset containing customer attributes, including churn status, subscribed services, account information, and demographic details. Data preprocessing will involve data cleaning, handling missing values, and transforming the data to ensure its quality and suitability for analysis.

Next, we will focus on feature selection and engineering, where we will identify relevant attributes, create new features, and aggregate data over different time intervals. Exploratory data analysis (EDA) will play a crucial role in uncovering insights from the dataset, using visualizations, statistical analysis, and data segmentation techniques.

For churn prediction, we will employ various machine learning algorithms, including logistic regression, AdaBoost, Support Vector Machine and neural networks. The performance of these models will be evaluated using metrics such as accuracy, precision, recall, F1-score, and ROC AUC. Given the potential class imbalance in the churn dataset, we will implement techniques like oversampling to address this issue.

One of the classification methods that we plan on using is **Adaboost** which is short for Adaptive Boosting. AdaBoost is an ensemble machine learning algorithm primarily designed for binary classification tasks. Its key concept centers on the dynamic boosting of misclassified data points. It works as follows

- (1) **Initialization:** Begin by assigning equal weights to all data points in the training dataset.
- (2) **Iterative Training:** During each iteration (t), train a weak classifier using the weighted training data. A weak classifier is a simple model that performs slightly better than random chance. The classifier prioritizes the misclassified data points from the previous round due to updated weights.
- (3) **Classifier Weight Calculation:** Calculate the weighted error (ϵ) of the weak classifier, which is the sum of weights of misclassified points. Compute the weight of the weak classifier in the ensemble using $\alpha = 0.5 \times \ln(\frac{1-\epsilon}{\epsilon})$. A higher α signifies better performance.
- (4) **Weight Update:** Update the weights of the training data points based on their classification correctness by the current weak classifier. Increase the weights of misclassified points and decrease the weights of correctly classified points to focus on the misclassified points for the next iteration.
- (5) **Ensemble Creation:** Combine the weak classifiers into an ensemble, assigning a weight (α) to each weak classifier. The ensemble's prediction is a weighted sum of individual weak classifier predictions.
- (6) **Normalization of Weights:** Normalize the weights of the data points to ensure they sum up to 1, maintaining a probability distribution.
- (7) **Final Classification:** Compute the final ensemble prediction by aggregating the predictions of all weak classifiers based on their weights. The sign of the sum of weighted predictions determines the final prediction (+1 or -1 in binary classification).
- (8) **Repeat Iterations:** Steps 2 to 7 are repeated for a predetermined number of iterations or until a desired level of accuracy is reached.
- (9) **Final Ensemble:** The final ensemble is a combination of multiple weak classifiers, with each classifier contributing based on its performance (α).
- (10) **Prediction:** Utilize the final ensemble to make predictions on new, unseen data..

Similarly we plan on using **Support Vector Machine** to improve the accuracy further. The objective of the support vector machine algorithm is to find a hyperplane in an N-dimensional space(N — the number of features) that distinctly classifies the data points. In the process of distinguishing between the two data point classes, various hyperplanes can be considered. Our goal is to identify a hyperplane with the largest margin, which is the greatest separation between data points from both classes.

- Hyperplanes serve as decision boundaries that aid in classifying data points. They separate data points into distinct classes based on which side of the hyperplane they fall. The dimension of the hyperplane depends on the number of features in the data.
- Support vectors are data points that are closer to the hyperplane and influence the position and orientation of the hyperplane. Using these support vectors, we maximize the margin of the classifier. Deleting the support vectors will change the position of the hyperplane.

- In SVM we aim to maximize the margin between the data points and the hyperplane. The loss function that helps maximize the margin is hing loss given by:

$$c(x, y, f(x)) = \begin{cases} 0, & \text{if } y * f(x) \geq 1 \\ 1 - y * f(x), & \text{else} \end{cases} \quad c(x, y, f(x)) = (1 - y * f(x))_+$$

- The cost is zero when the predicted value aligns with the actual value in terms of sign. In cases where they do not match, the loss value is computed. Additionally, we introduce a regularization parameter into the cost function to strike a balance between maximizing the margin and minimizing the loss.

$$\min_w \lambda ||w||^2 + \sum_{i=1}^n (1 - y_i \langle x_i, w \rangle)_+$$

- Using this loss function we take partial derivatives with respect to the weight to find the gradients. Using the gradients, we can update our weights/

$$\frac{\delta}{\delta w_k} ||w||^2 = 2\lambda w_k$$

$$\frac{\delta}{\delta w_k} (1 - y_i \langle x_i, w \rangle)_+ = \begin{cases} 0, & \text{if } y_i \langle x_i, w \rangle \geq 1 \\ -y_i x_{ik}, & \text{else} \end{cases}$$

If there is no misclassification, we only have to update the gradient from the regularization parameter.

$$w = w - \alpha \cdot (2\lambda w)$$

If there is misclassification, we include the loss along with the regularization parameter to perform the gradient update

$$w = w + \alpha \cdot (y_i \cdot x_i - 2\lambda w)$$

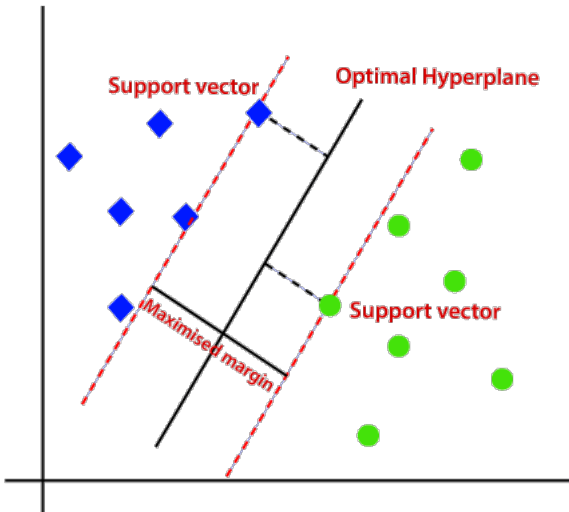


Figure 1: Support Vector Machine

Another classification method that holds significance is **Logistic Regression**, a widely used statistical method for binary classification and, with extensions, for multi-class classification as well. It

models the probability that a given instance belongs to a particular category. The key principles of logistic regression are as follows:

- (1) **Model Representation:** Logistic Regression models the probability $P(Y = 1)$ that an instance belongs to the positive class. The logistic function, denoted as $\sigma(z)$, where z is a linear combination of input features, is employed:

$$P(Y = 1) = \sigma(z) = \frac{1}{(1 + e^{-z})}$$

- (2) **Cost Function:** The logistic regression model aims to minimize a cost function, often the cross-entropy loss. For a single training instance, the cost is defined as:

$$J(\theta) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

where \hat{y} is the predicted probability, and y is the true class label (0 or 1).

- (3) **Gradient Descent:** Optimization techniques such as gradient descent are employed to minimize the cost function by adjusting the model parameters (weights θ). The update rule for gradient descent is:

$$\theta := \theta - \alpha \nabla J(\theta)$$

where α is the learning rate, and $\nabla J(\theta)$ is the gradient of the cost function with respect to the parameters.

- (4) **Decision Boundary:** Logistic Regression calculates a decision boundary that separates the instances into two classes. In a two-dimensional space, the decision boundary is a line, and in higher dimensions, it becomes a hyperplane.
- (5) **Probabilistic Interpretation:** The output of the logistic regression model can be interpreted as the probability of the instance belonging to the positive class. A threshold is chosen (commonly 0.5), and instances with predicted probabilities above this threshold are classified as positive.
- (6) **Prediction:** Once trained, the logistic regression model can be used to predict the probability of new instances belonging to the positive class.

Similarly, **Random Forest** is an ensemble learning method that operates by constructing a multitude of decision trees during training and outputs the mode of the classes (classification) or the mean prediction (regression) of the individual trees.

- (1) **Model Representation:** Logistic Regression models the probability $P(Y = 1)$ that an instance belongs to the positive class. The logistic function, denoted as $\sigma(z)$, where z is a linear combination of input features, is employed:

$$P(Y = 1) = \sigma(z) = \frac{1}{(1 + e^{-z})}$$

- (2) **Initialization:** Begin with a dataset containing features and their corresponding labels. Random Forest builds multiple decision trees during the training phase.
- (3) **Bootstrapping :** For each tree, a random subset of the training data (with replacement) is selected. This process is known as bootstrapping, and it creates diverse datasets for each tree.
- (4) **Feature Randomization:** At each node of the decision tree, a random subset of features is considered for splitting. This helps to decorrelate the trees and ensures that no single feature dominates the decision-making process.

- (5) **Tree Growing:** Each decision tree is grown by recursively splitting nodes based on the selected features. The splits are determined by evaluating various criteria, such as Gini impurity for classification or mean squared error for regression.
- (6) **Voting (Classification) or Averaging (Regression):** For classification tasks, the final prediction is determined by a majority vote from all the decision trees. For regression tasks, the final prediction is the average of the predictions from all the trees.

We plan to deploy the selected churn prediction models for both real-time and batch scoring, enabling telecom companies to take timely retention actions and gain periodic insights. Customer segmentation based on churn likelihood will guide the development of tailored retention strategies, which may include personalized offers, service upgrades, and customer support enhancements.

The overarching goal of this research is to minimize churn rates while maximizing customer lifetime value and overall profitability.

4 EVALUATION

We plan on using the following strategies to evaluate our work

- (1) **Confusion Matrix:** We chose confusion matrix as it is one of the most important measure while evaluating prediction analysis, it makes evaluation easy by categorizing the outcomes TP, TN, FP and FN; which in our case will be Customers correctly predicted as churners, Customers correctly predicted as non-churners, Customers incorrectly predicted as churners and Customers incorrectly predicted as non-churners. As it reports false positives and false negatives so it will be helpful for us in devising if there is anything wrong with the strategy in selecting the feature.
- (2) **Accuracy:** As the most fundamental metric for evaluating ML models, accuracy will calculate the ratio of both the churners and the non churners and give us a high level view of the data.

$$\text{Accuracy} = \frac{(TP + TN)}{(TP + FP + FN + TN)}$$

- (3) **Precision, Recall and F1:** We will also evaluate the models predicting the churners by using Precision and Recall as they will help us in knowing the proportion of correctly predicted churners in a detailed way. And after precision and recall, we can apply F1 to check on the prediction of both churners and non churners.

$$\begin{aligned} \text{Precision} &= \frac{TP}{TP + FP} \\ \text{Recall} &= \frac{TP}{TP + FN} \\ \text{F1-score} &= \frac{2 \cdot (\text{Precision} \cdot \text{Recall})}{\text{Precision} + \text{Recall}} \end{aligned}$$

If required, we are also thinking of using AUC-ROC to evaluate the models due to its characteristics of handling discrimination in data and ways of handling imbalanced data. We plan to use

Cross-validation as our validation strategy:

Cross-Validation - Dividing the dataset into a training set and a validation/test set.

The model is trained on the training set, and its effectiveness is evaluated on the validation/test set.

Iteratively, training and validating the model while dividing the data into different subsets, cross-validation (such as k-fold cross-validation), is a strong method that makes sure all data points are taken into account. To handle outliers, include outlier data points in different folds to ensure robust model evaluation. To handle the slight imbalance in classes we will be using resampling techniques like oversampling and undersampling. Oversampling involves increasing the number of instances in the minority class. Undersampling involves decreasing the number of instances in the majority class.

5 MILESTONES

We propose the following timeline for the Telecom Customer Churn Analysis Project. The project's first phase, encompassing 1-2 weeks focuses on comprehensive data pre-processing, a pivotal step ensuring data quality and suitability for analysis. Following this, 1-2 weeks are allocated to exploratory data analysis (EDA), offering insights into underlying patterns and trends within the data set. Subsequently, 3-4 weeks will be invested in the application of diverse machine learning models, with an emphasis on model selection and optimization. A further 2-3 weeks will be devoted to rigorous model evaluation, employing various metrics and techniques to ensure robust and accurate results. The project's conclusion is allocated 1 week for comprehensive reporting, documentation, and presentation of findings.

6 RESULTS

After doing a thorough exploration of the data, the below key insights were found about the customer behavior.

- **Overall Churn Rate:** Approximately 26.6% of customers decided to switch to another telecom service provider, indicating a significant portion of the customer base in transition. (Figure2)

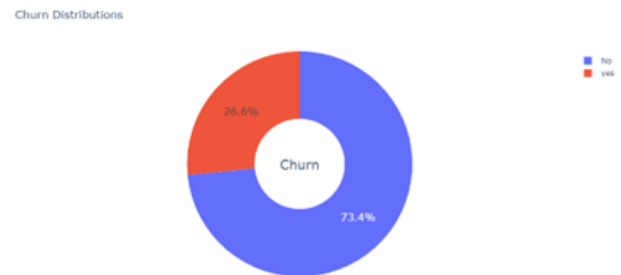


Figure 2: Overall Churn Rate

- **Churn Distribution by Gender:** When examining churn distribution based on gender, we observed a negligible difference. Both male and female customers exhibited similar tendencies when it came to migrating to a different service provider.(Figure3)

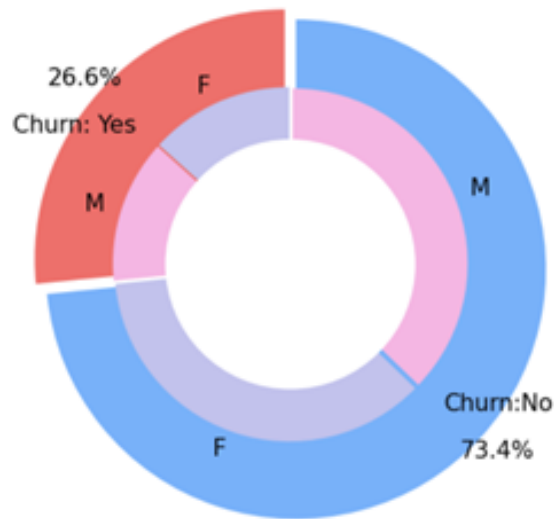


Figure 3: Churn Distribution by Gender

- **Customer Contract Preferences:** The analysis revealed that customers with a Month-to-Month contract were more likely to opt for a change in service provider compared to those with longer-term contracts. Specifically, around 42% of customers with Month-to-Month contracts moved out, in contrast to 13% with One-Year contracts and 3% with Two-Year contracts.
- **Impact of Tech Support:** A noteworthy insight emerged regarding the influence of tech support. Customers without tech support were found to be the most likely to migrate to another service providert.
- **Payment Method Preferences:** Payment method choices also played a role in customer retention. Customers using Electronic Check as their payment method were more prone to churn. In contrast, those utilizing Credit-Card automatic transfer, Bank Automatic Transfer, or Mailed Check were less likely to switch to a different service provider.
- **Distribution of monthly charges by turn:** It is seen from the data that customers with higher monthly charges are more likely to churn. (Figure4)

After analyzing the customer behavior, we did k fold cross validations on our models such as Logistic Regression, SVM, Random Forest and AdoBoost and then tested them for their prediction evaluation.

From all the above evaluation metrics it is seen that AdaBoost has an accuracy of 81%, and it correctly identified 302 correct customers who would churn and 1412 correct customers who would

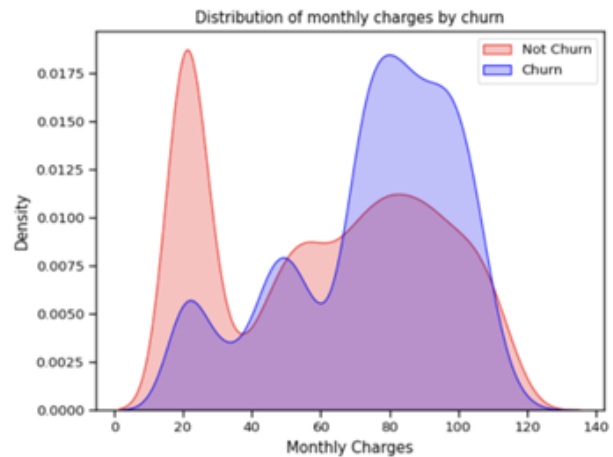


Figure 4: Distribution of monthly charges by turn

Model	Accuracy	Precision	Recall	F1
Logistic Regression	0.8057	0.6583	0.5597	0.6050
SVM	0.7838	0.6599	0.4635	0.5445
Random Forest	0.7967	0.6507	0.5080	0.5445
Adaboost	0.8123	0.6879	0.5383	0.6040

Table 1: Results

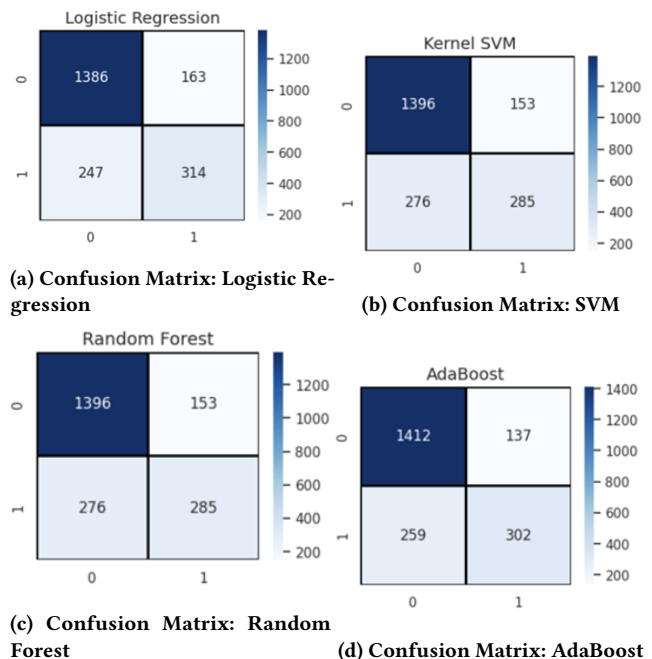


Figure 5: Confusion Matrices

not churn. But it also incorrectly classified 137 instances of customers incorrectly who would churn and it missed 259 instances of

customers who had been churned. So the precision of this model is 68% and the recall is 53%.

In order to further optimize the results of the AdaBoost model in future, we could explore the concepts of hyperparameter tuning. Fine-tuning the model's hyperparameters can lead to improvements in accuracy, precision, and recall, ultimately enhancing its predictive capabilities.

REFERENCES

- [1] Amal M. Almana, Mehmet Sabih Aksoy, and Rasheed AlZahrani. 2014. A Survey On Data Mining Techniques In Customer Churn Analysis For Telecom Industry. <https://api.semanticscholar.org/CorpusID:1283811>
- [2] Preeti K. Dalvi, Siddhi K. Khandge, Ashish Deomore, Aditya Bankar, and V. A. Kanade. 2016. Analysis of customer churn prediction in telecom industry using decision trees and logistic regression. <https://doi.org/10.1109/CDAN.2016.7570883>
- [3] B. Prabadevi, R. Shalini, and B.R. Kavitha. 2023. Customer churning analysis using machine learning algorithms. *International Journal of Intelligent Networks* 4 (2023), 145–154. <https://doi.org/10.1016/j.ijin.2023.05.005>
- [4] Irfan Ullah, Basit Raza, Ahmad Kamran Malik, Muhammad Imran, Saif Ul Islam, and Sung Won Kim. 2019. A Churn Prediction Model Using Random Forest: Analysis of Machine Learning Techniques for Churn Prediction and Factor Identification in Telecom Sector. *IEEE Access* 7 (2019), 60134–60149. <https://doi.org/10.1109/ACCESS.2019.2914999>
- [5] Wouter Verbeke, Karel Dejaeger, David Martens, Joon Hur, and Bart Baesens. 2012. New insights into churn prediction in the telecommunication sector: A profit driven data mining approach. *European Journal of Operational Research* 218 (2012), 211–229. <https://doi.org/10.1016/j.ejor.2011.09.031>