

# *СОРЕВНОВАНИЕ* *DATA* *FUSION*

*| 15 ФЕВРАЛЯ – 5 АПРЕЛЯ 2024 |*

**Meetup #1**

29.02.2024

задача  
**Модели оттока**

1 000 000 руб



**Задача 2**  
Предсказание оттока

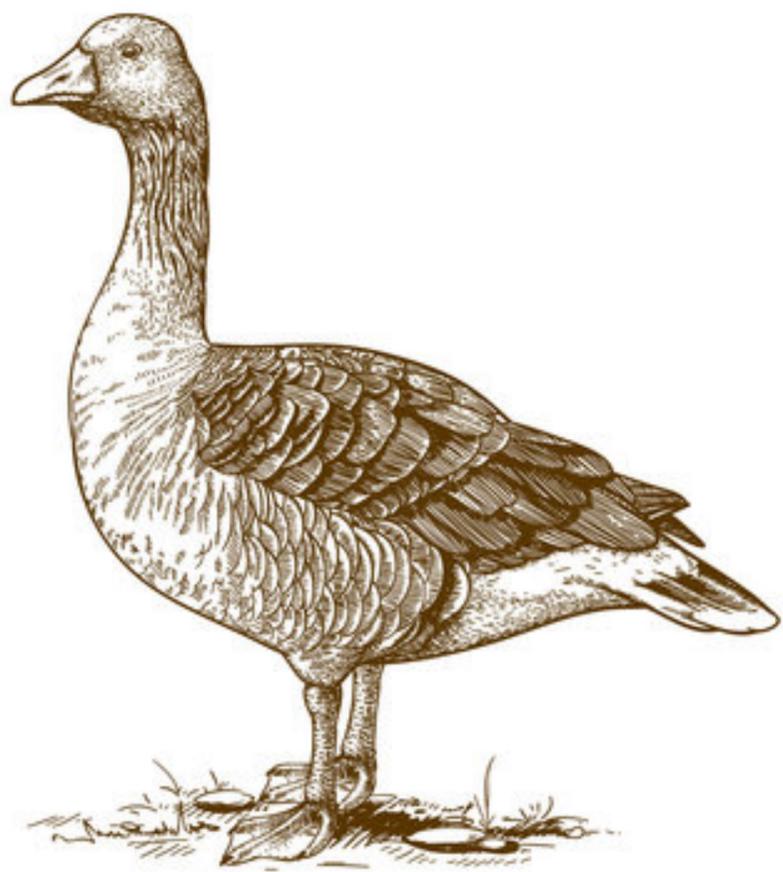
задача  
**Модели оттока**

1 000 000 руб

# Часть 0

0 таблетках и гусях





задача  
**Модели оттока**

1 000 000 руб

# Часть 1

Что происходит в табличных данных?

# Данные

Для решения задачи “Отток” предлагается несколько групп данных и материалов:

## 1. Табличные клиентские данные:

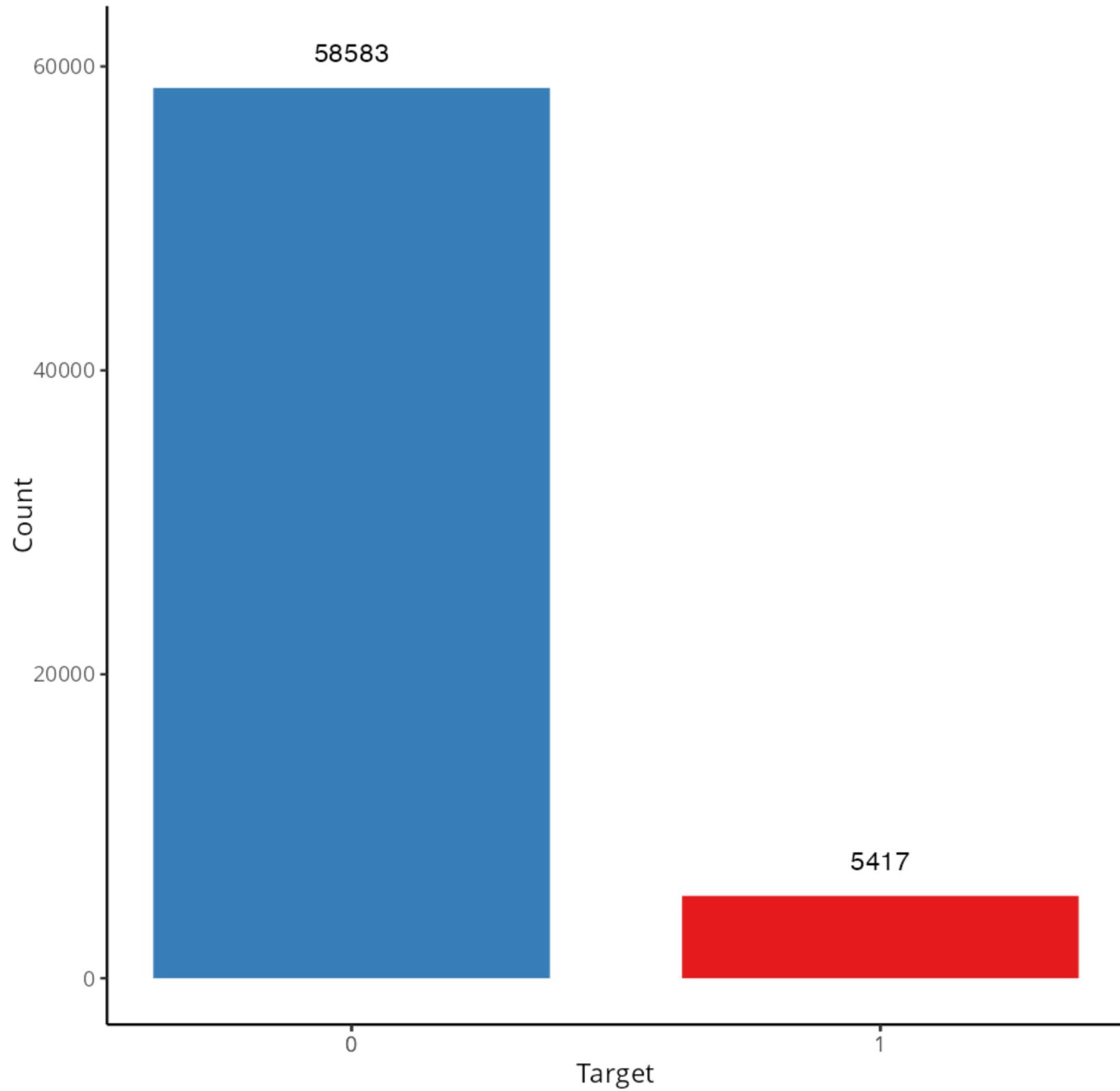
Участникам доступны несколько наборов данных и артефактов:

1. Базовая информация про все 96,000 клиентов в табличном `.csv` формате: `clients.csv` (2.5 MB)
2. Тренировочные данные по целевой переменной и времени последней транзакции для 64,000 клиентов: `train.csv` (745 KB)
3. Информация об отчетах, в рамках которых клиенты сгруппированы по времени: `report_dates.csv` (288 KB)

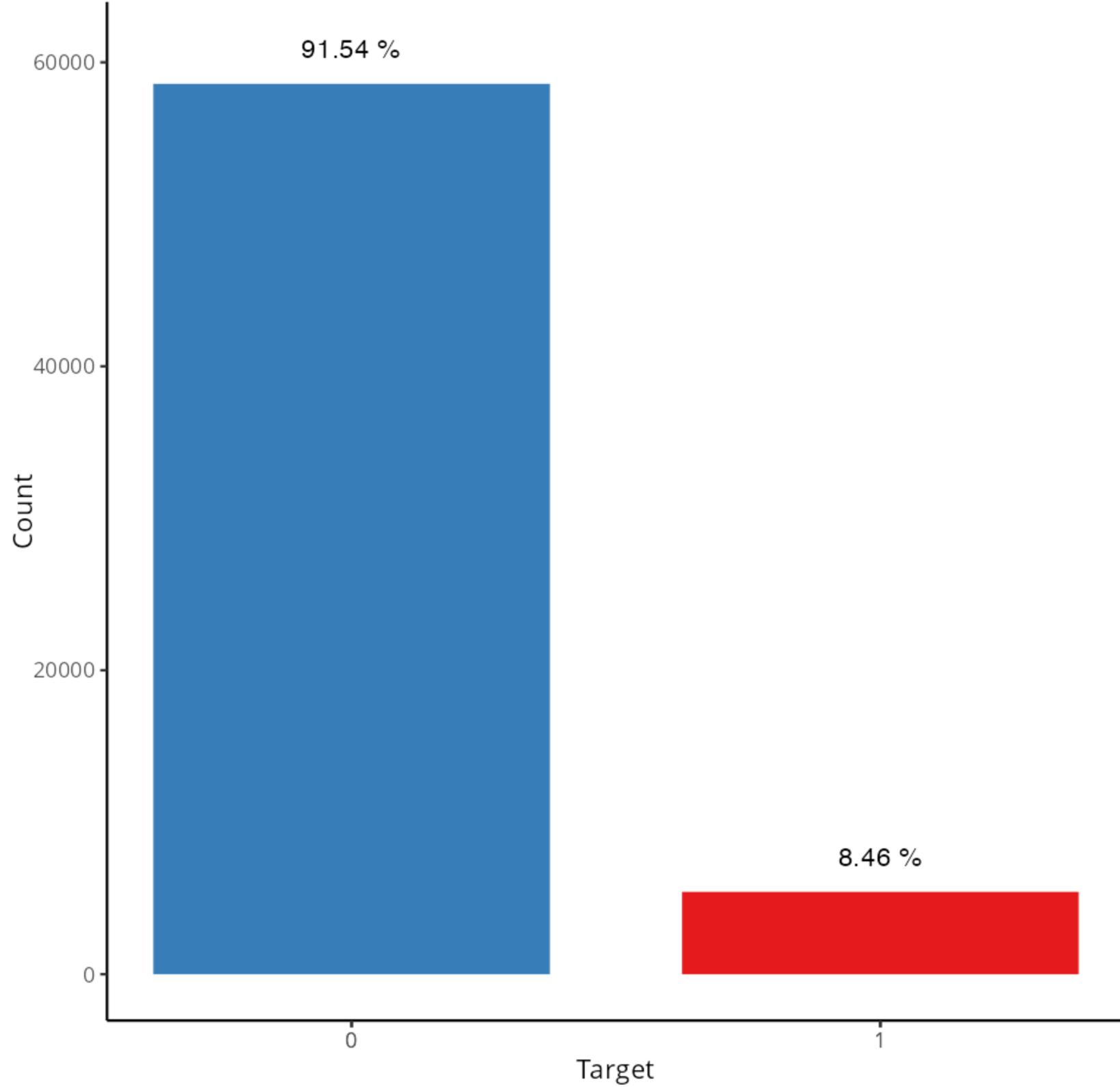
## 2. Данные клиентских транзакций:

1. Клиентские транзакции для всех 96,000 клиентов (~13M) в табличном `.csv` формате: `transactions.csv.zip` (197 MB)

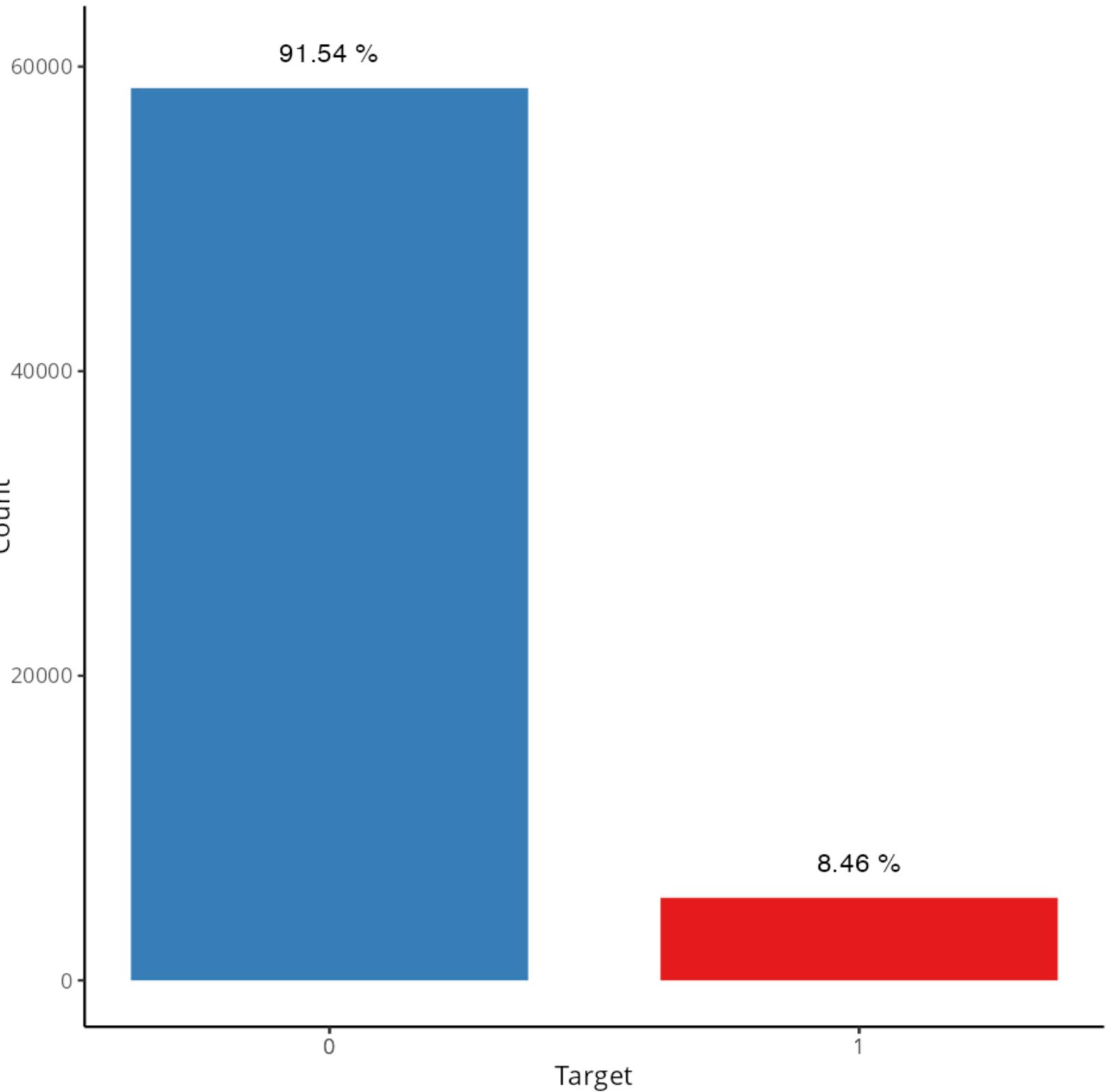
Target distribution in train.csv



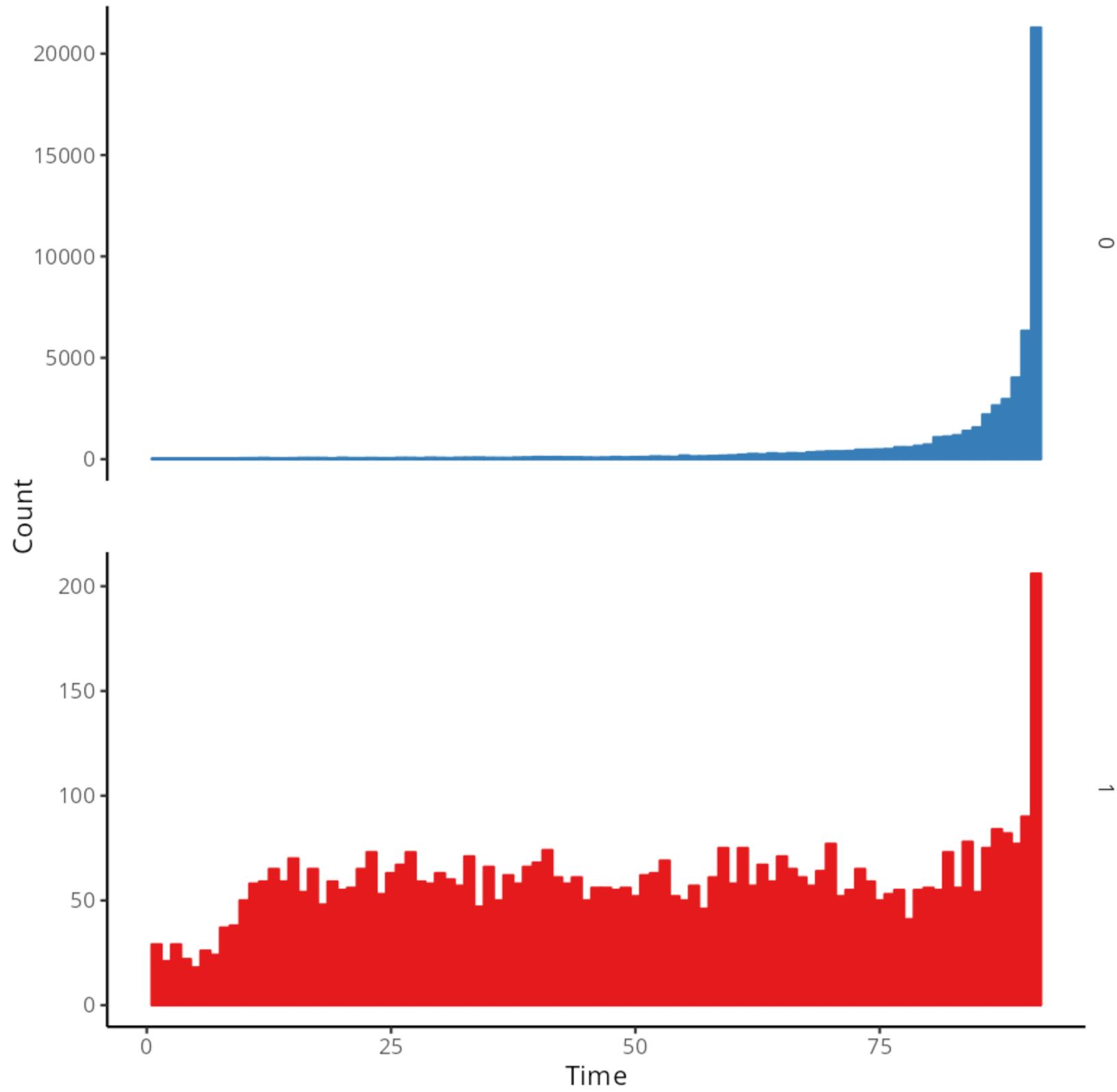
Target distribution in train.csv



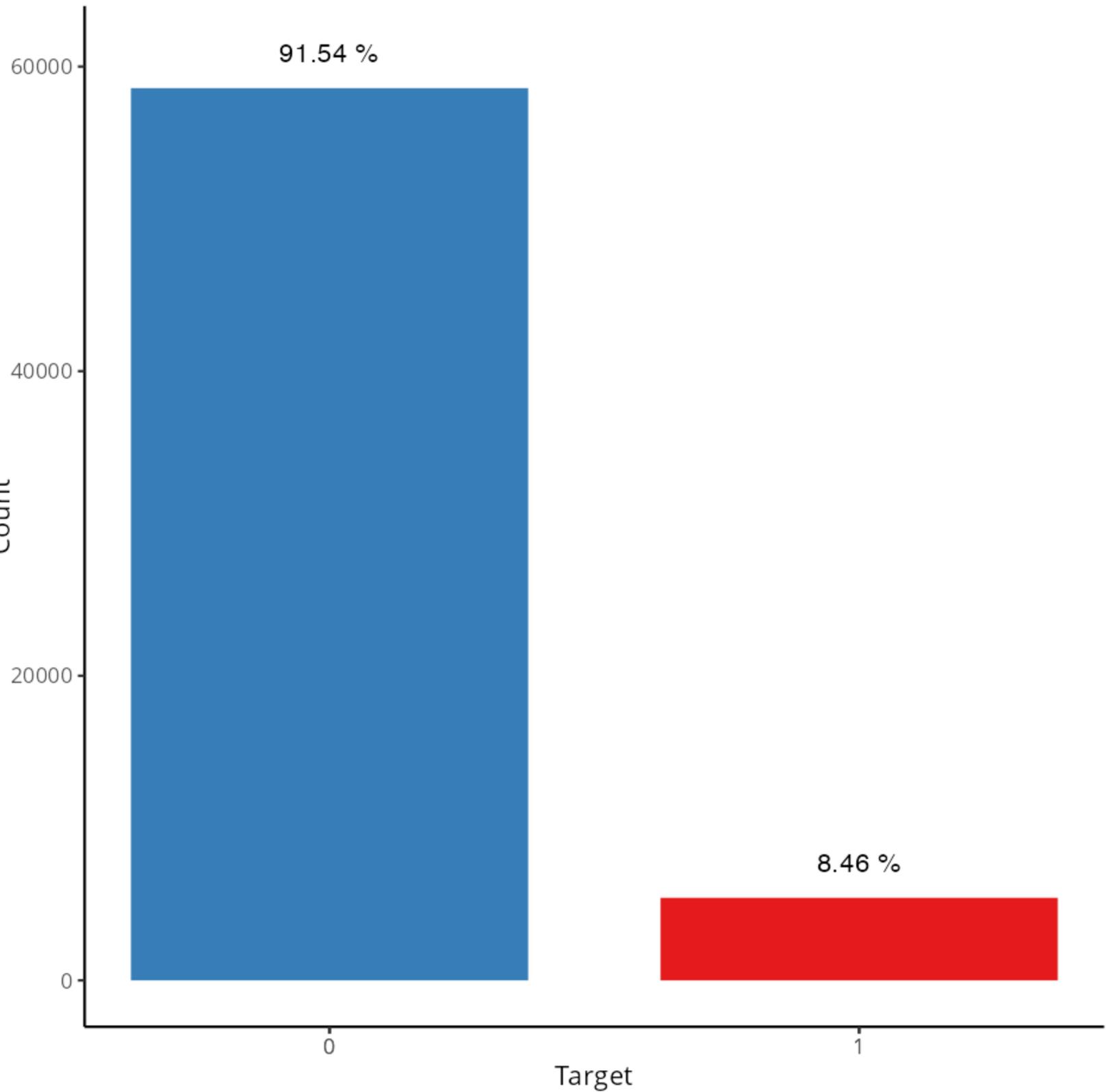
Target distribution in train.csv



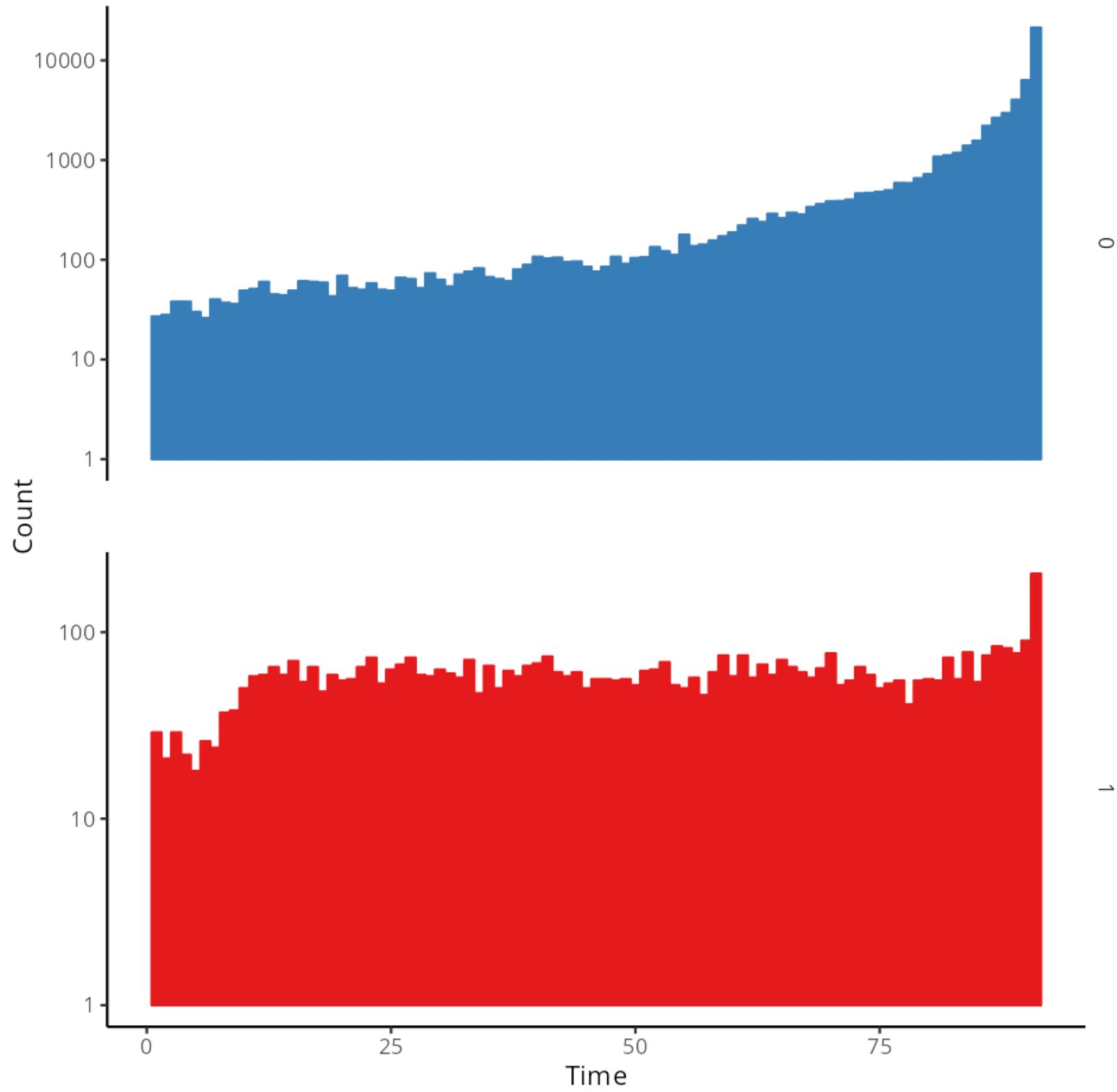
Time \* target distribution in train.csv?



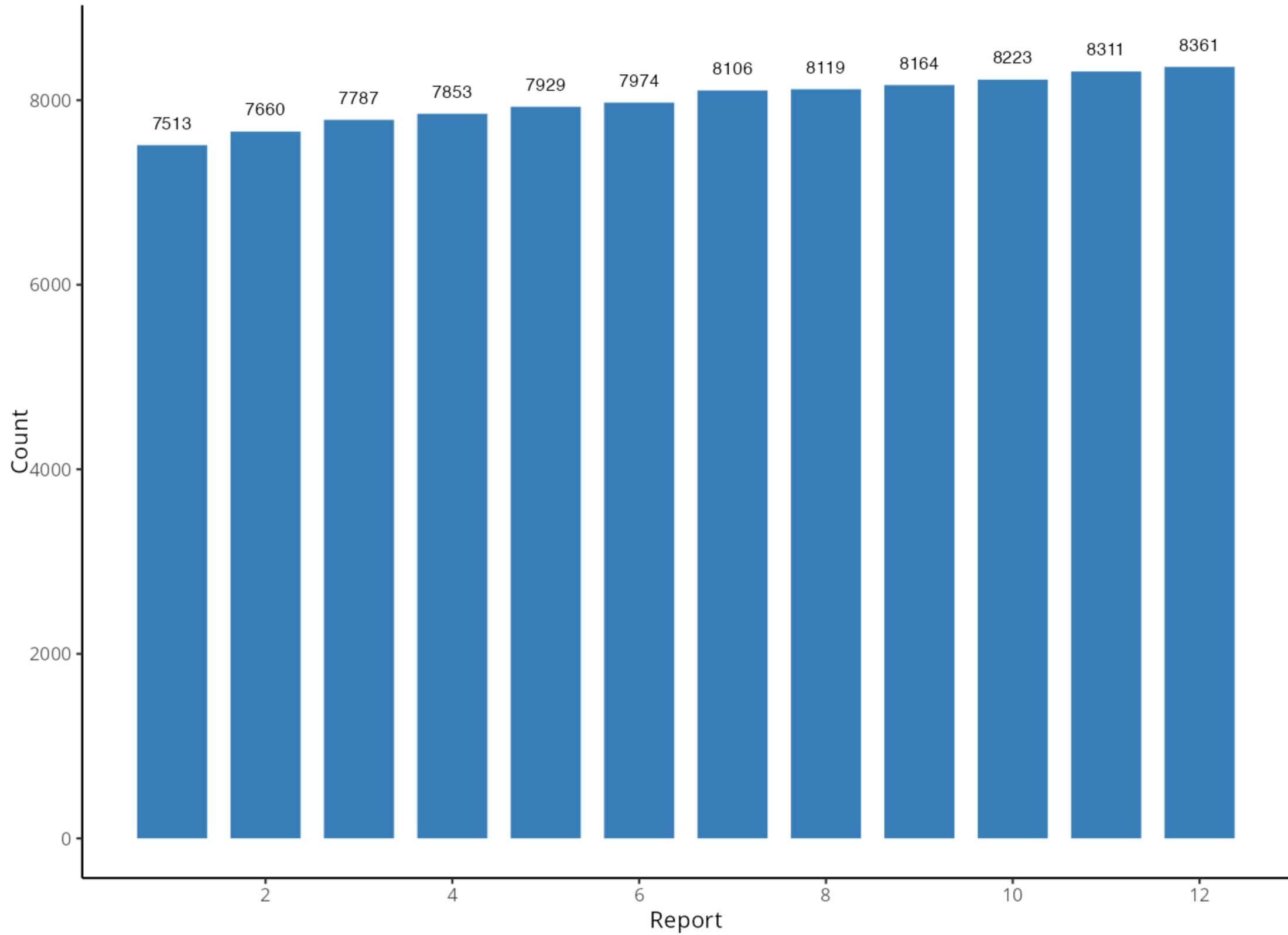
Target distribution in train.csv



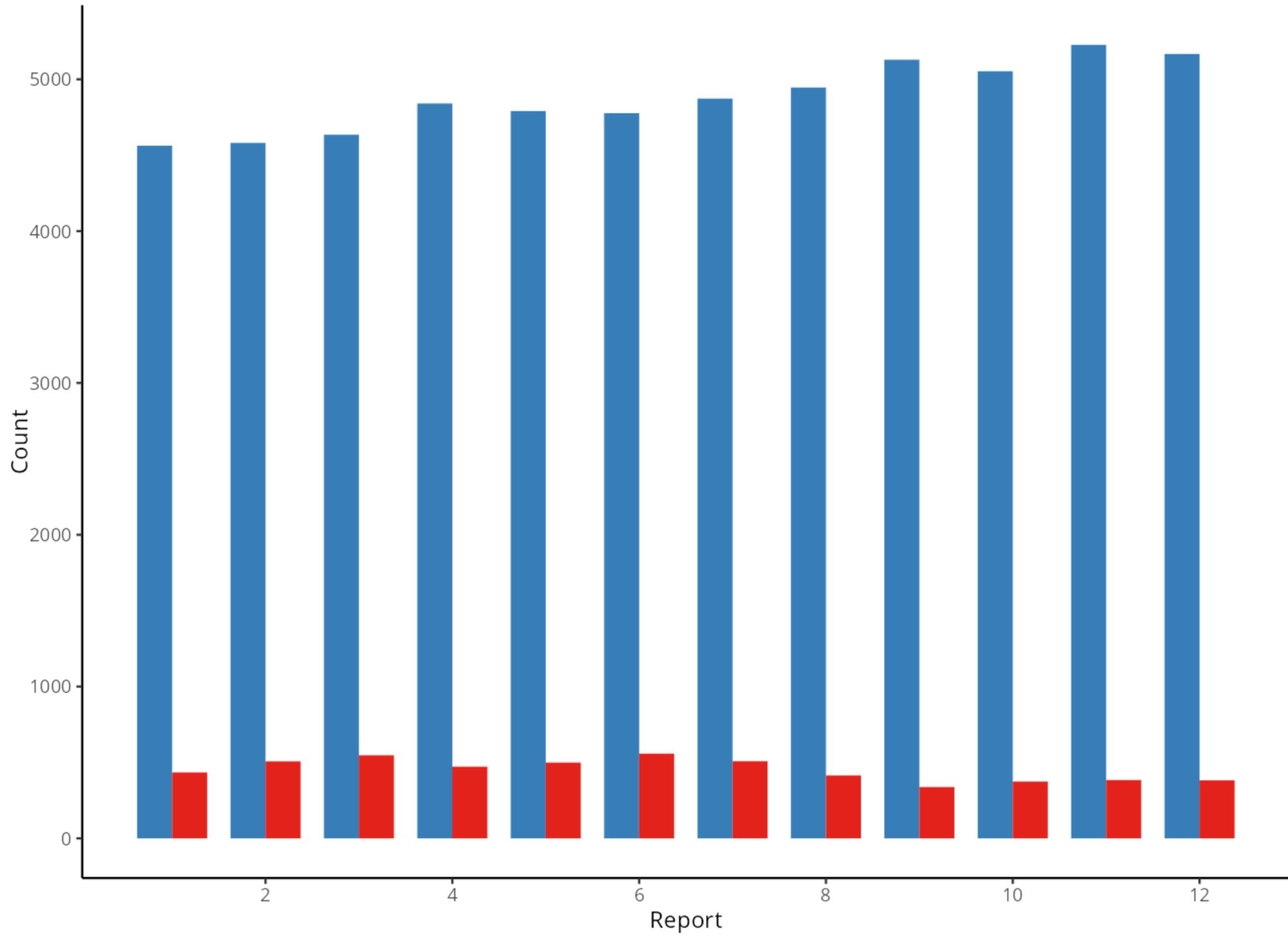
Time \* target distribution in train.csv?



Users by report distribution in clients.csv



Users and target across reports: train.csv \* clients.csv



задача  
**Модели оттока**

1 000 000 руб

## Часть 2

Что происходит в транзакционных данных?

# Данные

Для решения задачи “Отток” предлагается несколько групп данных и материалов:

## 1. Табличные клиентские данные:

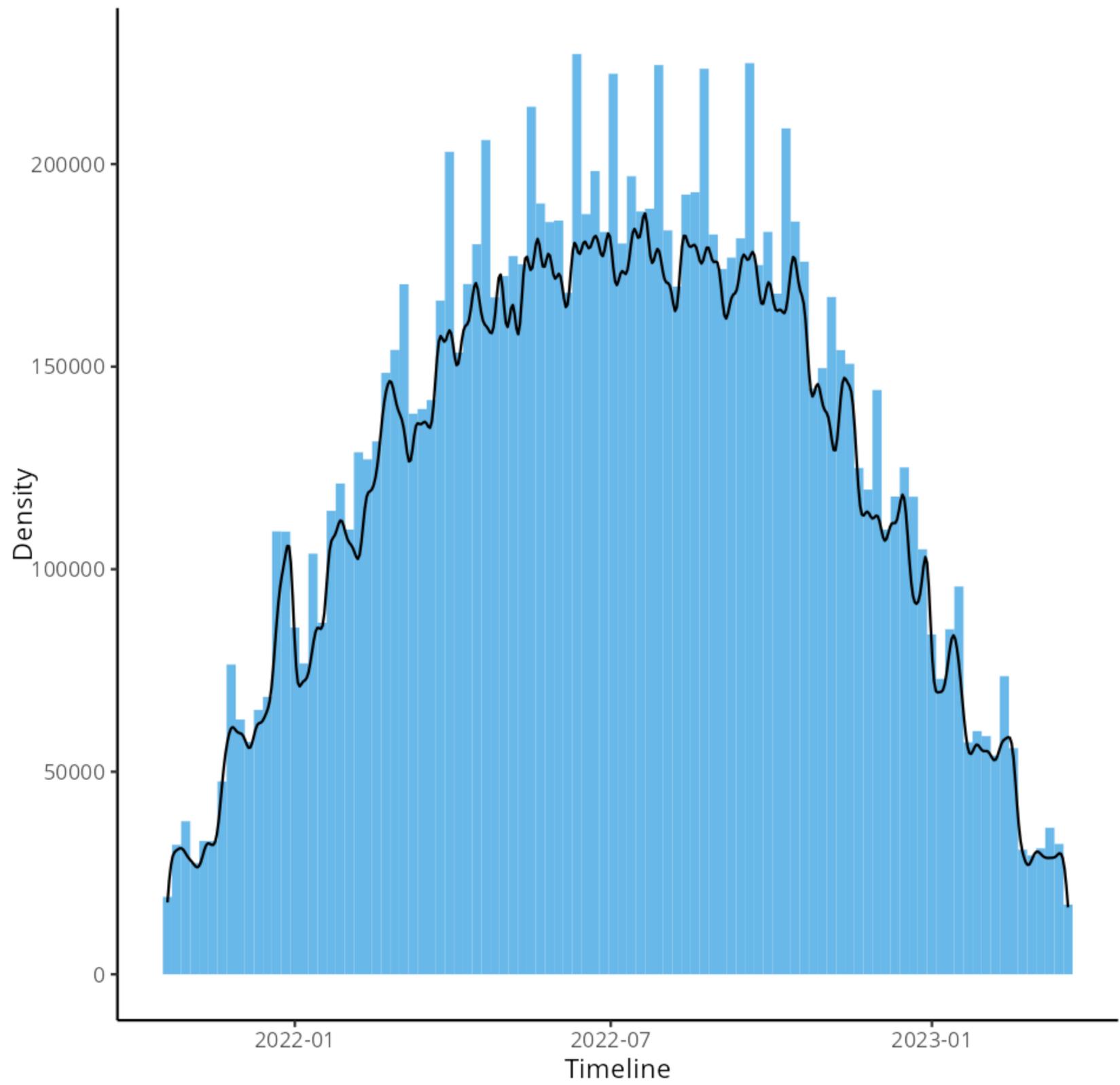
Участникам доступны несколько наборов данных и артефактов:

1. Базовая информация про все 96,000 клиентов в табличном `.csv` формате: `clients.csv` (2.5 MB)
2. Тренировочные данные по целевой переменной и времени последней транзакции для 64,000 клиентов: `train.csv` (745 KB)
3. Информация об отчетах, в рамках которых клиенты сгруппированы по времени: `report_dates.csv` (288 KB)

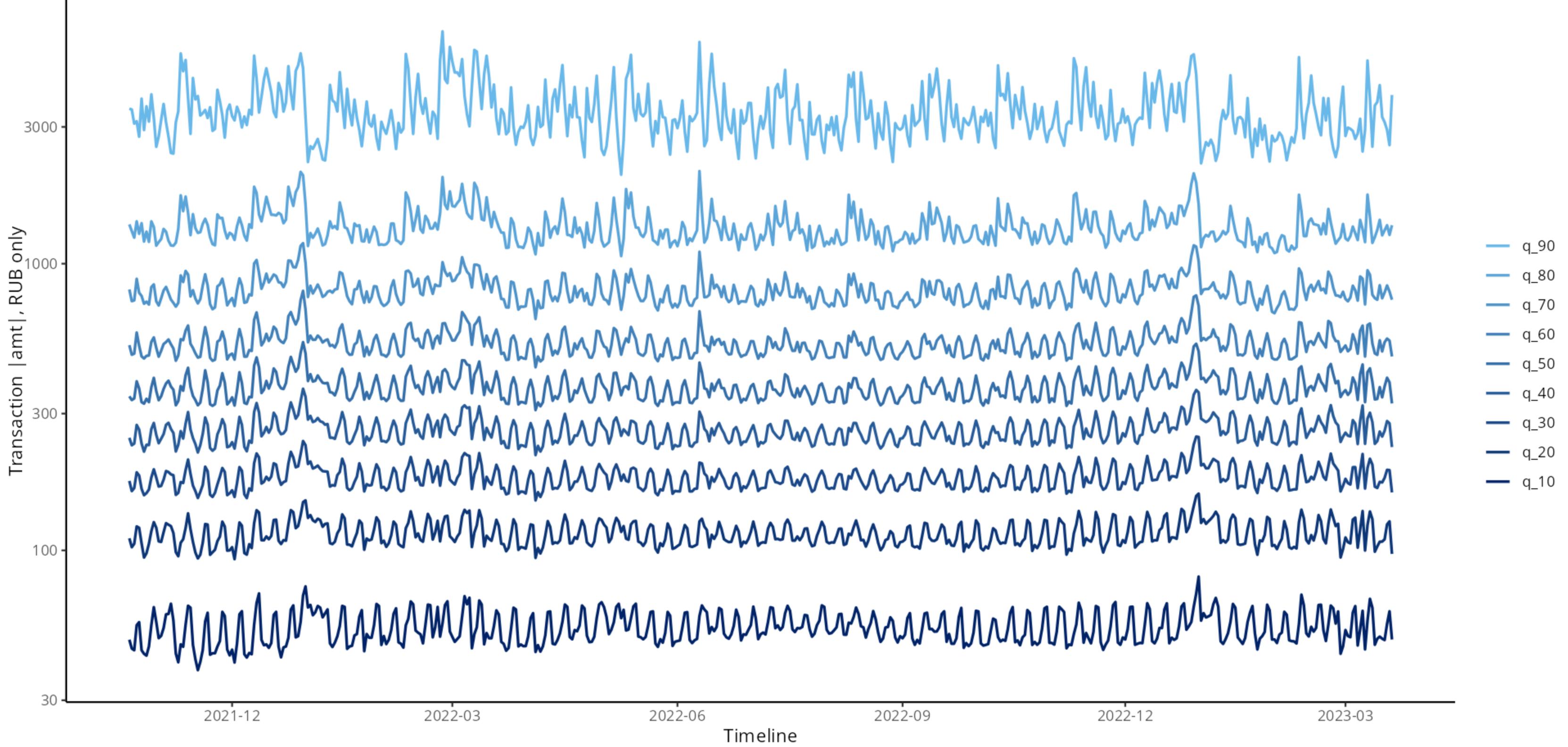
## 2. Данные клиентских транзакций:

1. Клиентские транзакции для всех 96,000 клиентов (~13M) в табличном `.csv` формате: `transactions.csv.zip` (197 MB)

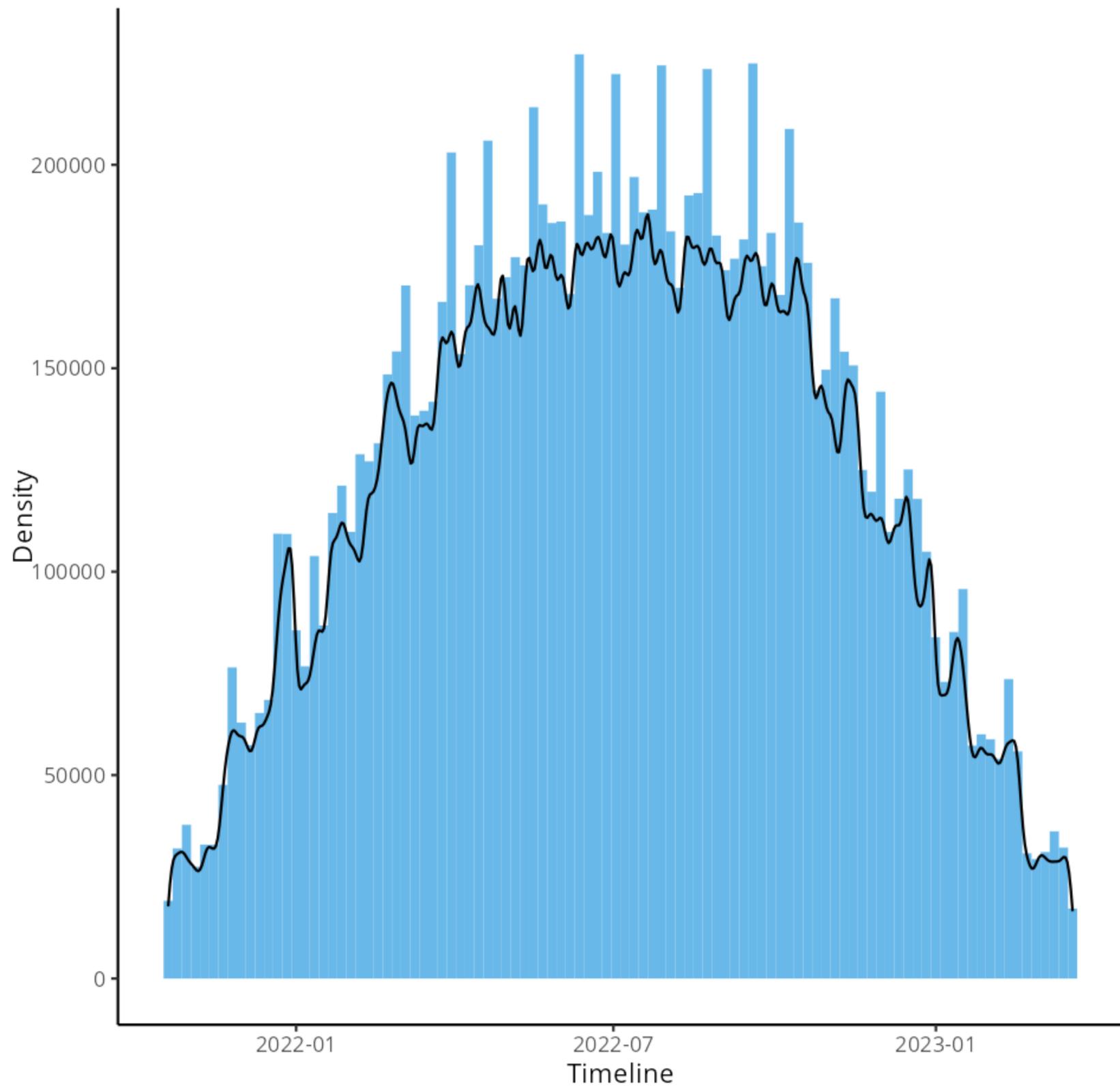
Transaction distribution across time in transaction.csv



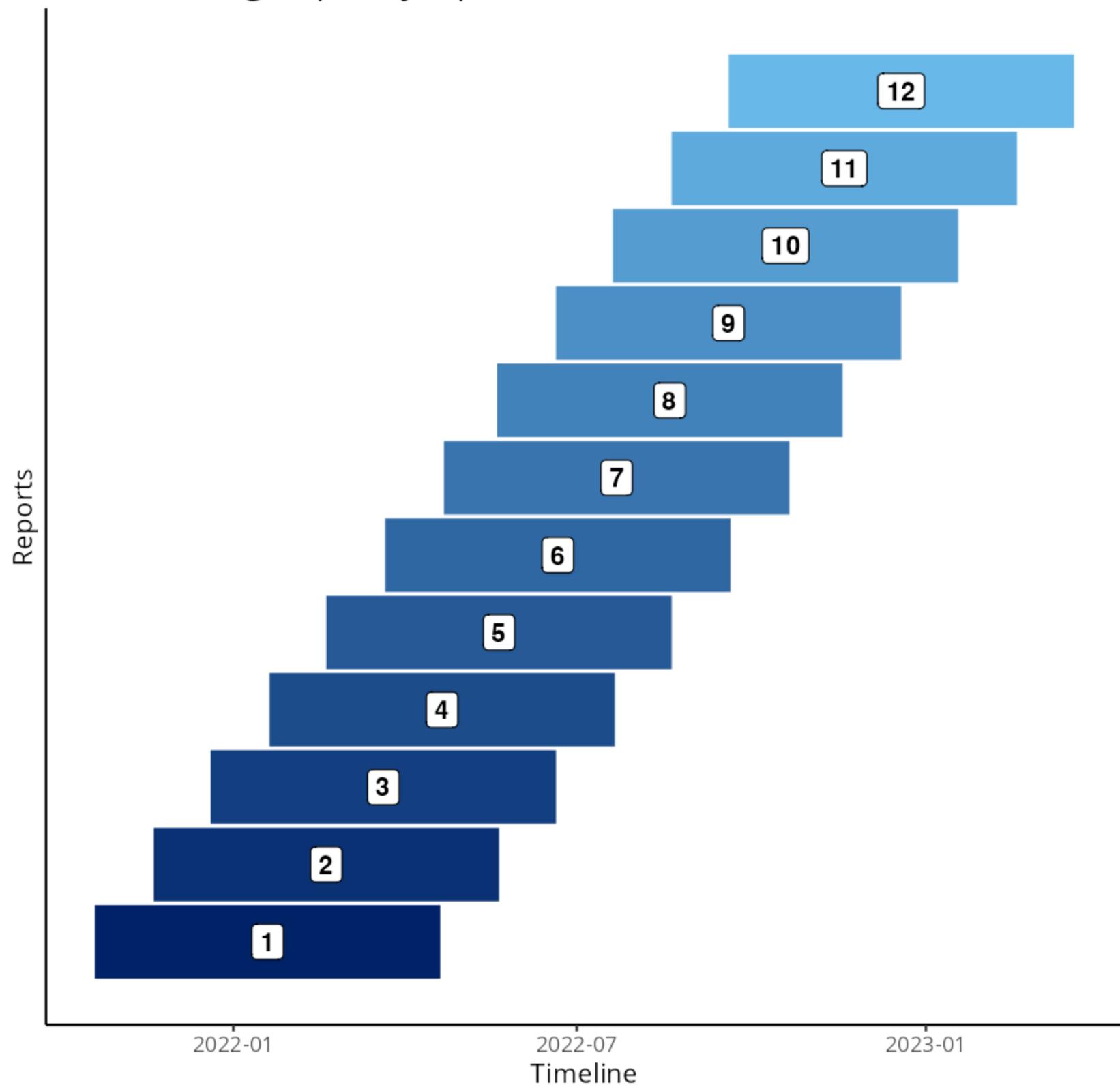
Transaction |amt| quantiles in transaction.csv



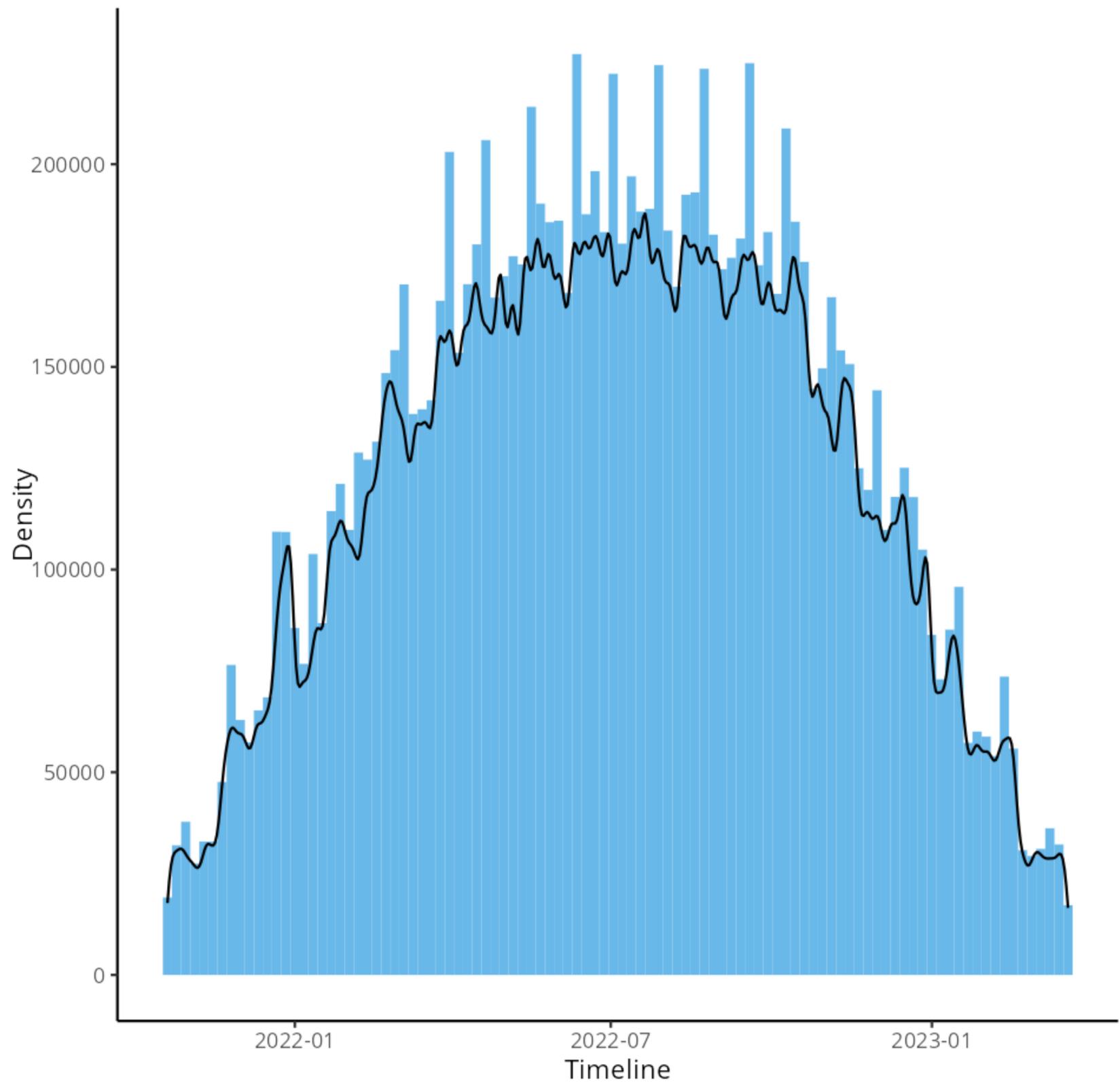
Transaction distribution across time in transaction.csv



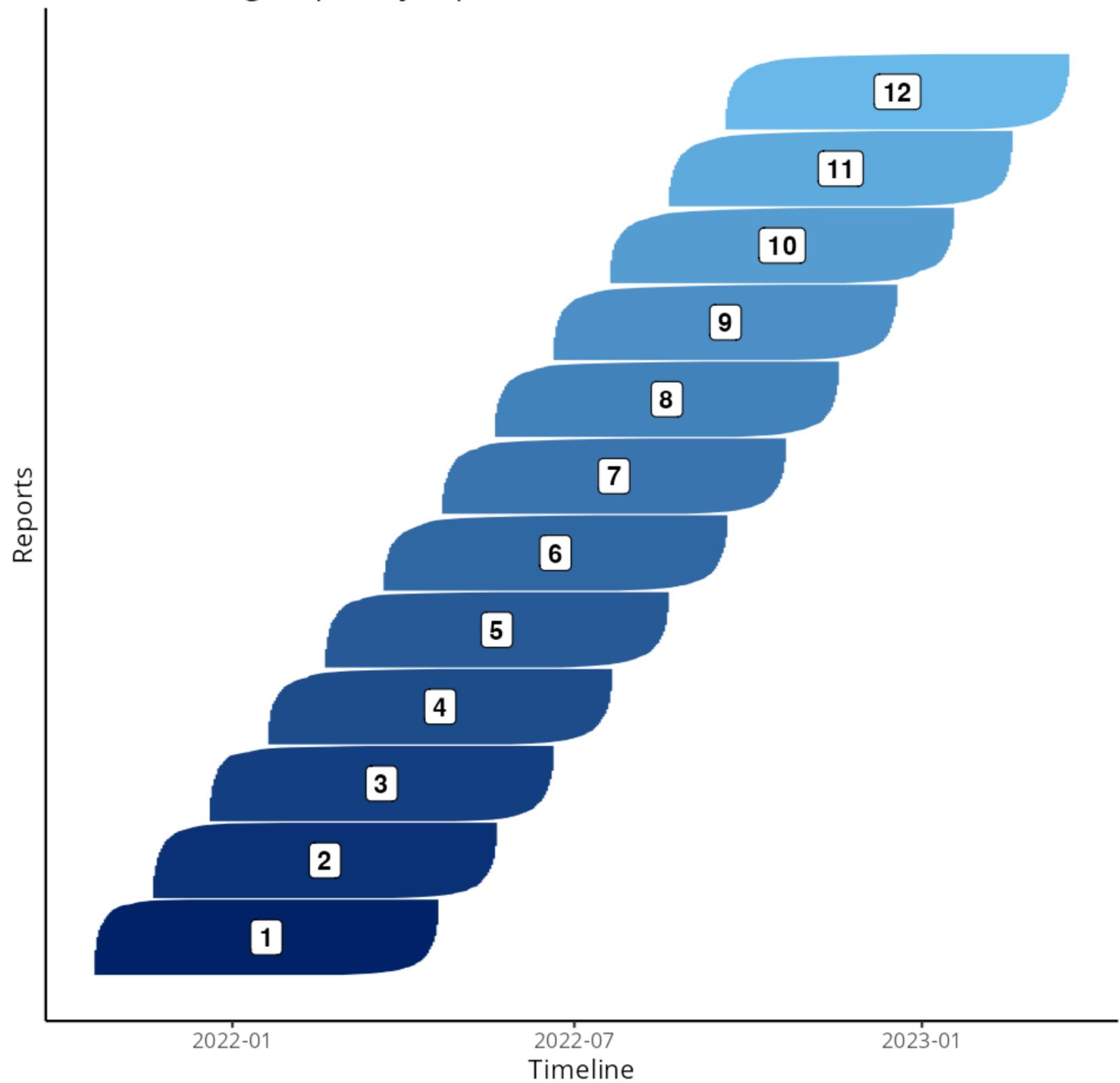
Transactions grouped by reports



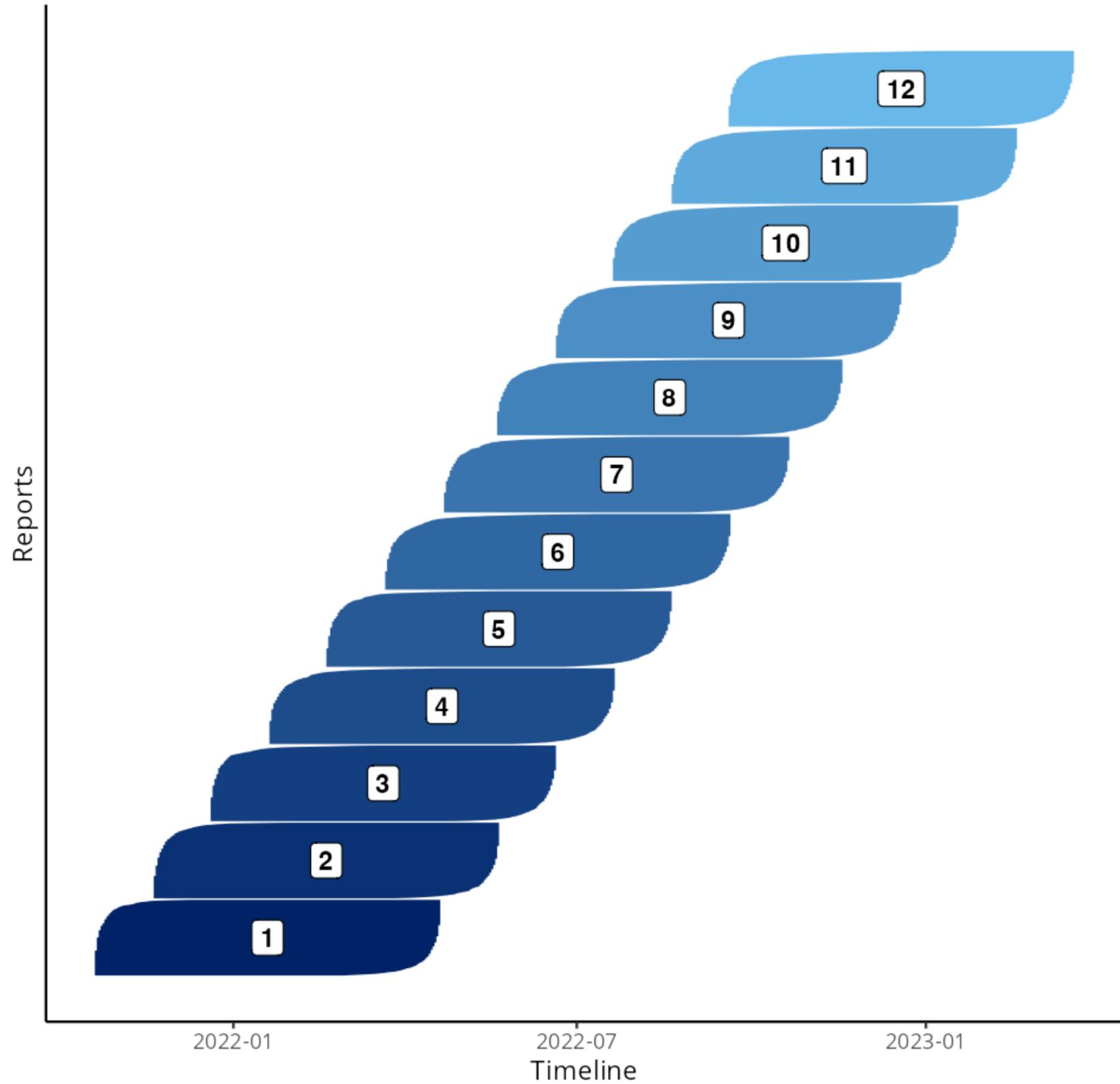
Transaction distribution across time in transaction.csv



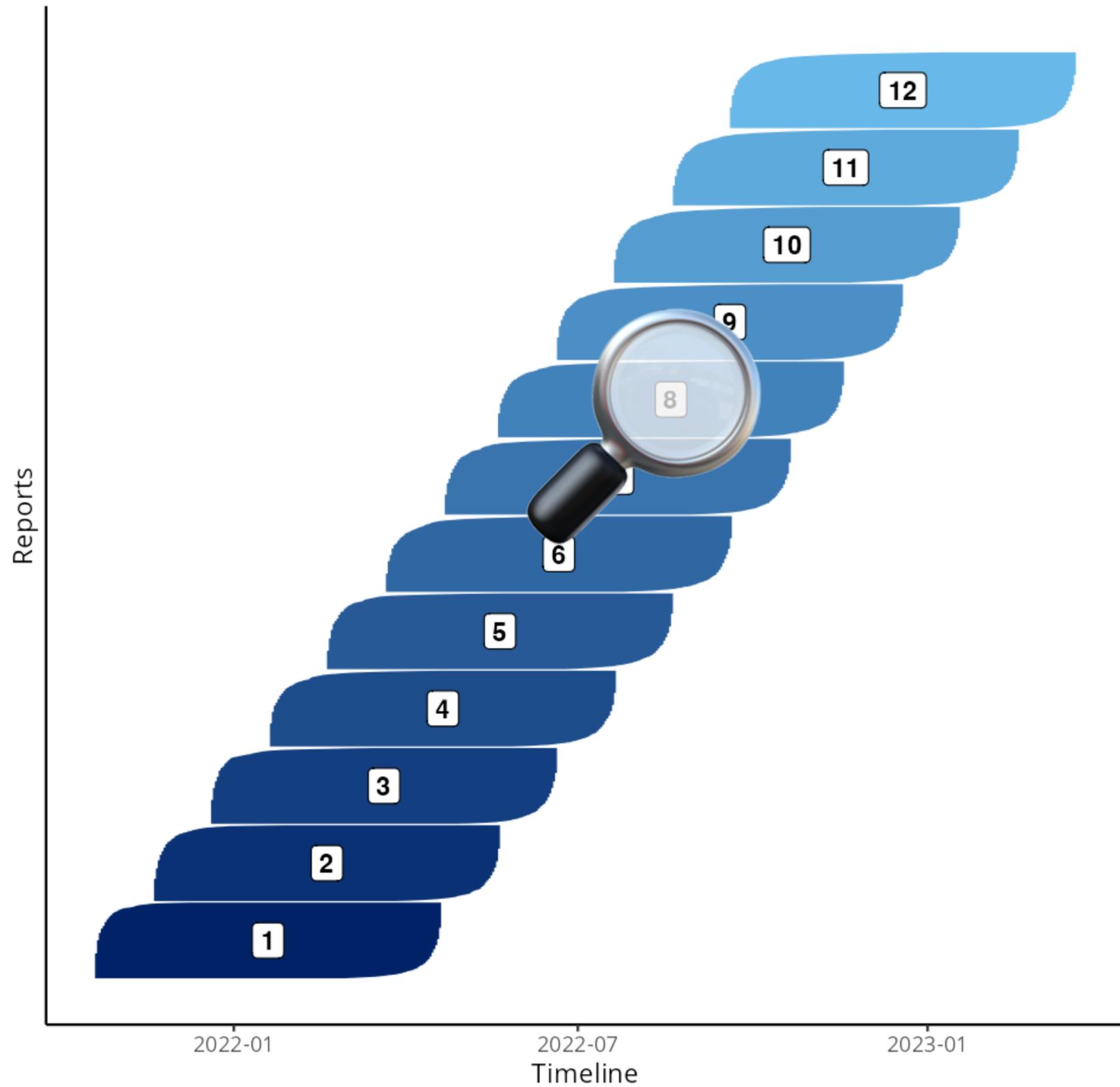
Transactions grouped by reports, twister edition



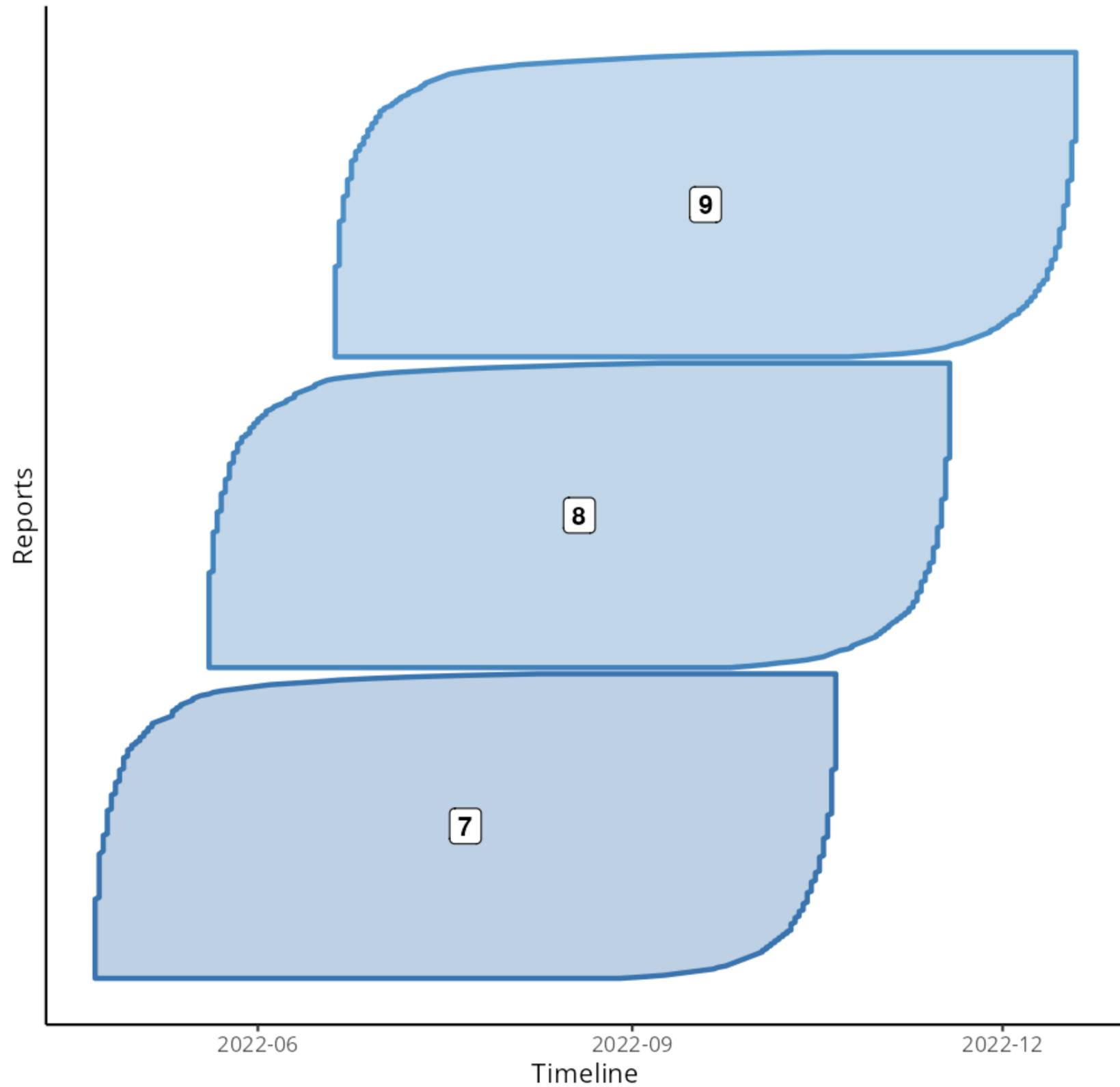
Transactions grouped by reports, twister edition



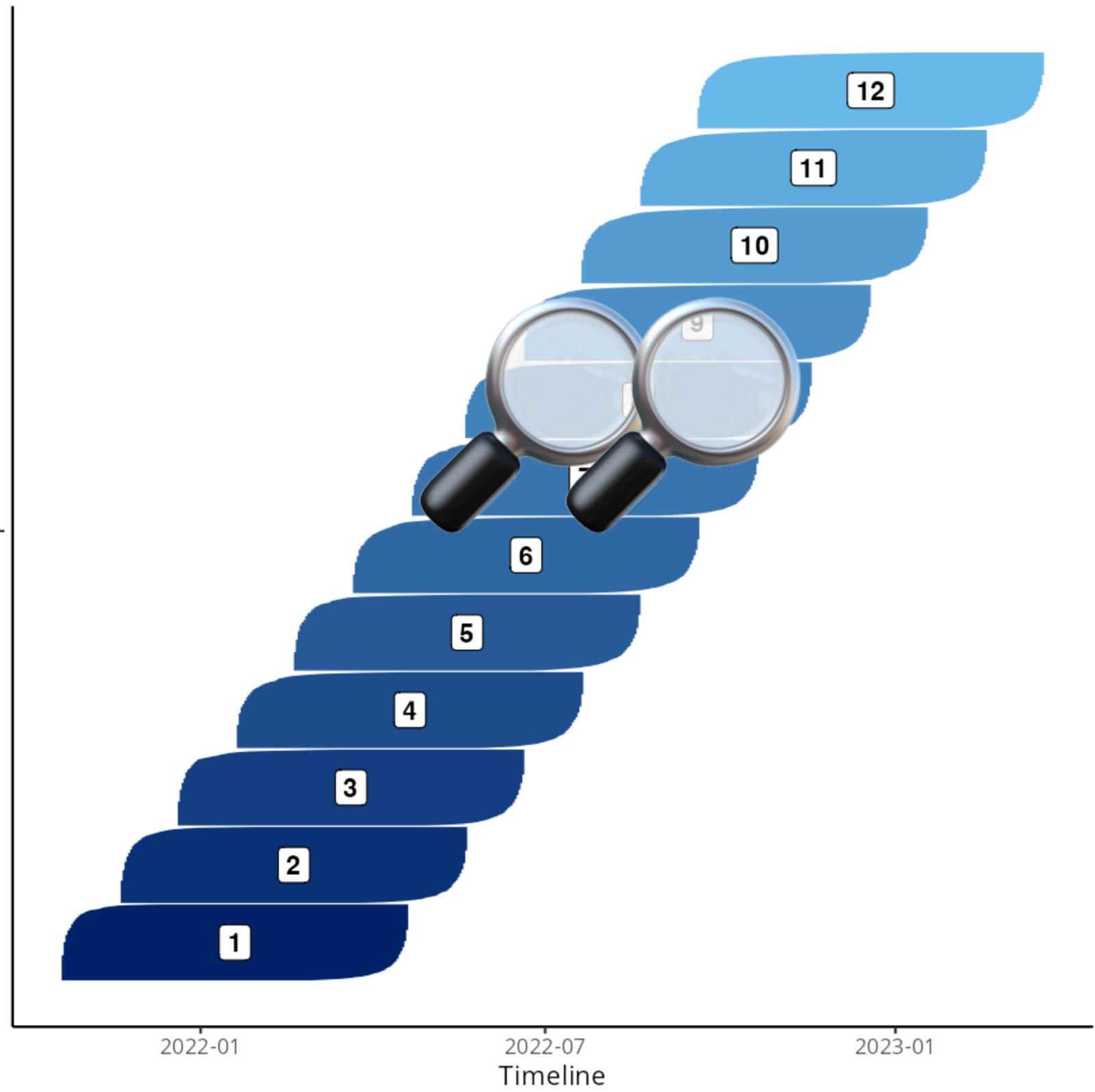
Transactions grouped by reports, twister edition



Transaction distribution within reports #7, #8, #9



Transactions grouped by reports, twister edition



Transaction distribution within reports #7, #8, #9



задача  
**Модели оттока**

1 000 000 руб

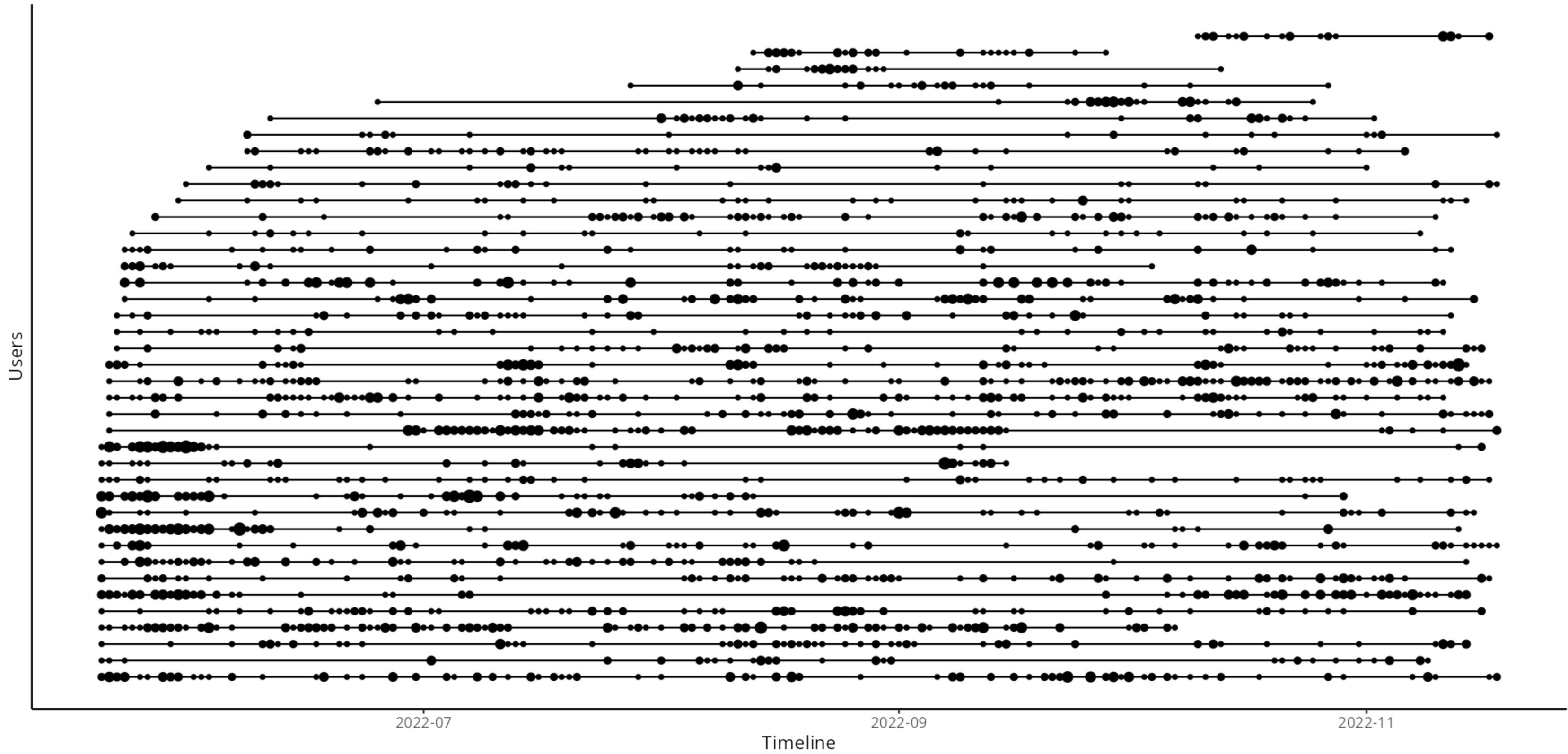
## Часть 3

Что на самом деле происходит в транзакциях? ~~Что предсказываем?~~

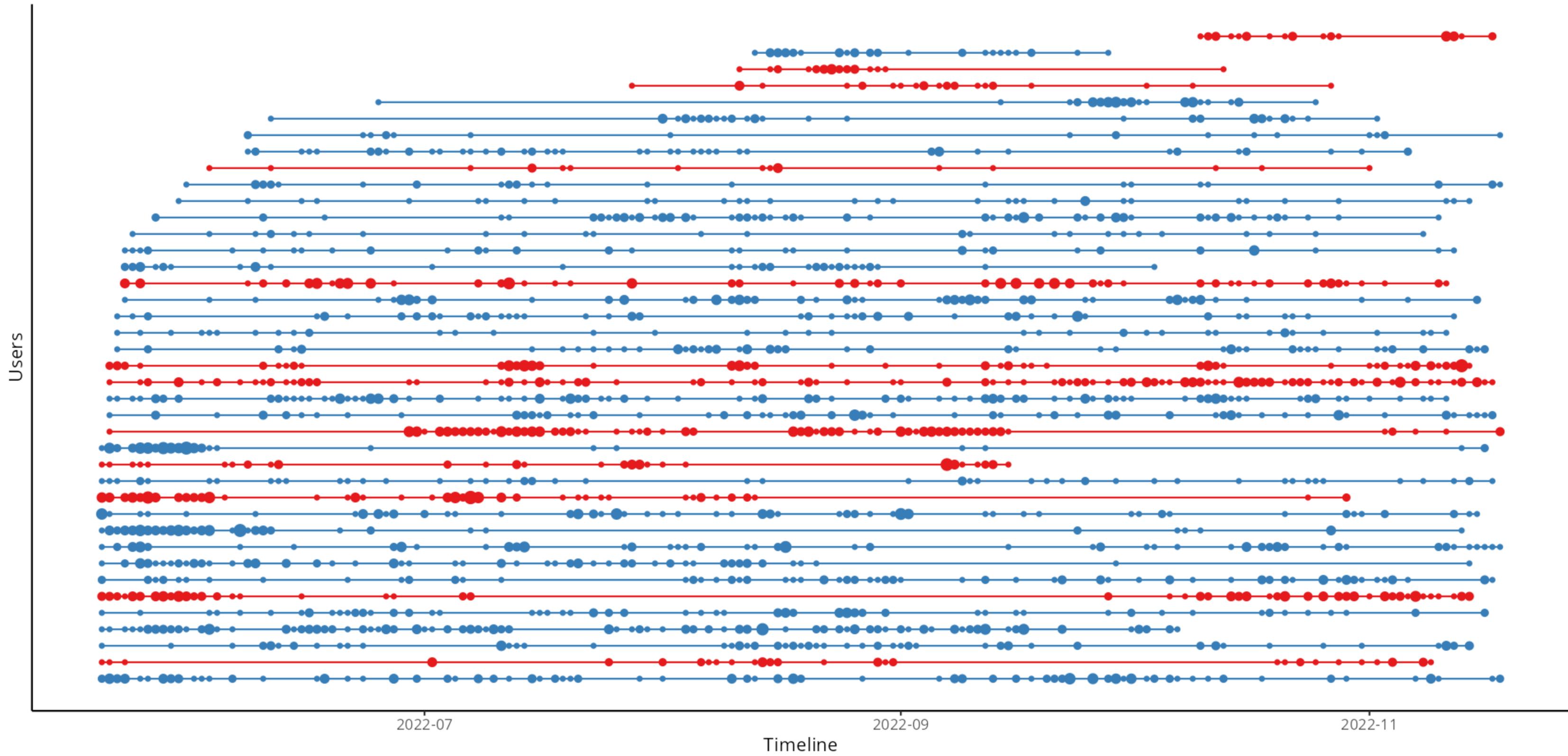
Transaction distribution for report #8, random 40 users



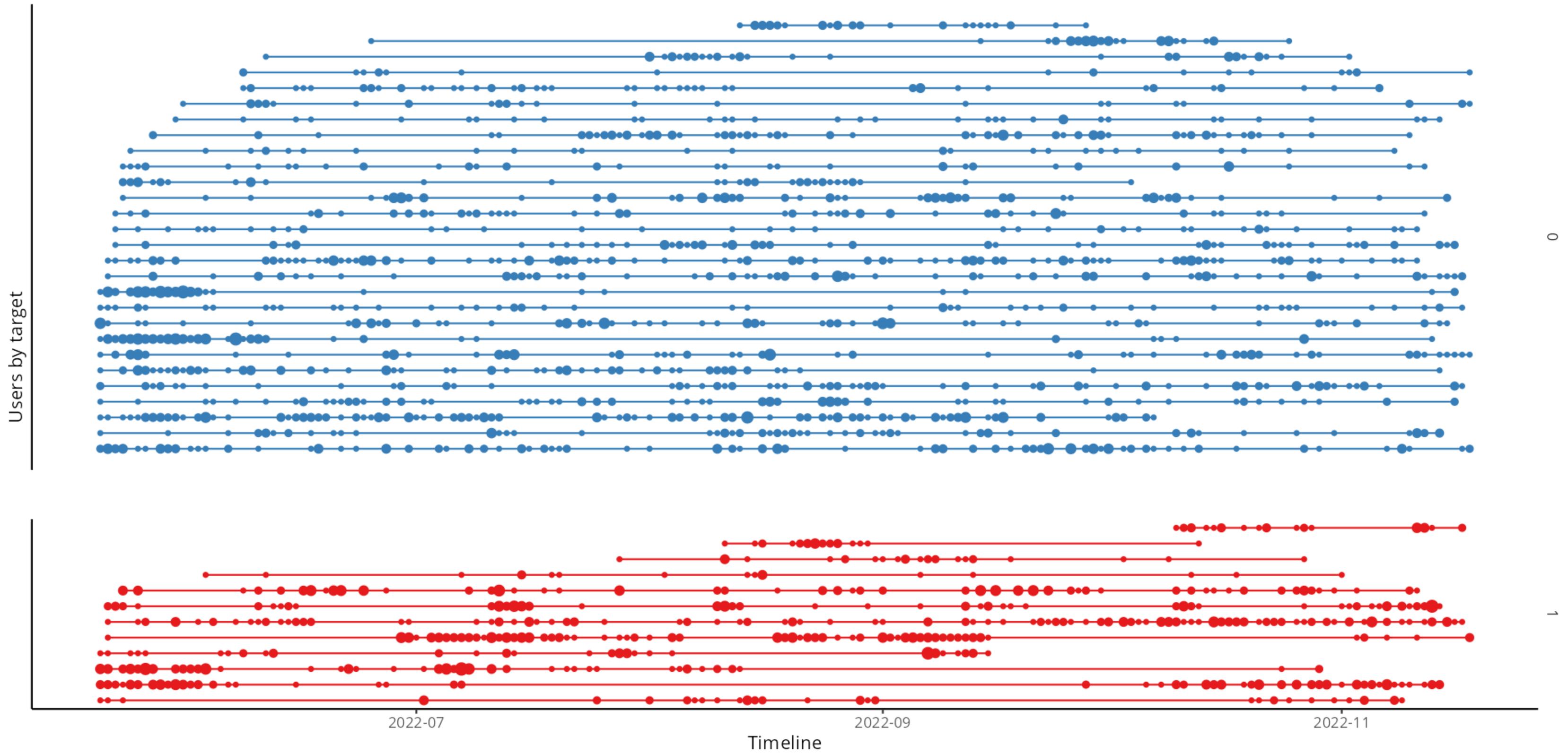
Transaction distribution for report #8, random 40 users



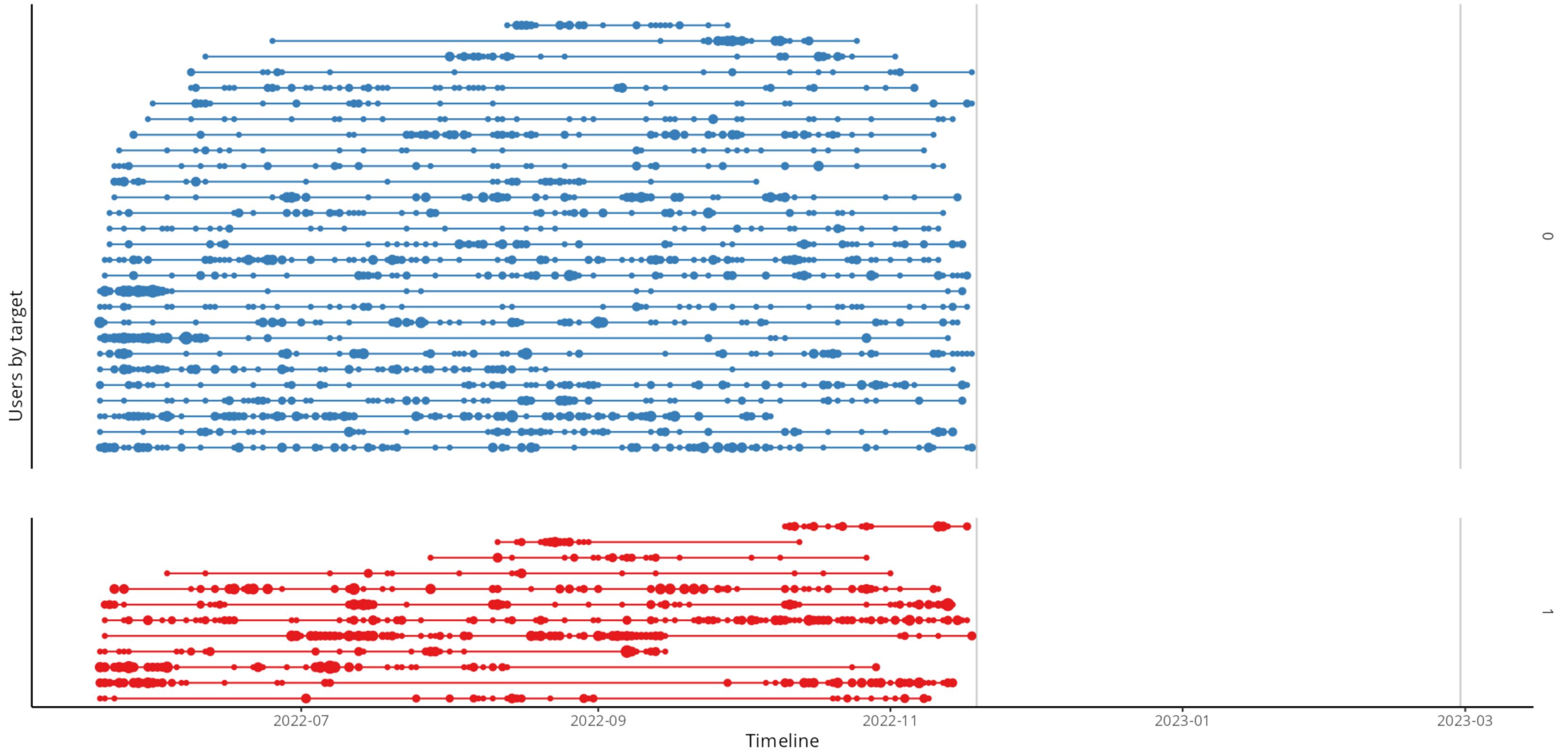
Transaction distribution for report #8, random 40 users



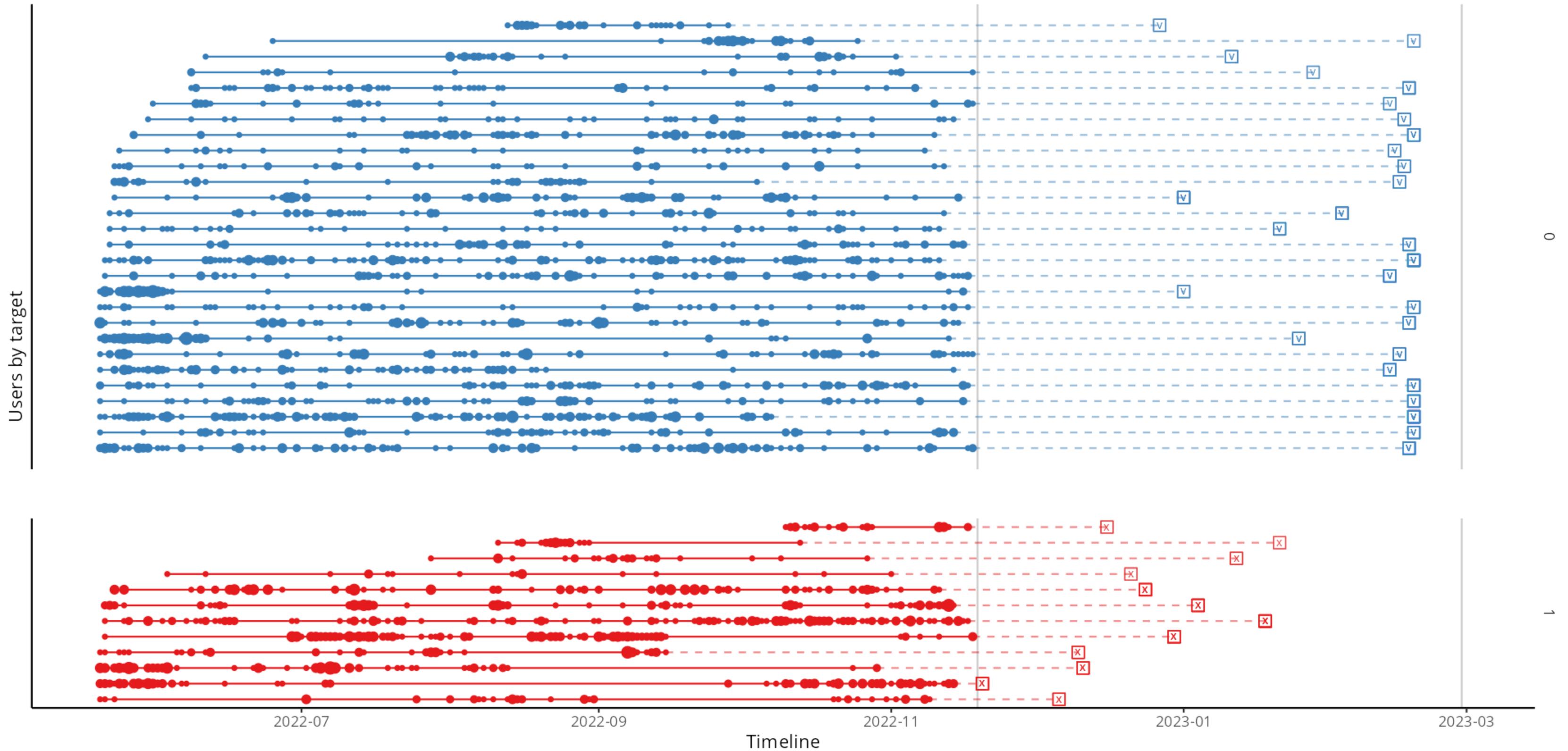
Transaction distribution for report #8, random 40 users



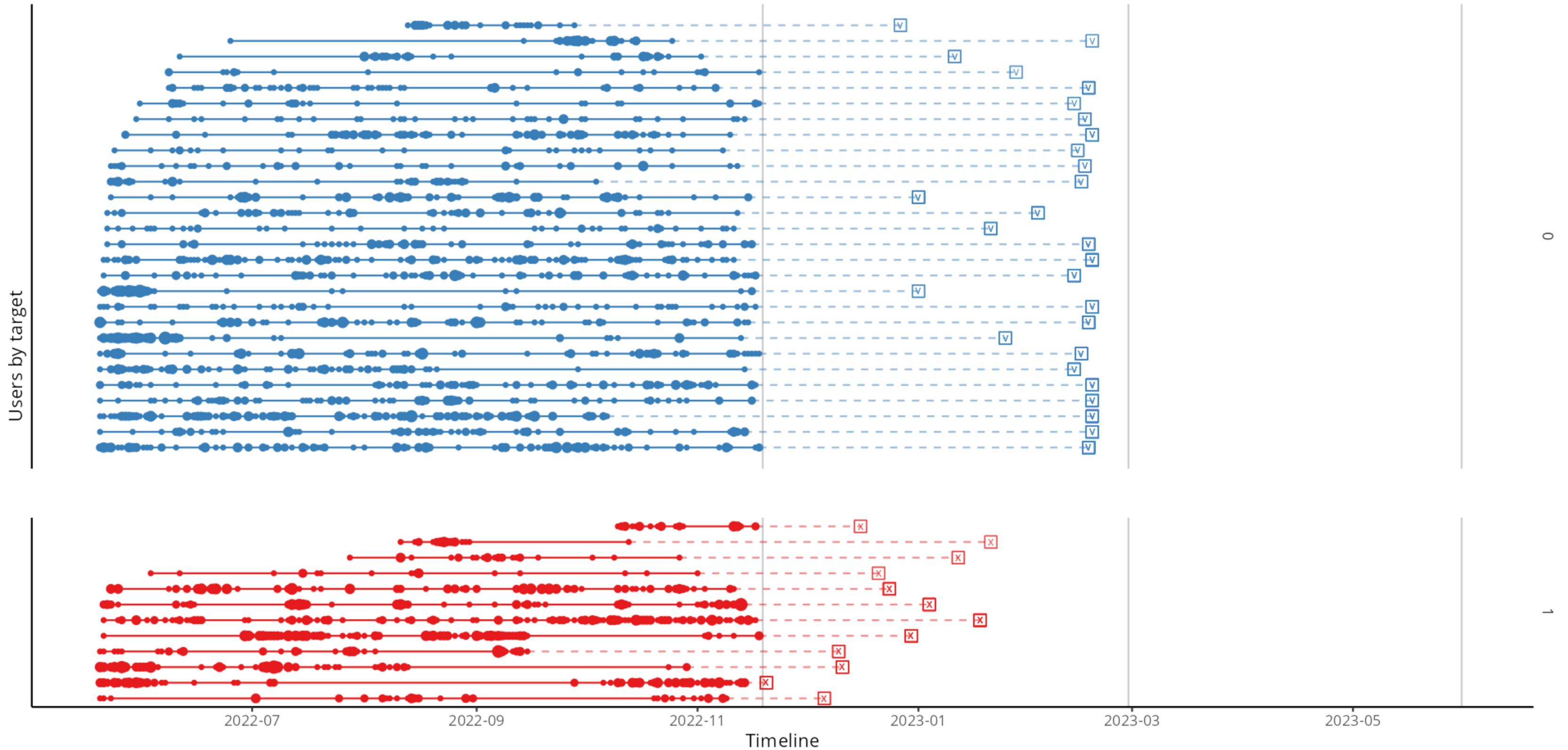
Transaction distribution for report #8, before report



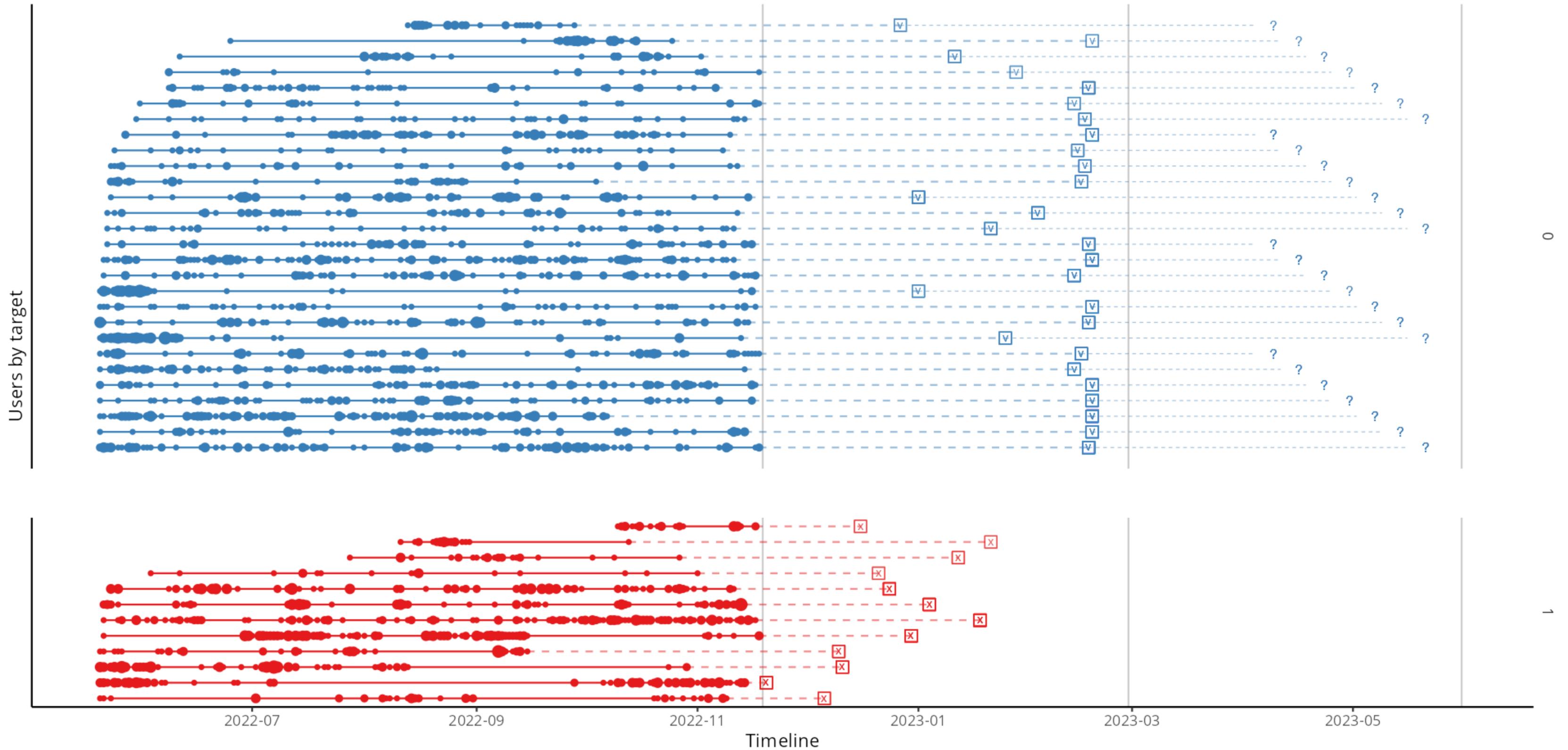
Transaction distribution for report #8, before report



Transaction distribution for report #8, after report

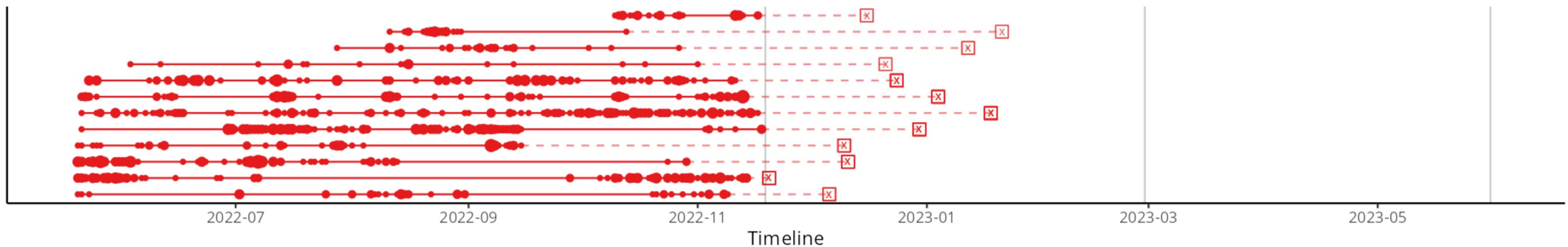


Transaction distribution for report #8, after report



Исходные метки оттока, используемые в этом соревновании, получены именно с помощью этого процесса. Как не трудно видеть, это определение подразумевает, что последняя транзакция клиента должна была состояться до даты отчета. В таком случае, можно считать что у нас есть 2 определения оттока:

- **Документальный** отток: вызревшие метки, по которым срабатывает бизнес-правило (через 60-90 дней после отчета).
- **Клинический** отток: событие, когда у клиента “остановился пульс” транзакций, и которое наступило не раньше даты последней транзакции (до отчета). А спустя 60-90 дней отсутствия транзакций, по этому клиенту запишут статус документального оттока.





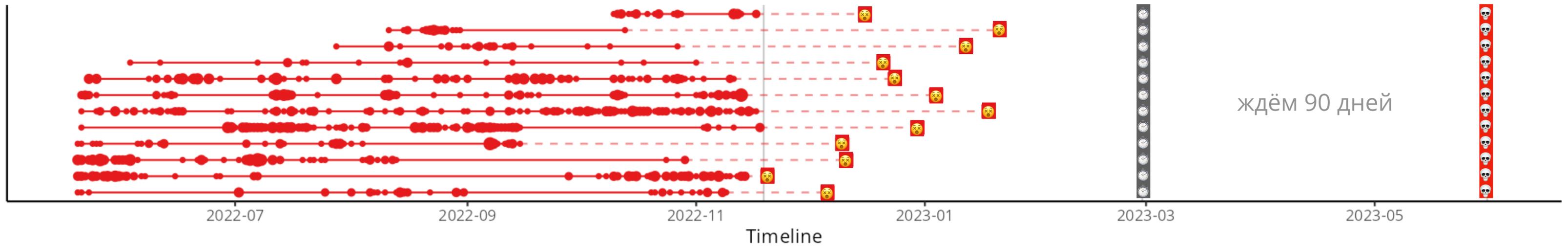
**“клинический”  
ОТТОК**



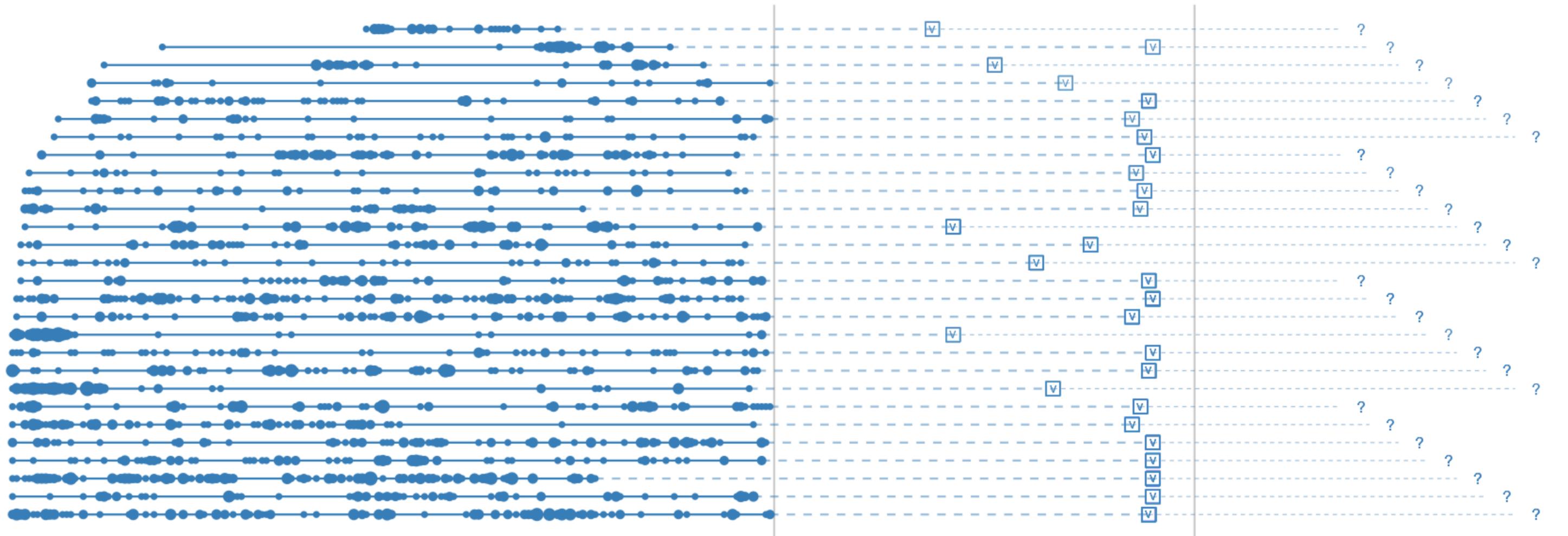
**дата  
отчета**



**“документальный”  
ОТТОК**



Users by target



последние  
транзакции  
до отчета



дата  
отчета



цензурированные  
данные

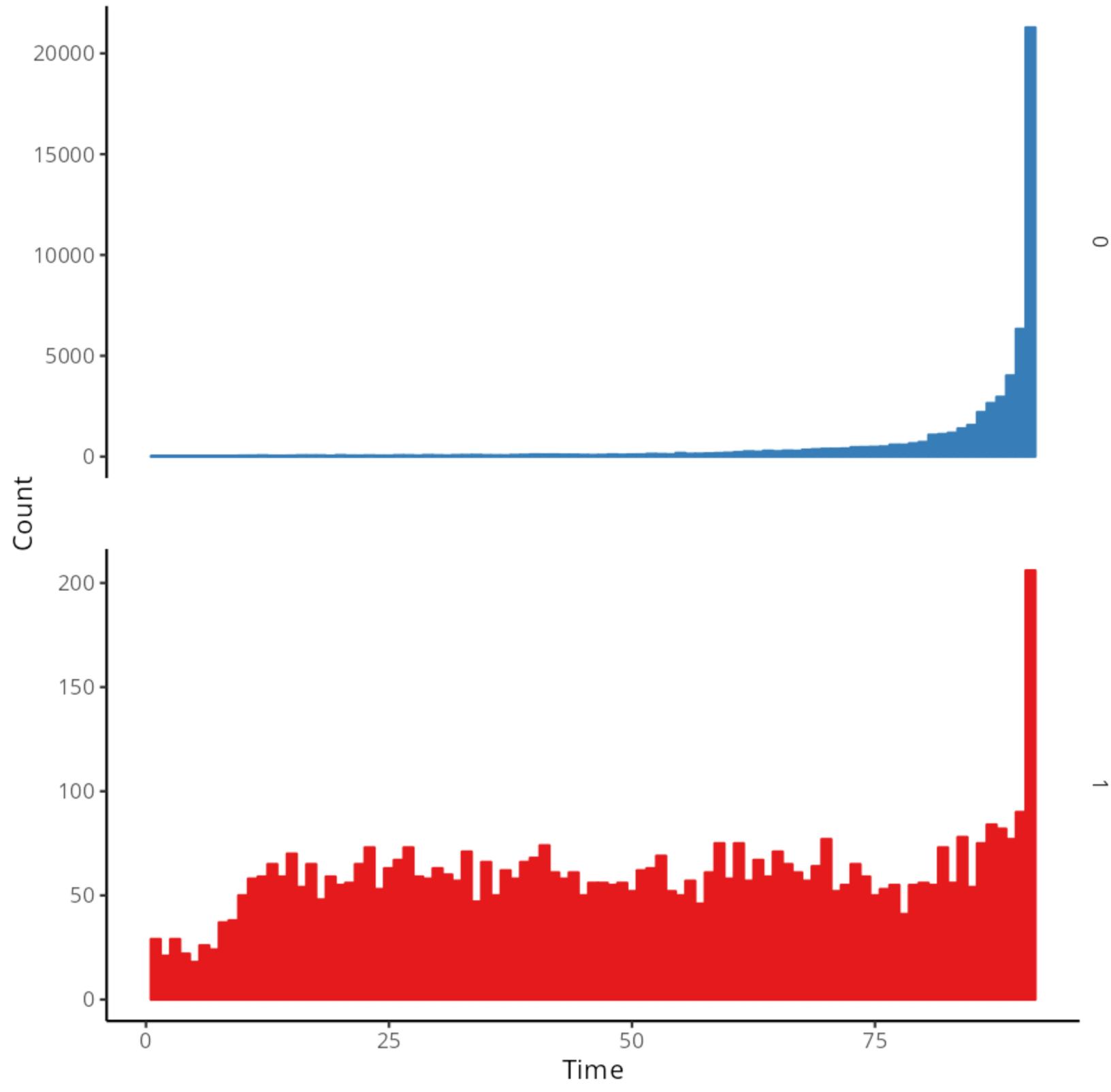
задача  
**Модели оттока**

1 000 000 руб

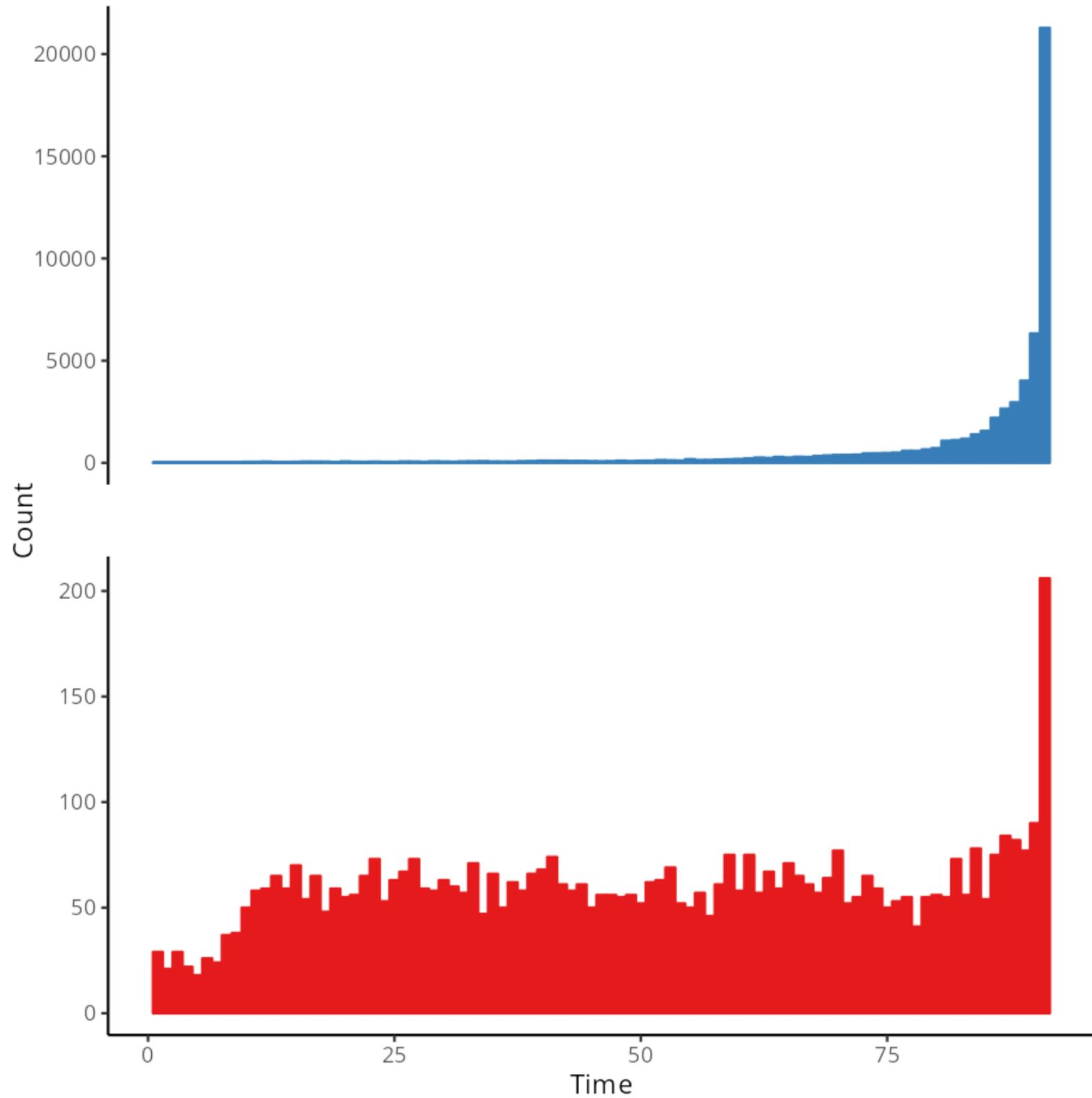
# Часть 4

Что еще мы можем предсказывать?

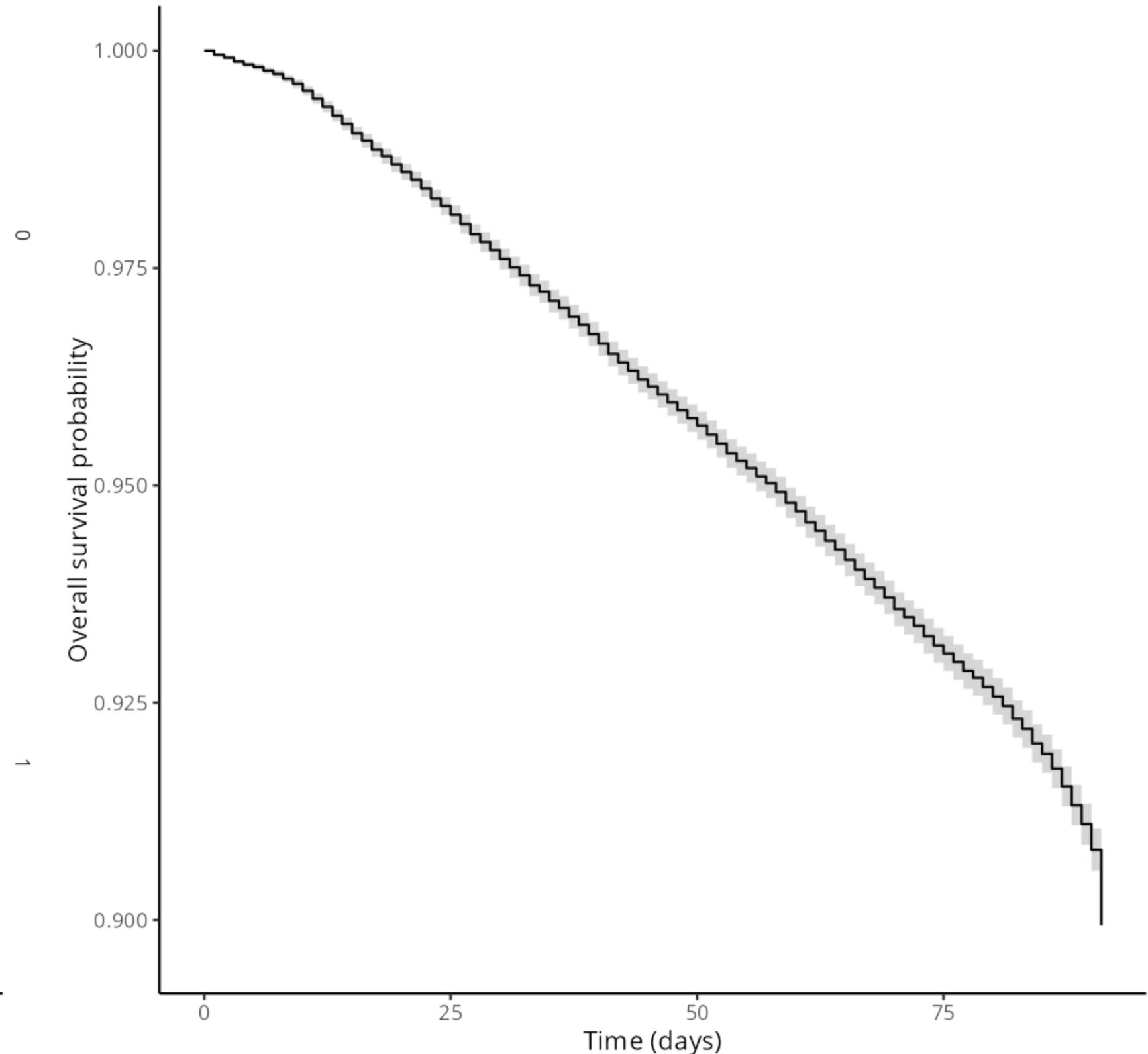
Time \* target distribution in train.csv?



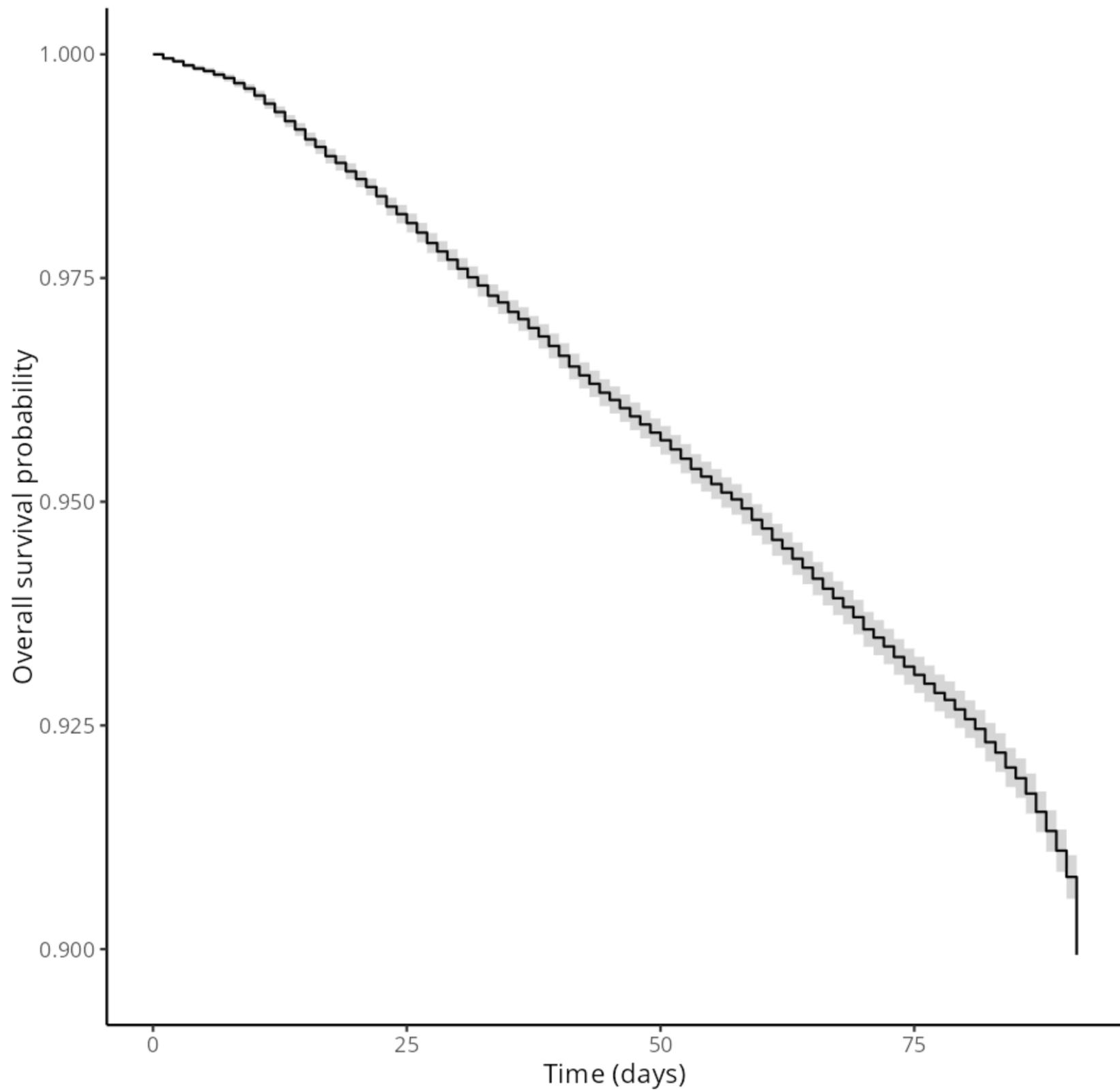
Time \* target distribution in train.csv?



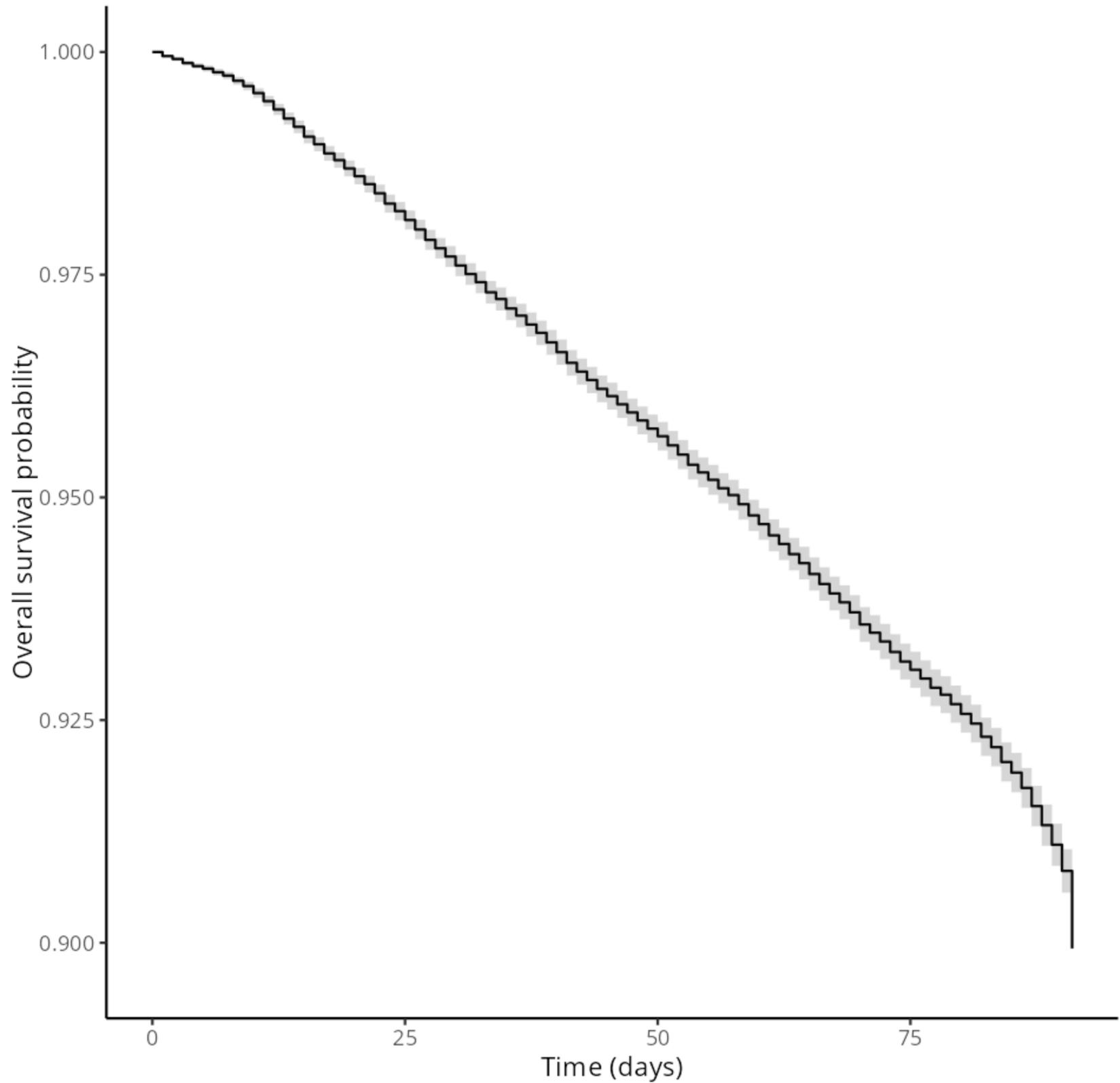
Global survival curve based on train.csv



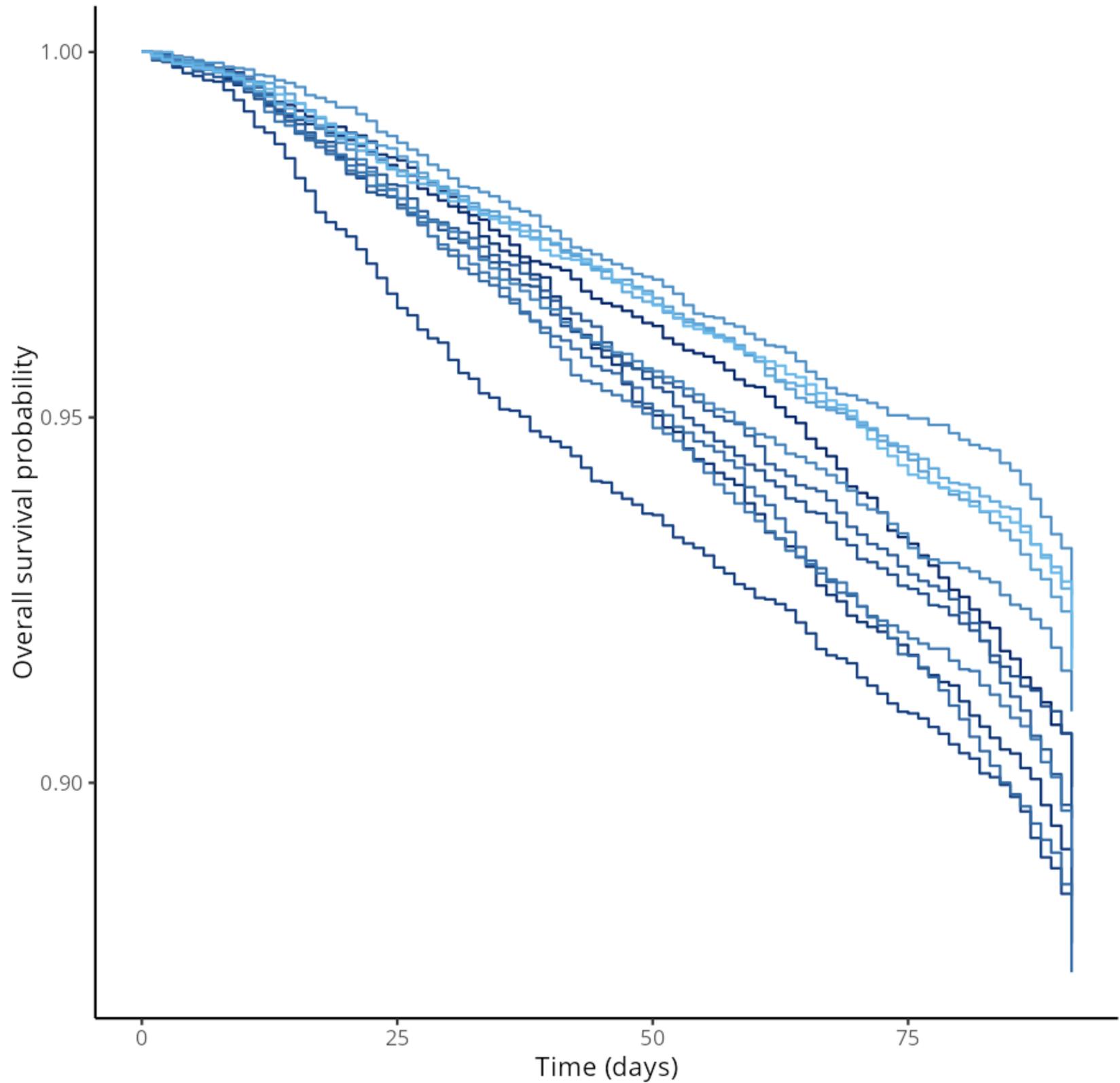
Global survival curve based on train.csv



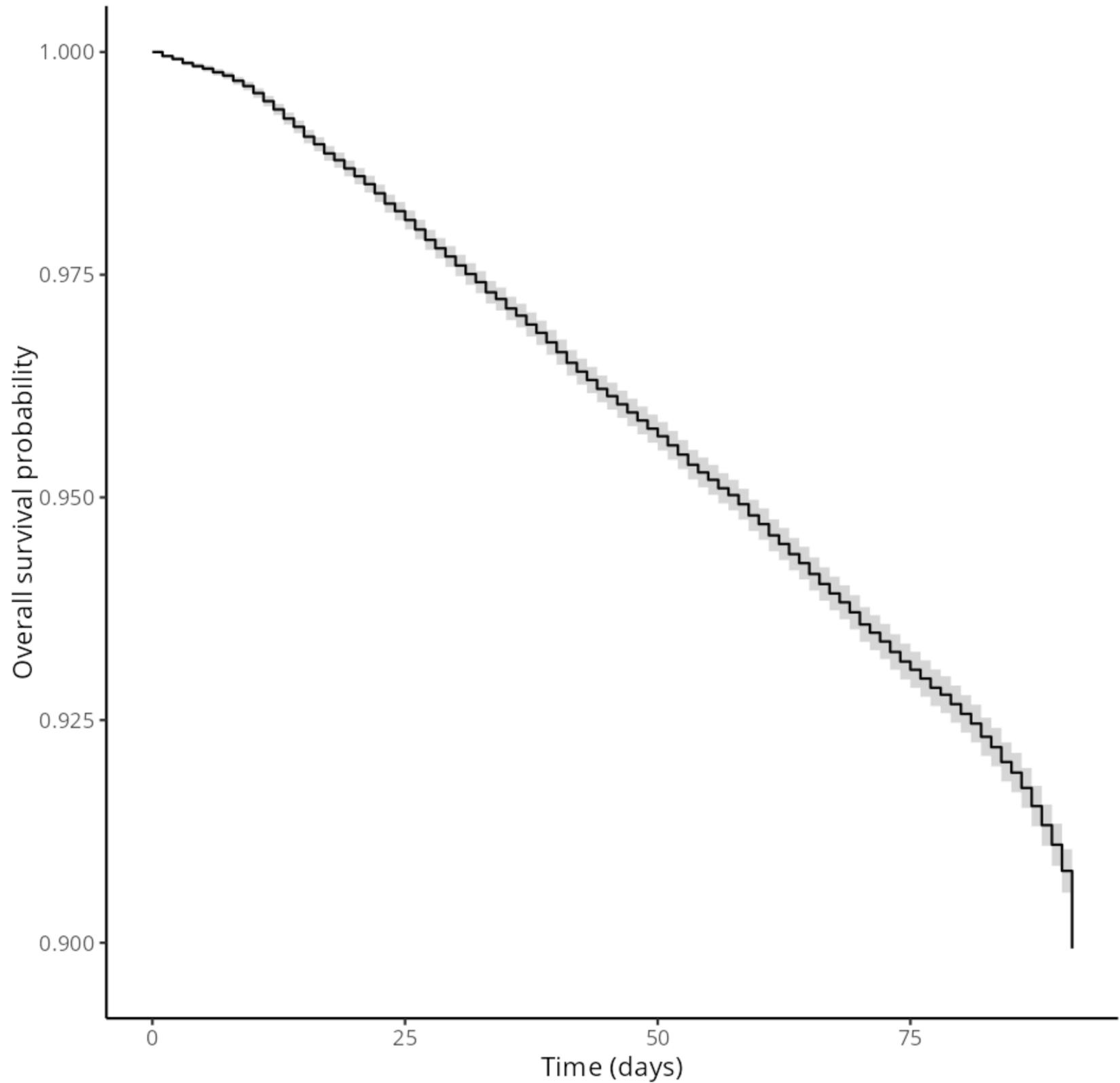
Global survival curve based on train.csv



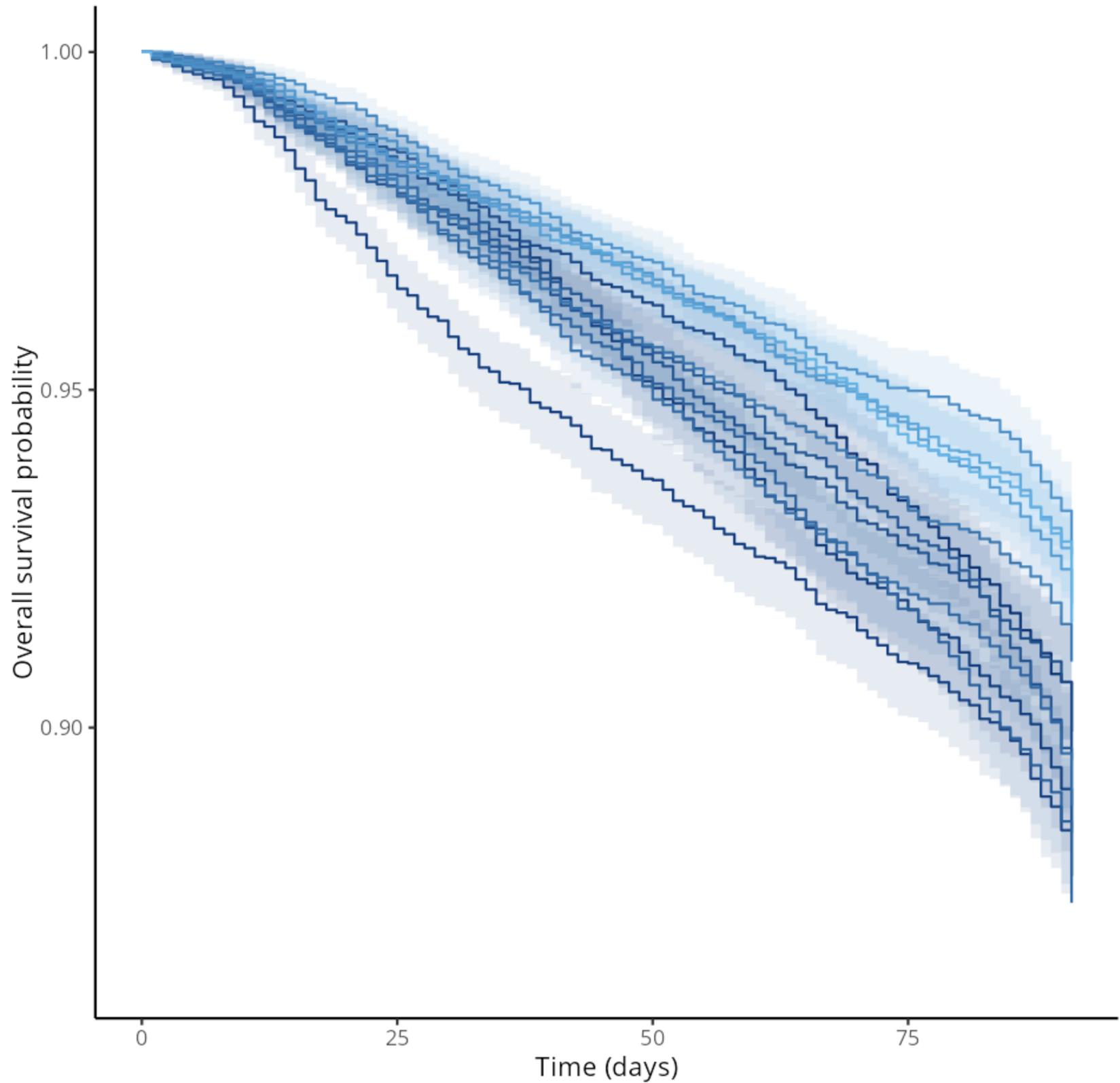
Report-wise survival curves: train.csv \* clients.csv



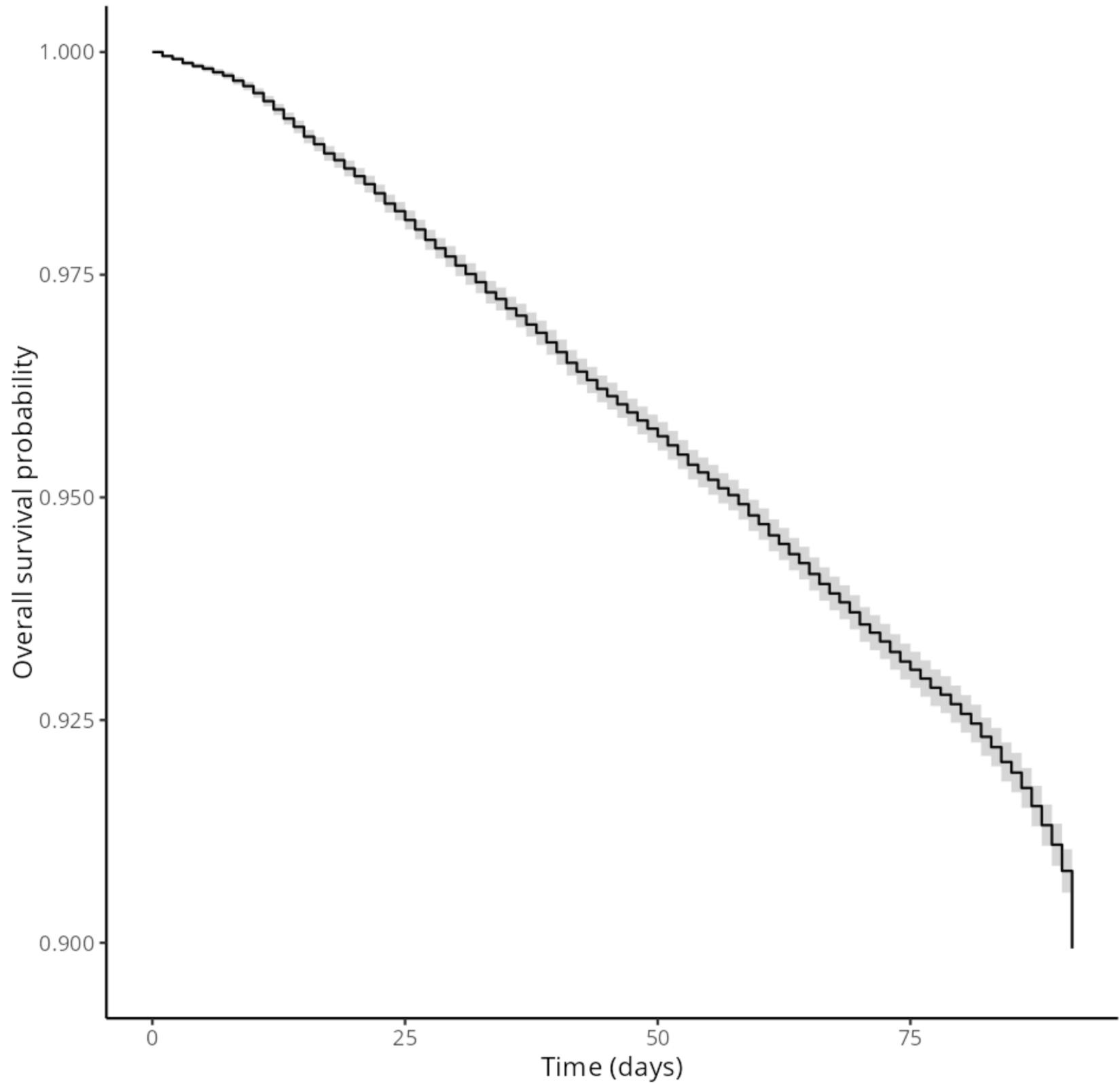
Global survival curve based on train.csv



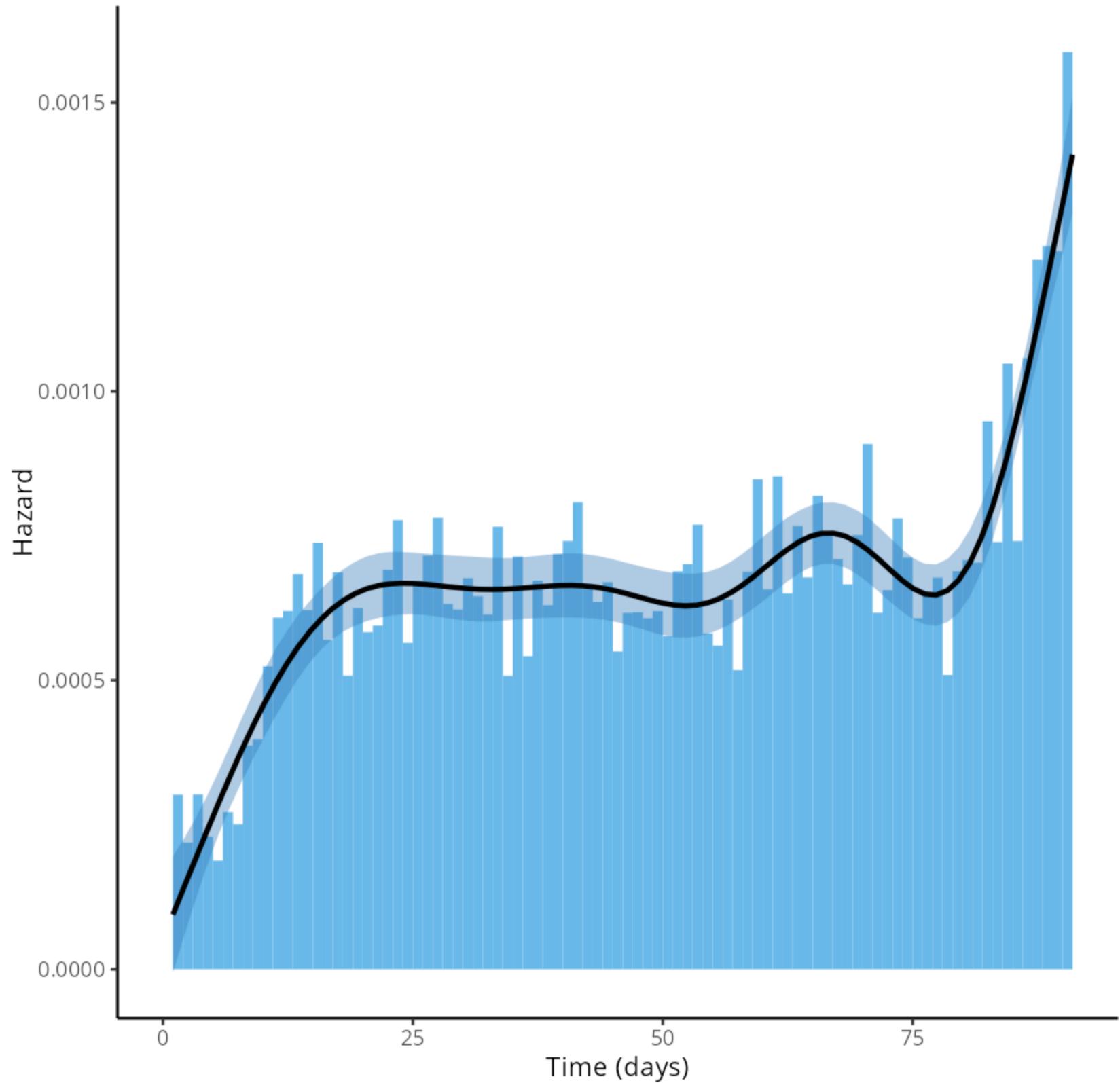
Report-wise survival curves: train.csv \* clients.csv



Global survival curve based on train.csv



Empirical hazard function estimate (GP)



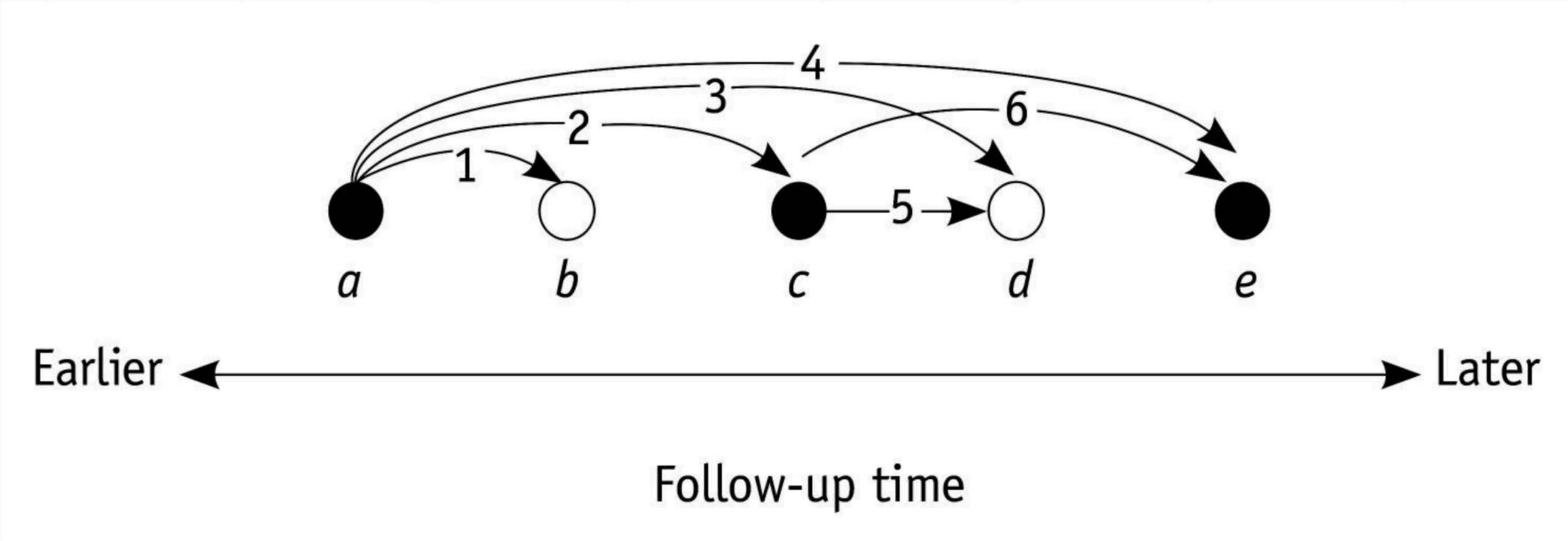
задача  
**Модели оттока**

1 000 000 руб

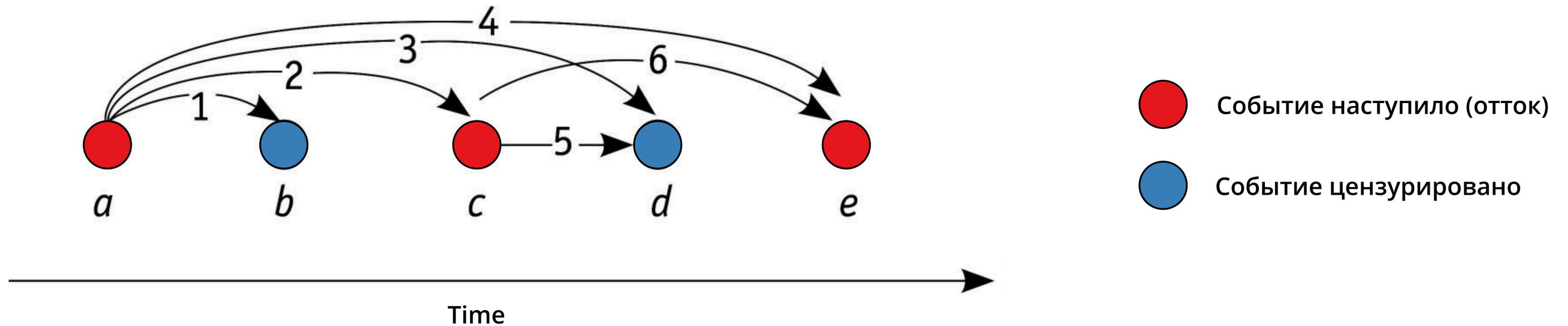
# Часть 5

Про метрику

# Concordance Index



# Concordance Index



- 5 клиентов => 10 возможных пар для сравнения
- Только 6 пар из 10 можем сравнить из-за цензурирования
- Наглядно: **b** и **d** не сравнимы так как оба цензурированы

задача  
**Модели оттока**

1 000 000 руб

# Часть 6

Куда копать?

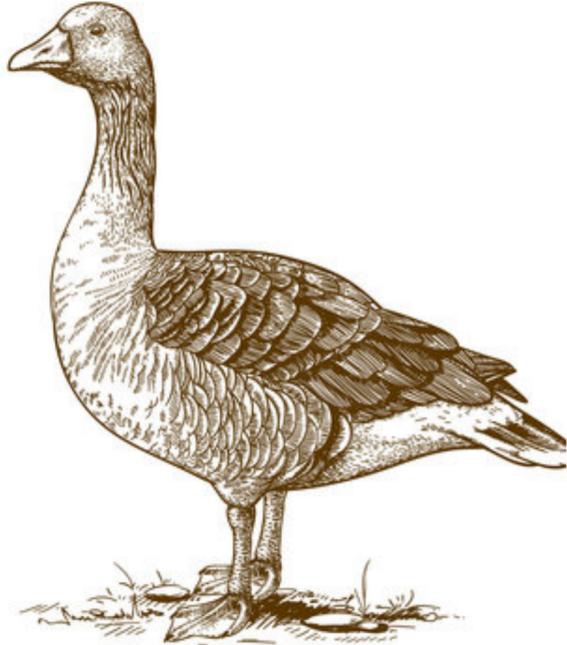
## Статьи:

- **[Arxiv, 2023]** *Deep Learning for Survival Analysis: A Review* — качественный обзор современных нейросетевых архитектур и постановок в Time-to-Event задачах.
- **[NCBI, 2021]** *Review of Statistical Methods for Evaluating the Performance of Survival or Other Time-to-Event Prediction Models (from Conventional to Deep Learning Approaches)* — обзор метрик, используемых в Time-to-Event задачах.
- **[ACM, 2019]** *Machine Learning for Survival Analysis: A Survey* — обзор применения машинного обучения в Time-to-Event задачах, может использоваться как пособие для начинающих.
- **[Git, 2016]** *WTTE-RNN - Less hacky churn prediction* — легендарный пост-статья про проблемы задачи предсказания оттока, и как к задаче стоило бы подходить. **(выбор редакции ♥)**
- **[Chalmers Thesis, 2017]** *WTTE-RNN : Weibull Time To Event Recurrent Neural Network A model for sequential prediction of time-to-event in the case of discrete or continuous censored data, recurrent events or time-varying covariates* — "сопроводительный" диссер, по материалам которого была написана статья выше.

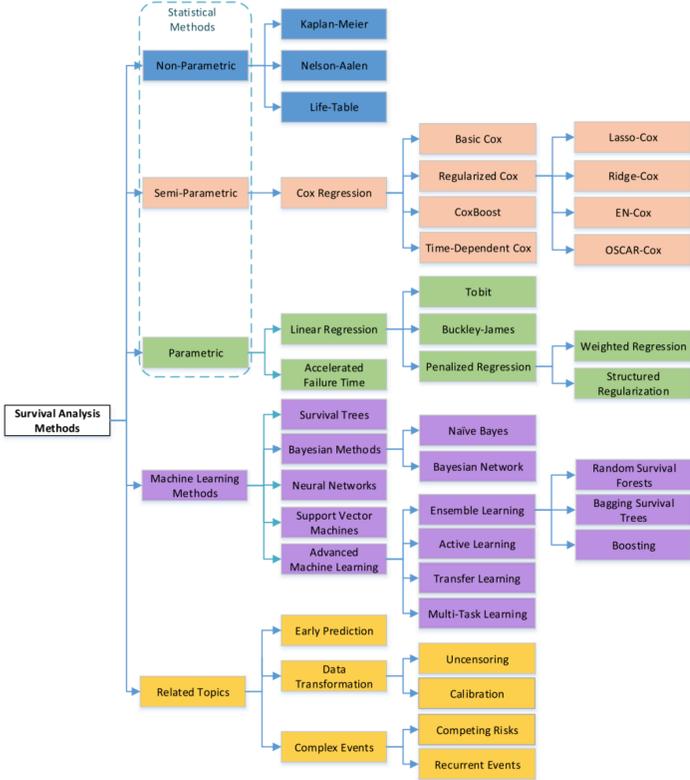
## Библиотеки:

- **Survival analysis with Catboost** — почти не секретный tutorial по тому, как запускать Cox и AFT модели на Catboost
- **scikit-survival** — библиотека по Survival / Time-to-Event задача на Python, на основе scikit-learn
- **auton-survival** — актуальная библиотека по Survival / Time-to-Event анализу на Python, включая нейросетевые модели от CMU
- **lifelines** — общая библиотека по классическому Survival Analysis на Python с основным статистическим инструментарием
- **(давно не обновлялась) PySurvival** — библиотека по Survival / Time-to-Event задачам на Python. В том числе с tutorialом по оттоку (churn)
- **(давно не обновлялась) PyCox** — библиотека по нейросетевым моделям для Time-to-Event задач на основе pytorch

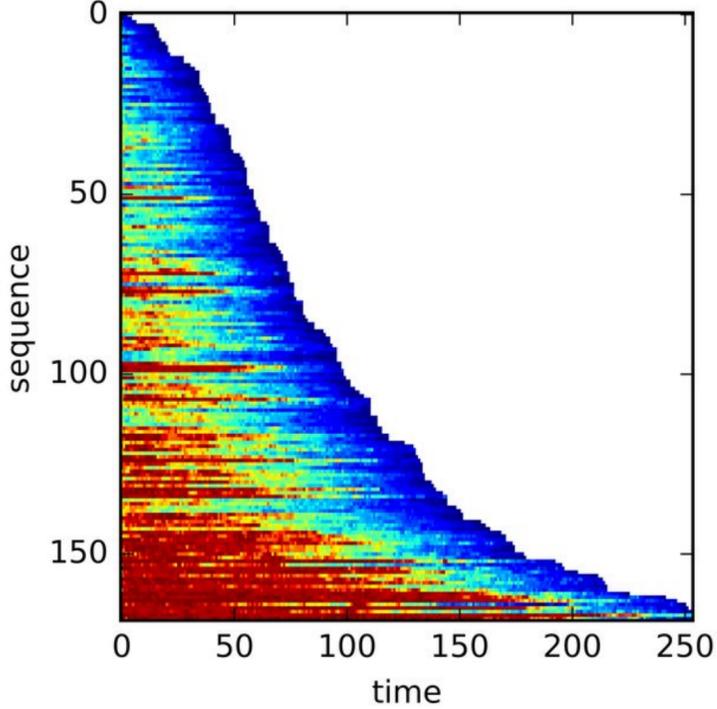
# Куда копать?



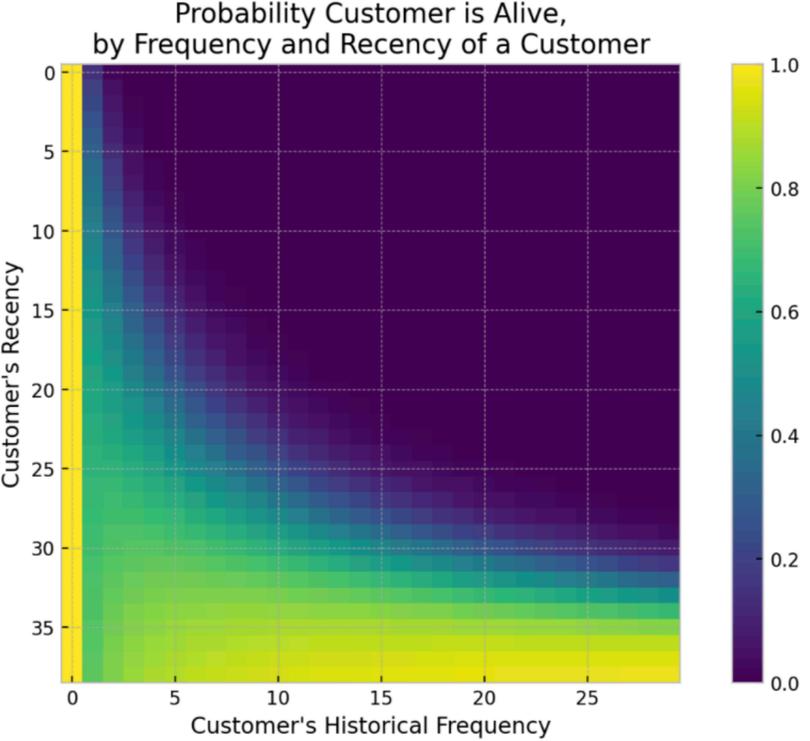
1. Включить "гуся"



2. Survival ML



3. Survival DL



4. Bayes++ BTYD

В ТРЯСИНЕ ЗАБОТ

# Включить гуся

Найди гуся, чтобы включить заботу



202 m

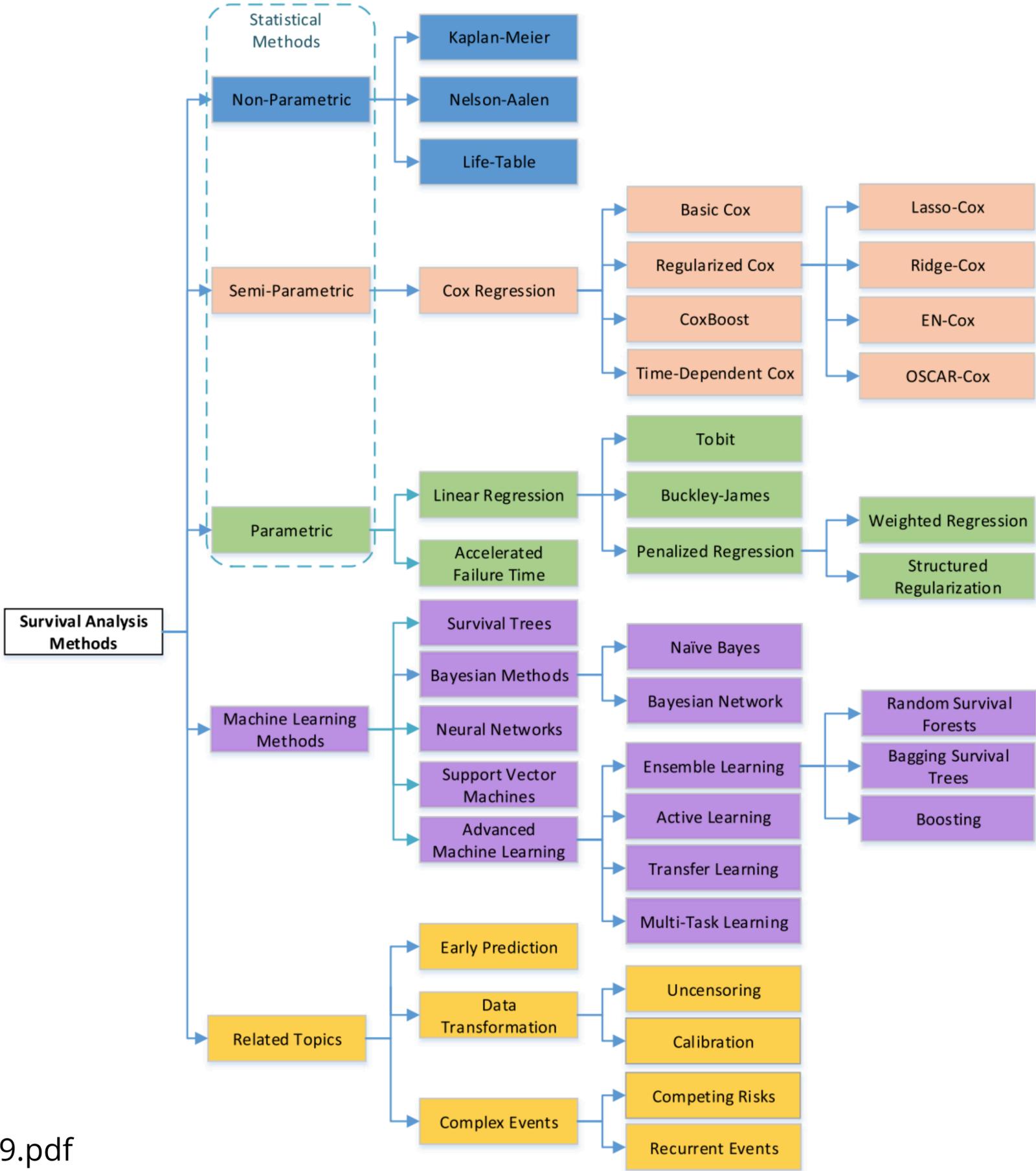
18



0



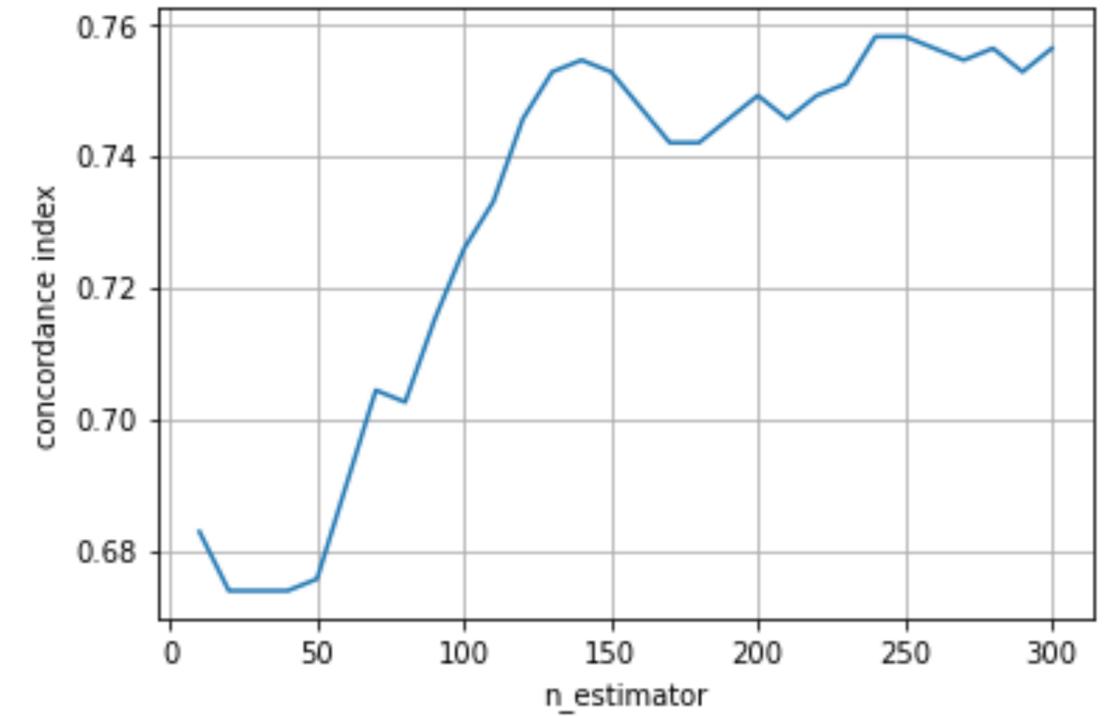
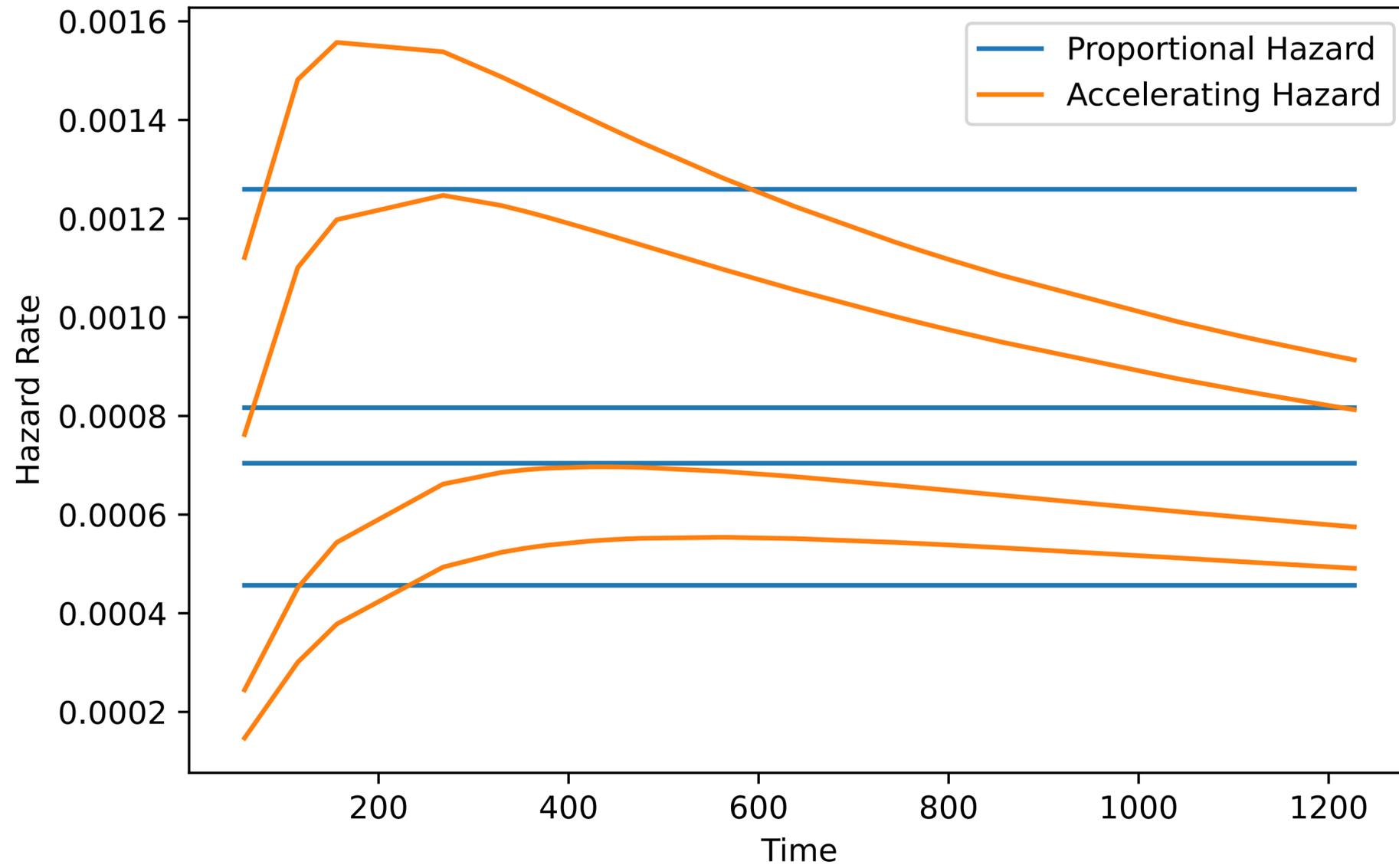
# Survival ML



# Survival ML



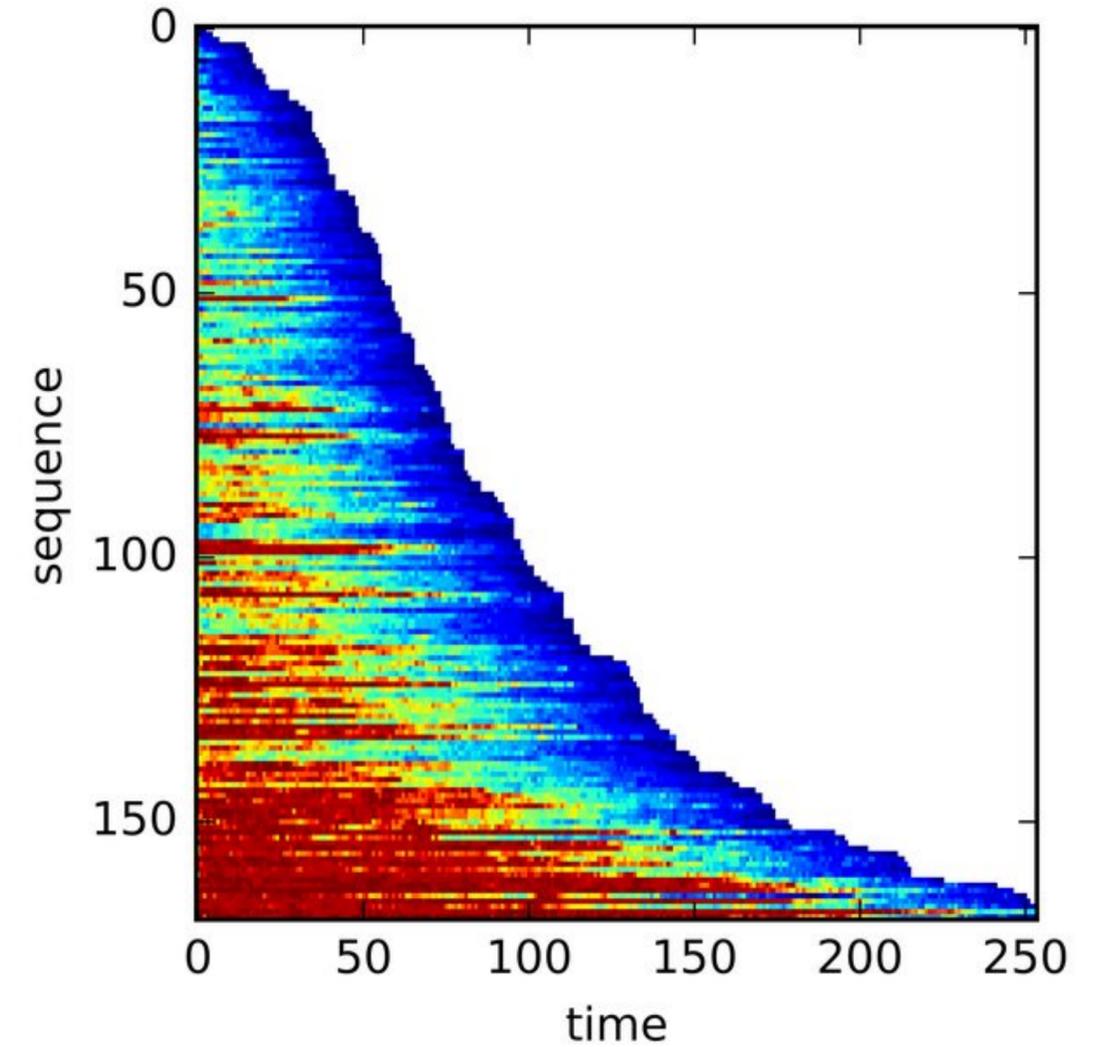
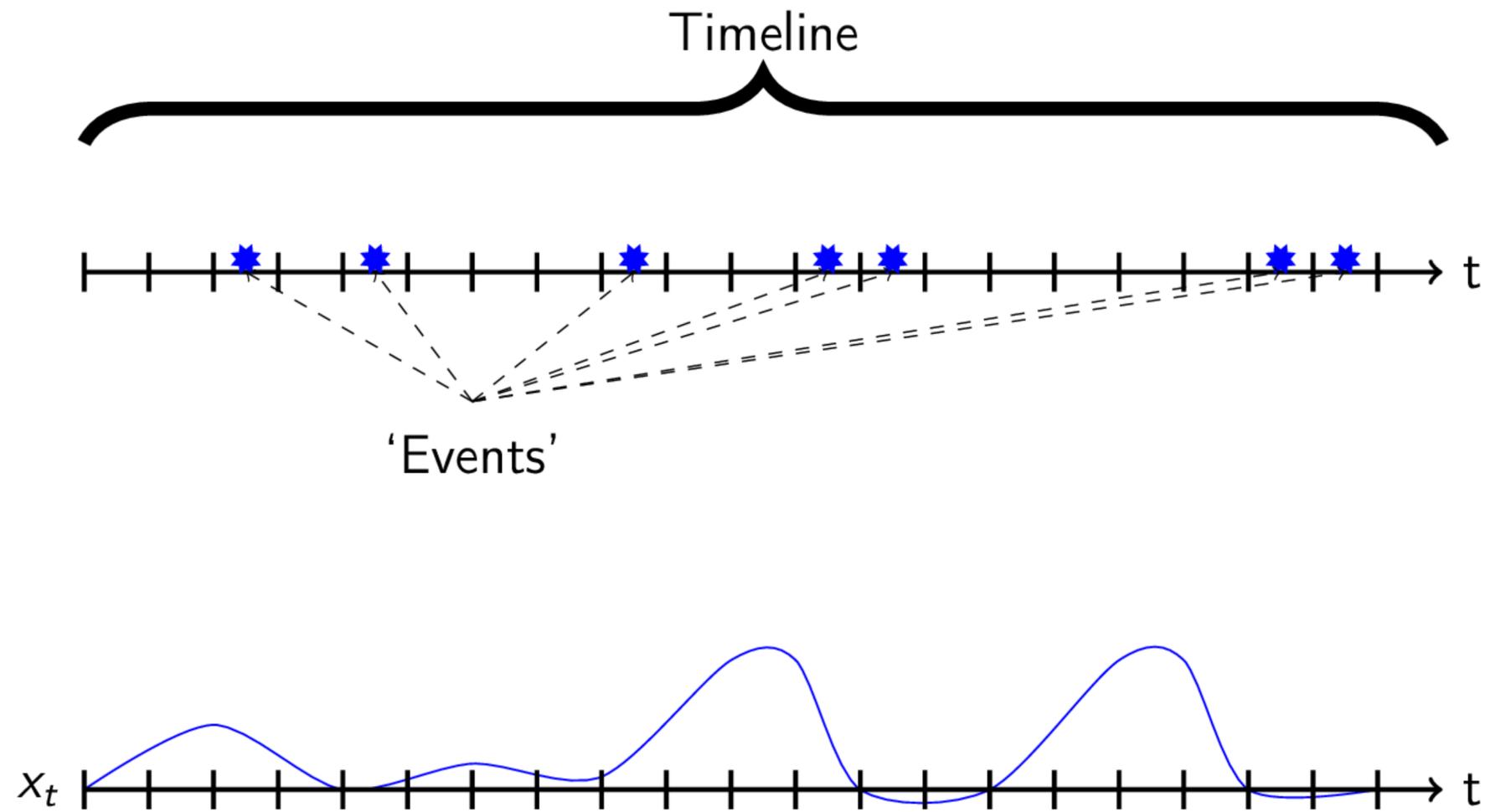
# Survival ML



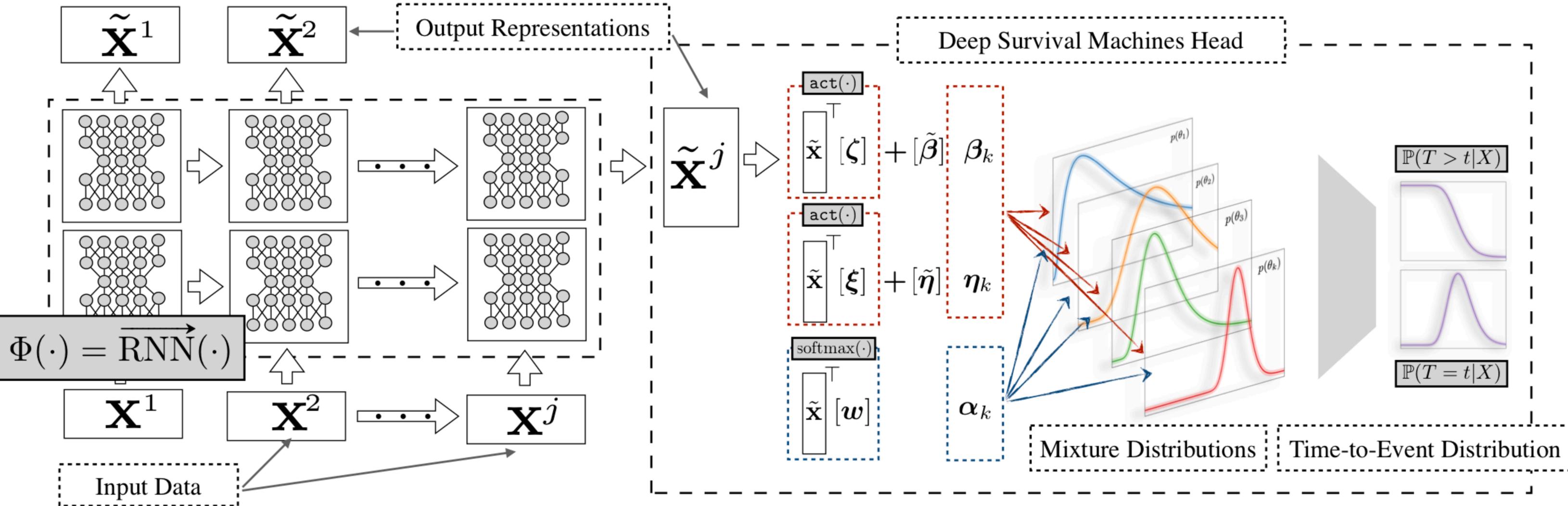
# Survival DL

## WTTE-RNN - Less hacky churn prediction

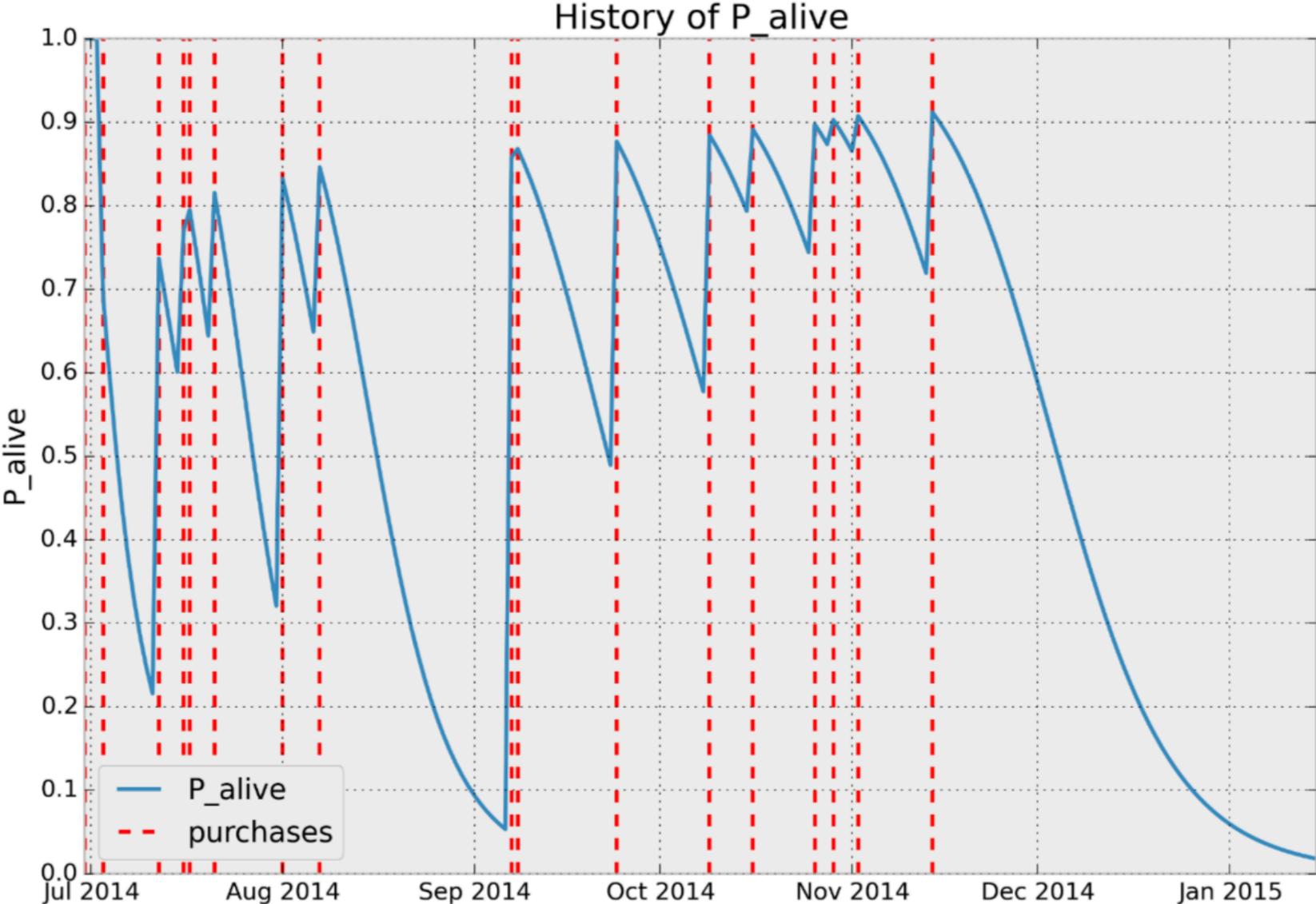
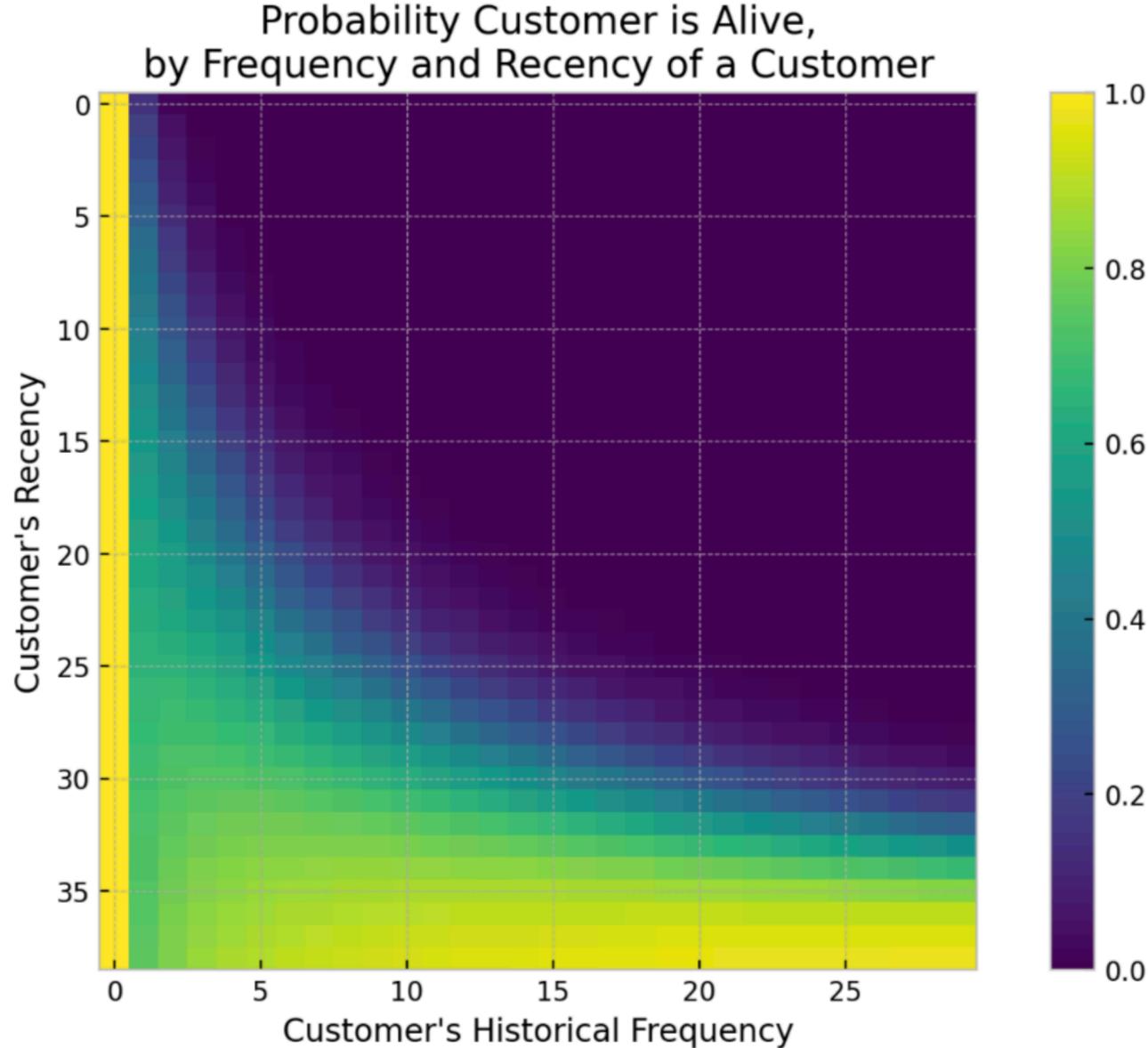
22 Dec 2016



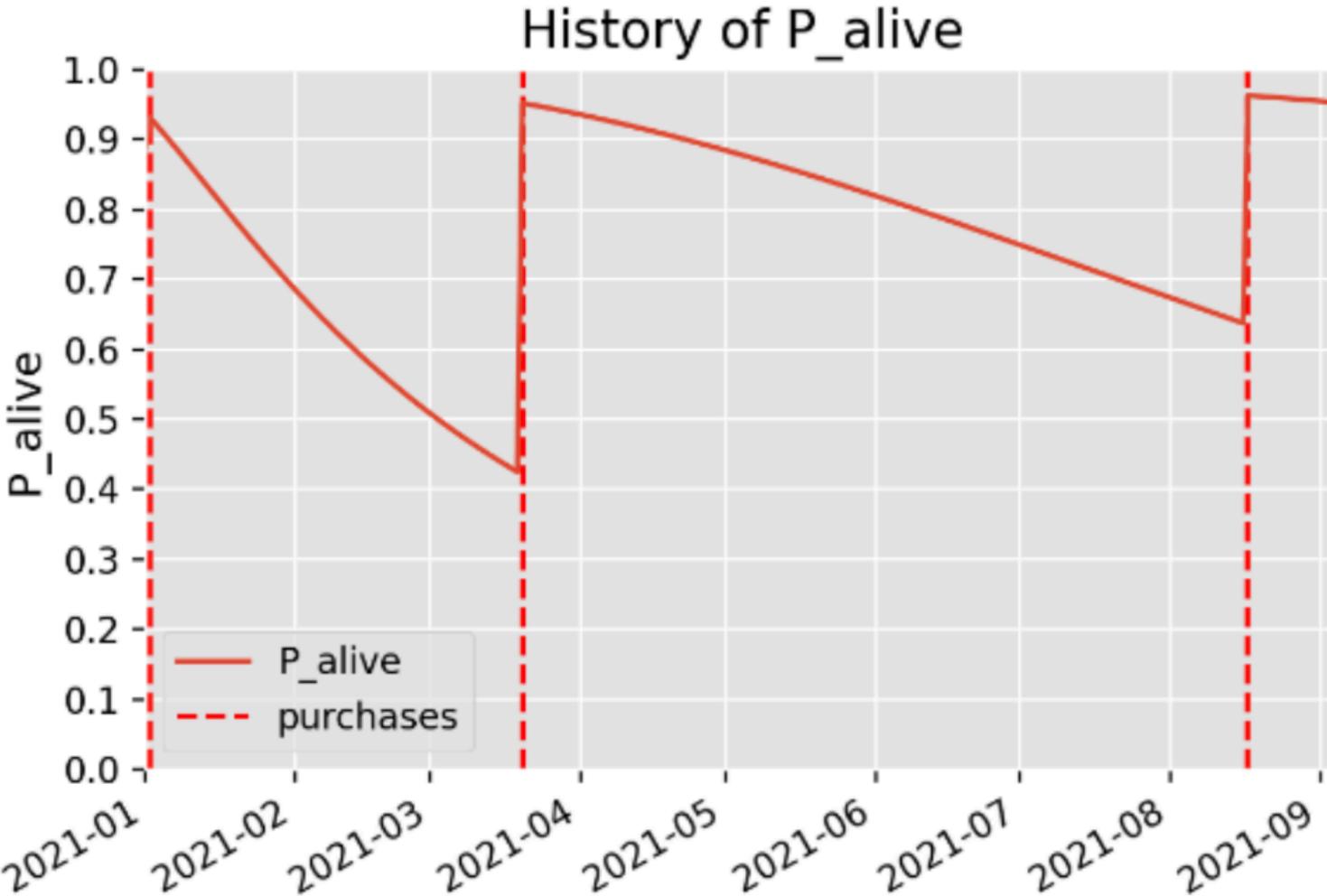
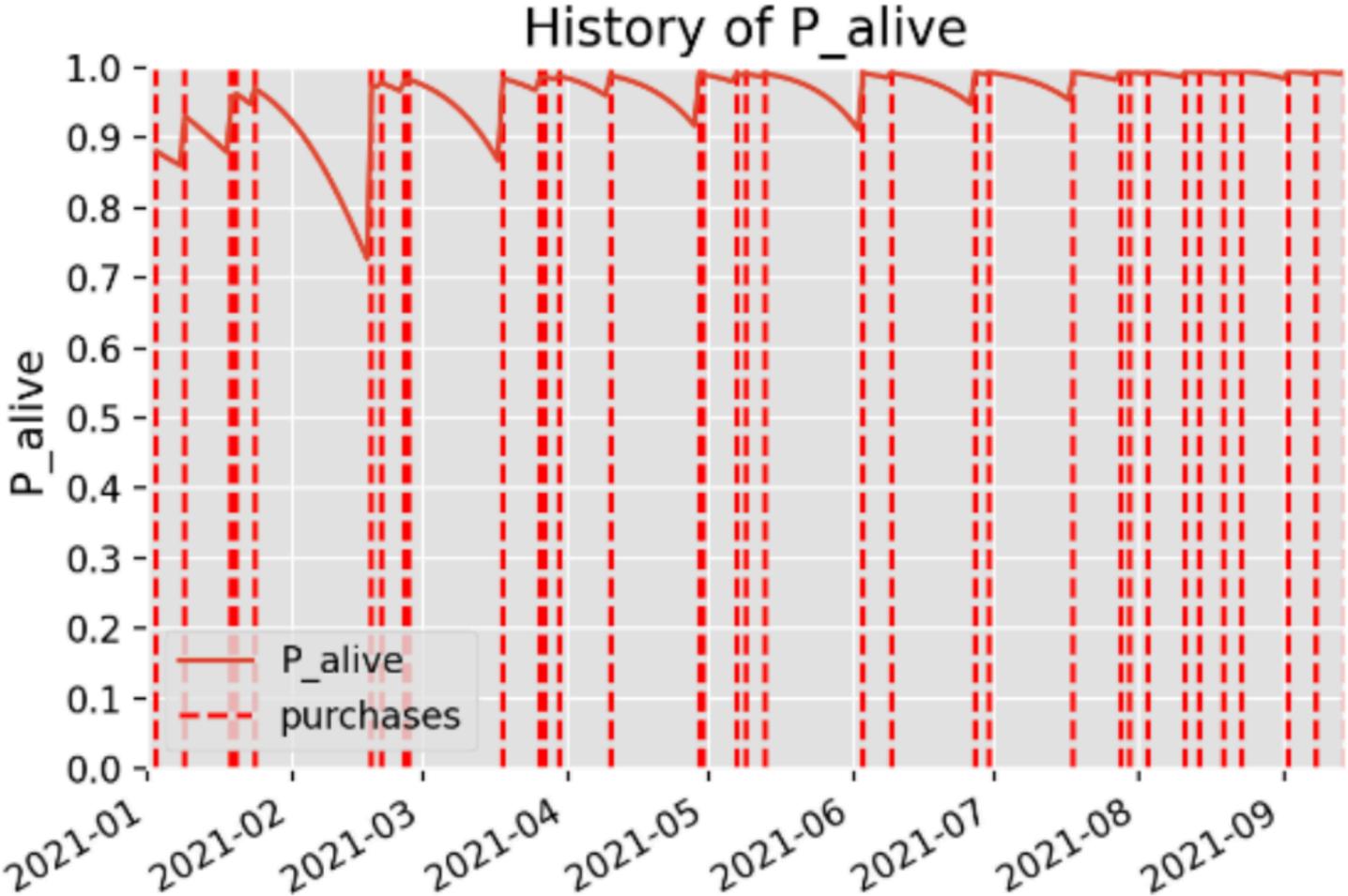
# Survival DL



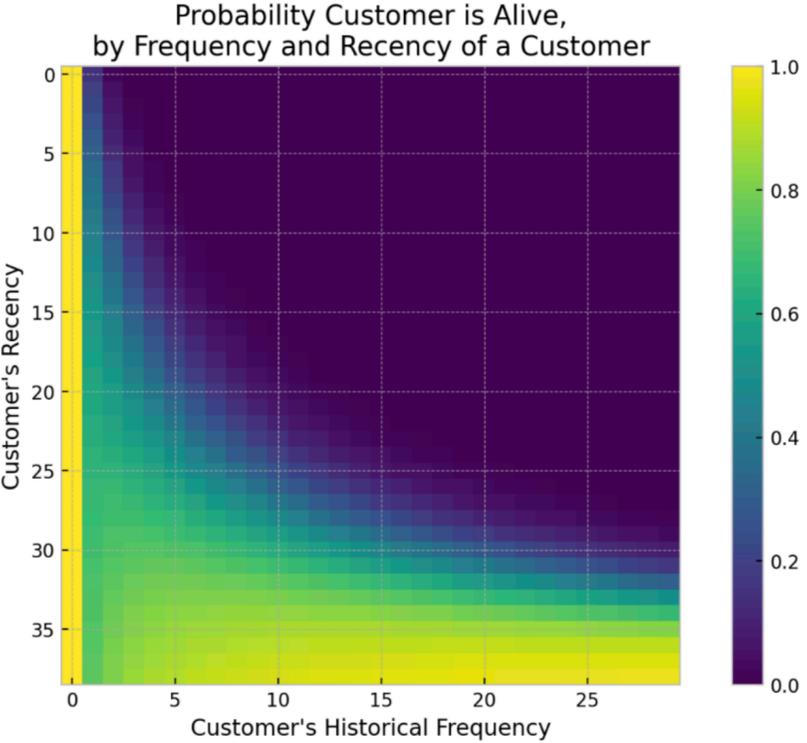
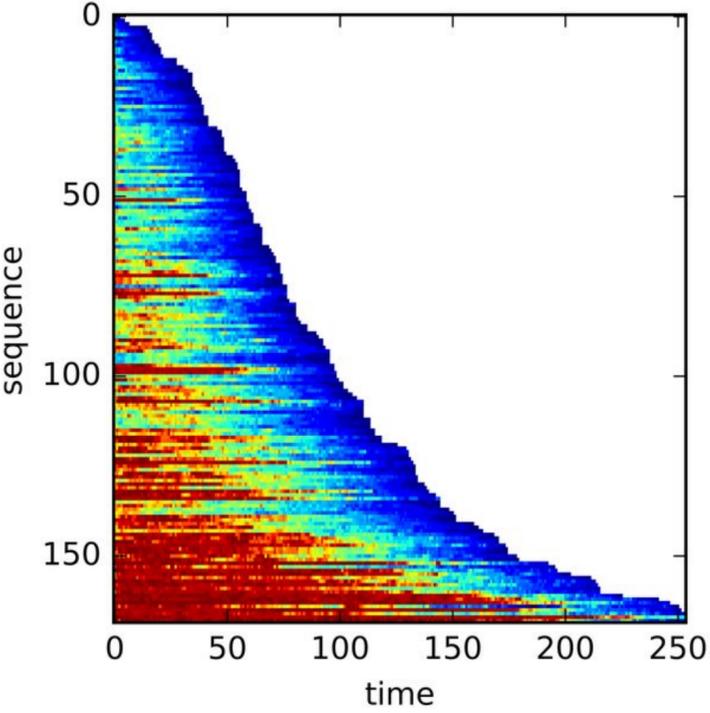
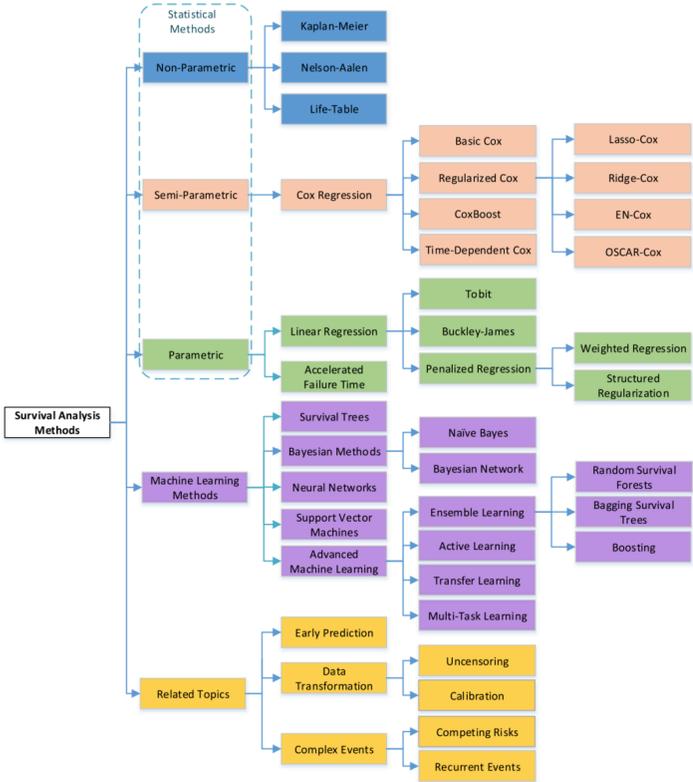
# Buy Till You Die (BTYD)



# Buy Till You Die (BTYD)



# Куда копать?



1. Включить "гуся"

2. Survival ML

3. Survival DL

4. Bayes++ BTYD

# Спасибо!

Трек соревнований на ODS:  
**Data Fusion Contest 2024**



Telegram:  
**Data Fusion Contest 2024 Chat**



# Q&A time

