

About the User Classification Problem Based on Analyzing the Odnoklassniki Friendship Graph

Alexey Zinoviev, PhD student, OmSU

Social Network

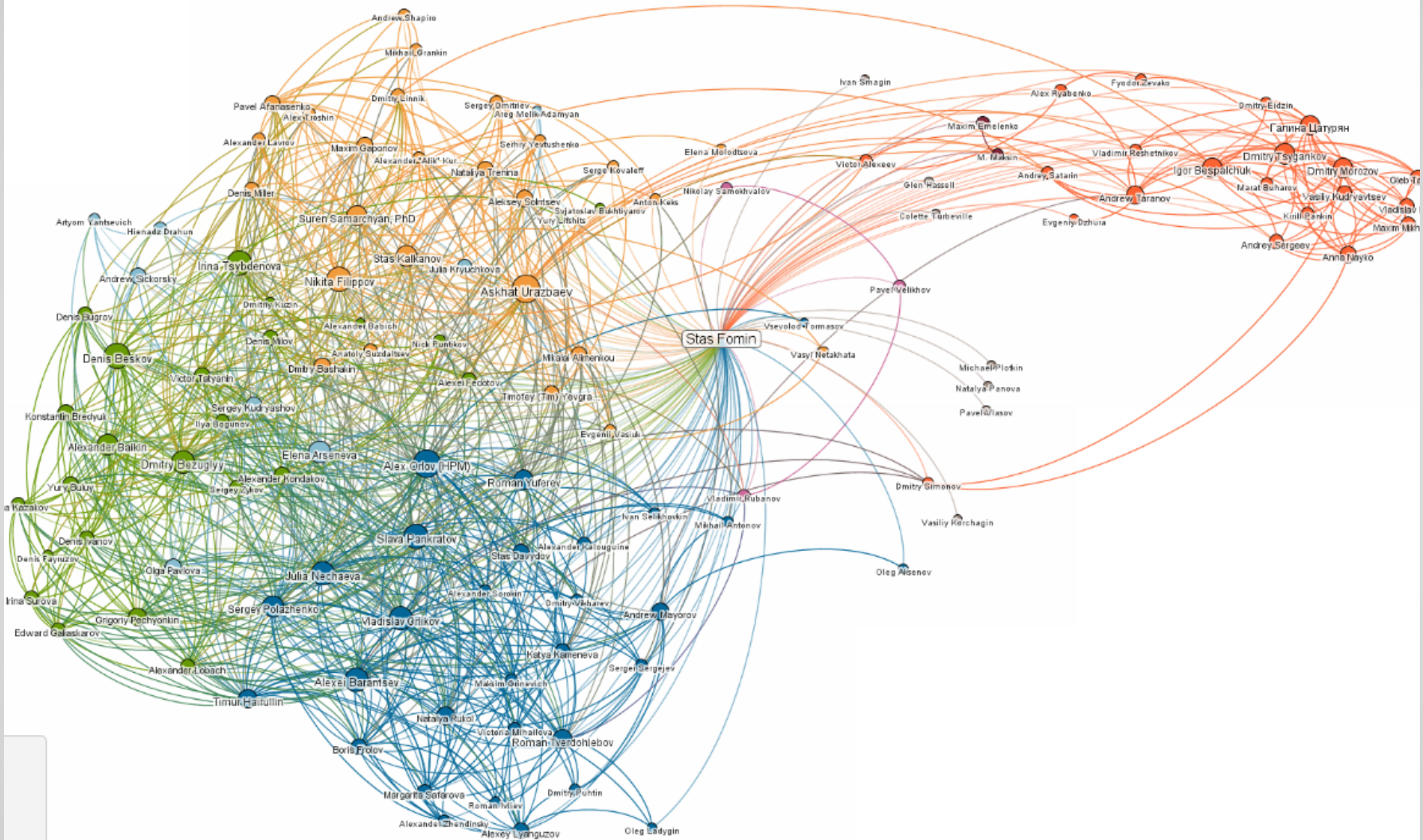


In common:

- 200 000 000 users
- 8 500 000 communities

Per day:

- 40 000 000 users
- 250 000 000 messages
- 8 000 000 posts
- 12 000 000 photos
- 7 000 000 new links (friendships)



The malicious activity

- Offense against the laws of ethics, morality, and articles
RF Criminal Code
- Creation of hidden subnetwork with spam accounts
- Hacking profiles of actual users
- Spam attack from hacked profiles
- Attraction of user's attention by user's page visiting

The benefits of social network

- Prevent the spread of profiles breaking "epidemic" and leakage of personal data
- Prevent spam before it arrives
- Reduce the number of complaints
- Reduce the burden on moderators
- Reduce the moderator staff

Dataset

- Graph ($\sim 9 * 10^6$, 39 Gb)
- Demography
- User likes
- History of logging ($\sim 3,2 * 10^8$, 12 Gb)
- Community posts
- Complaints about spam

Tools

- R 3.0.3 (for prototyping only)
- python 2.3 + scypi + numpy + pandas (data mining)
- Hadoop 2.6 (cluster infrastructure)
- Pig 14 (for user's features calculating)
- Giraph 1.1 (for graph-related features calculating)

The Problem

It should offer mathematical model makes prediction with high reliability to determine that the user is an attacker. It should be based on the number of friends, history of logging and analysis of other activities (type I error is not more than 1% and a type II error $< 10\%$) .

The model

Set of objects - social network users

Each object should be classified as User or Spamer.

Training set is produced on complaints of actual users.

Features

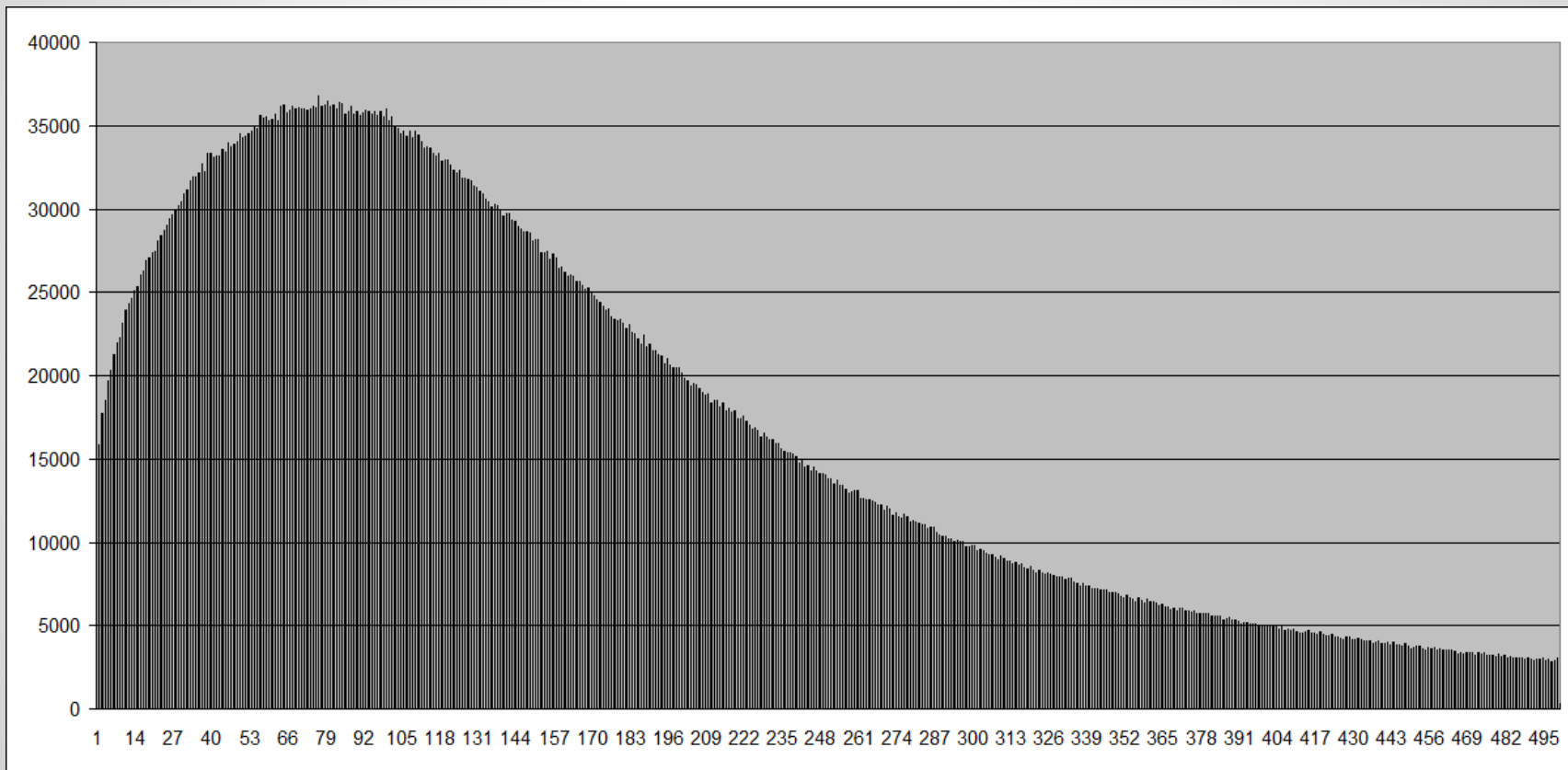
- Local feature: vertex degree
- Global feature: PageRank for each vertex
- Global-local feature: local clustering coefficient value (LCC)
- Number of successful logins
- Demography
- Geography

Training set

Features were calculated for 10000 users:

- age, is_male, is_female
- degree, lcc, page_rank, geo_lcc
- good_auth_per_week, bad_auth_per_week
- dist_from_Moscow, dist_from_borders

Vertex degree distribution



Computational experiment

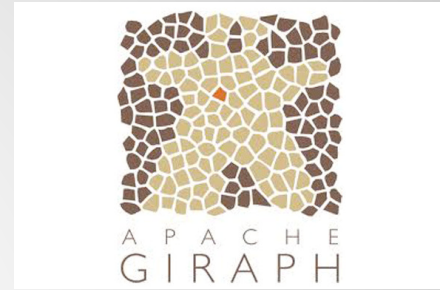
4 servers with 8 cores and 30 Gb RAM, in Google Compute Engine.

Hadoop Cluster + Pig for feature calculation.

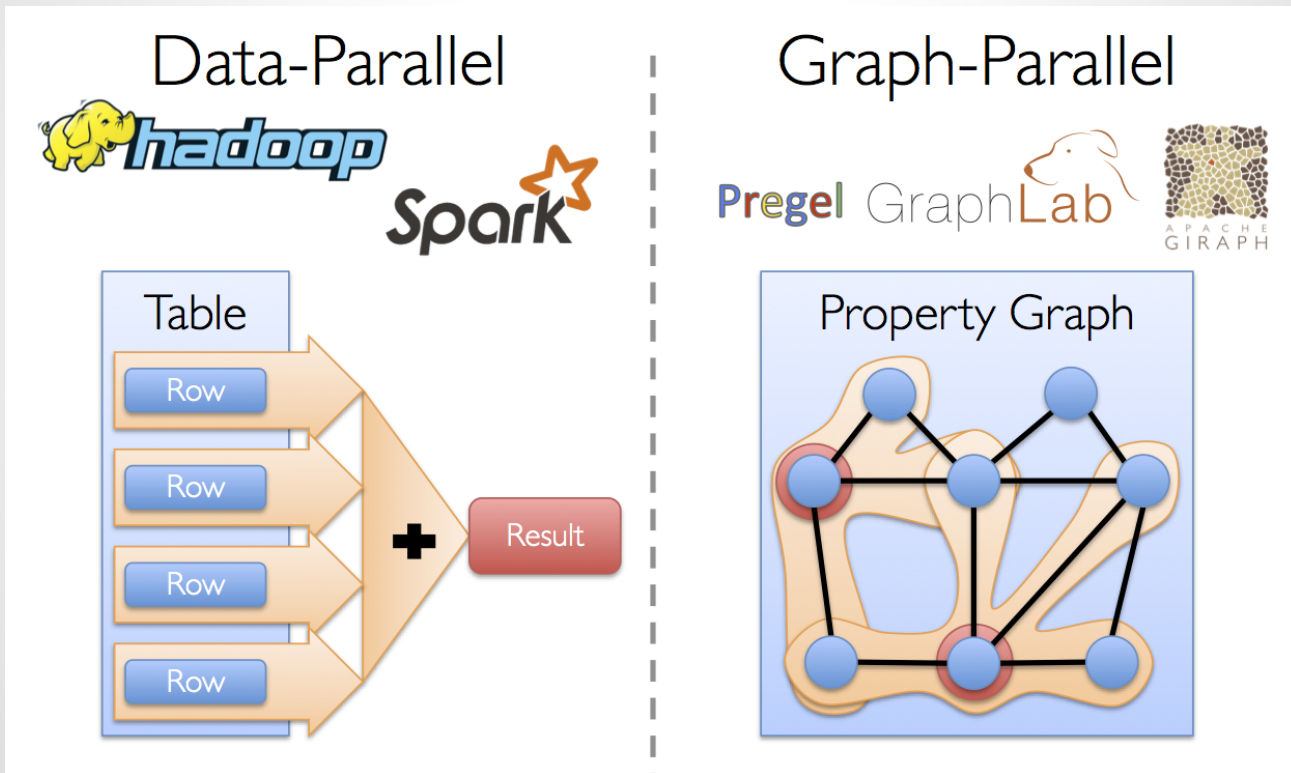
Giraph, above Hadoop cluster for calculating of PageRank and lcc.

Why Giraph?

- Open-source Pregel implementation
- Works on existing Hadoop infrastructure
- Calculations in memory
- Simple organized iterative calculations (it's important for PageRank)



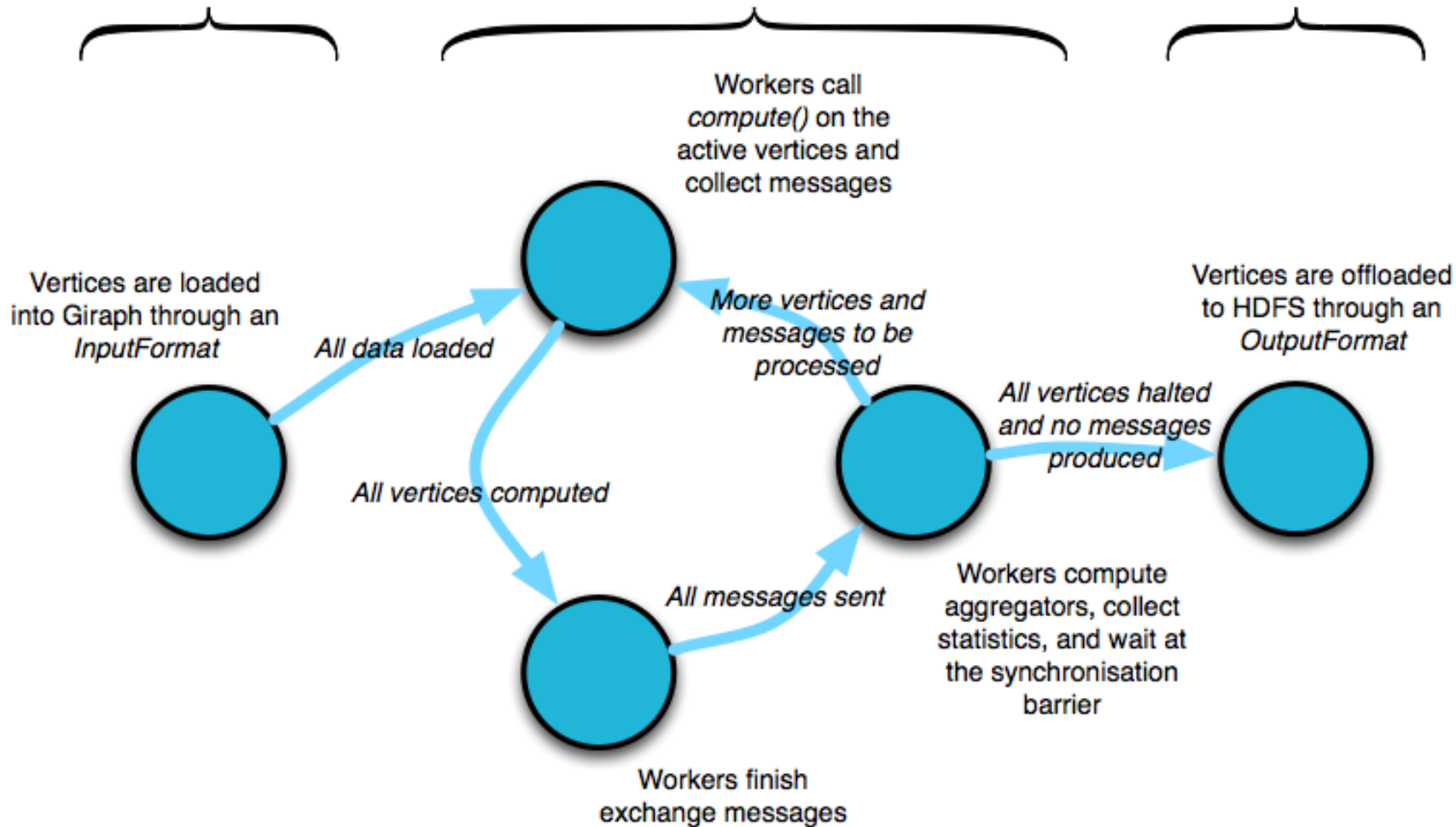
Думай вершинами, а не строками...



Loading phase

Compute phase

Offloading phase



Time of experiment

Iterative execution of PageRank, written in Pig was finished in 25 iterations, 123 minutes (~ 5 minutes per iteration)

Giraph implementation of PageRank cost 45 iterations and 25 minutes (~ 35 seconds per iteration) with running condition 1 worker per 1 core

Model

For model creation it used kNN, polynomial regression and decision trees(Random Forest, C4.5).

The best results had kNN ($n = 7$) and C4.5 with type I error 5% and 3%, type II error 12% and 19%, respectively.

Feature's importance

geo_lcc and **degree** are most important features, after followed in order of importance **lcc**, **dist_from_Moscow**, **good_auth_per_week** and **page_rank**.

But social-demography data provided by each user in his personal profile had a worst importance in decision trees and low importance for kNN.

In conclusion

- Calculation of graph features for big dataset is very difficult for MapReduce approach and needs in Pregel approach.
- Features derived from the analysis of the relationship structure are important in solving the problem of spam accounts searching.
- Hadoop + Pig + Giraph in Google Compute Engine - easy scalable infrastructure for implementing SNA models and algorithms.

Haec habui, quae dixi

