

Определение места проживания пользователей социальных сетей

на основе социального графа

Юля Трофимович, Илья Козлов, Денис Турдаков

1 декабря 2016 г.

ИСП РАН

Table of contents

1. Определение местоположения: о задаче

2. Что было сделано до нас

Первые попытки

Label propagation для определения местоположения

Machine learning для определения местоположения

3. Наш подход

Что такое graph node embeddings

Подход на основе graph node embeddings

Результаты, таблички, сравнения

Определение местоположения: о задаче

Определение местоположения: о задаче

Дано:

- социальный граф (ВКонтакте), 156 057 700 вершин
- информация о месте проживания указанная пользователями (известна для 63% пользователей)

Найти:

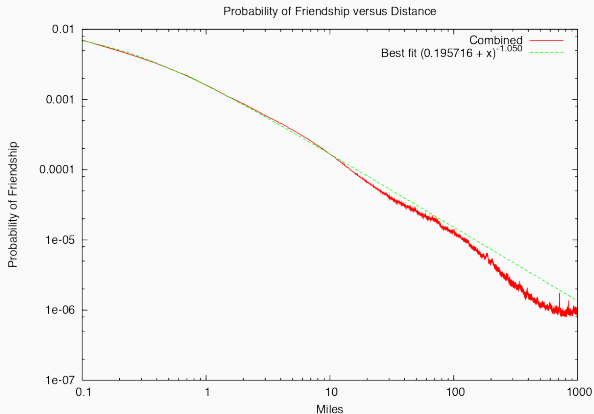
- основное место проживания пользователей, которые не указали его в профиле

Точность:

- субъект РФ/страна

Что было сделано до нас

Backstrom, Sun и Marlow (2010):



Probability of friendship as a function of distance

Backstrom, Sun и Marlow (2010):

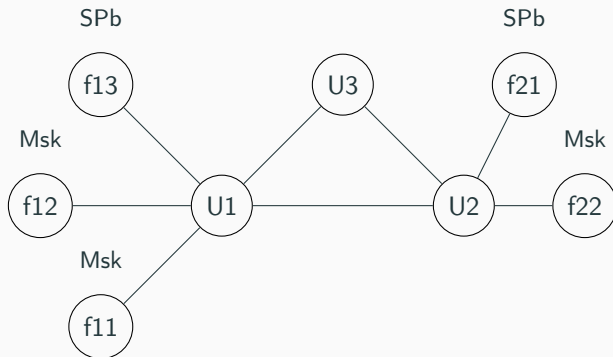
$$p(|l_u - l_v|) = 0.0019(|l_u - l_v| + 0.196)^{-1.05} \quad (1)$$

$$\prod_{(u,v) \in E} p(|l_u - l_v|) \prod_{(u,v) \notin E} 1 - p(|l_u - l_v|) \quad (2)$$

Jurgens (2013) “That’s What Friends Are For: Inferring Location in Online Social Media Platforms Based on Social Relationships.”:

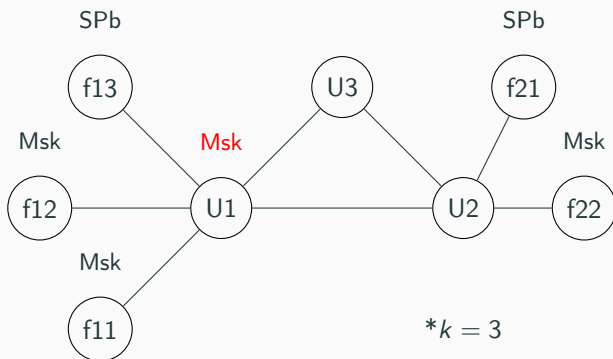
- Инициализация: присвоение меток вершинам
- Итеративно:
 - Каждая вершина отправляет свою метку соседним вершинам
 - Каждая вершина обрабатывает полученные метки, и, при необходимости, обновляет свою
- Окончание работы: если в ходе итерации ни одна вершина не обновила свою метку, работа алгоритма заканчивается

Label propagation для определения местоположения



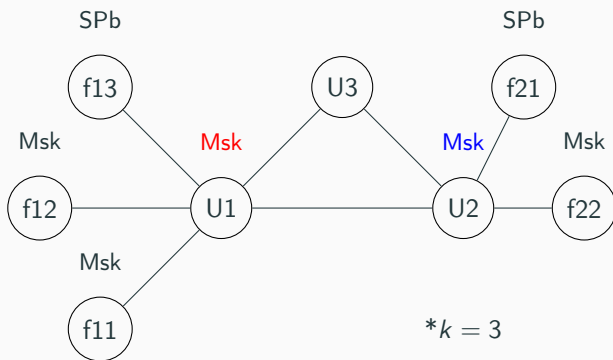
Пример социального графа

Label propagation для определения местоположения



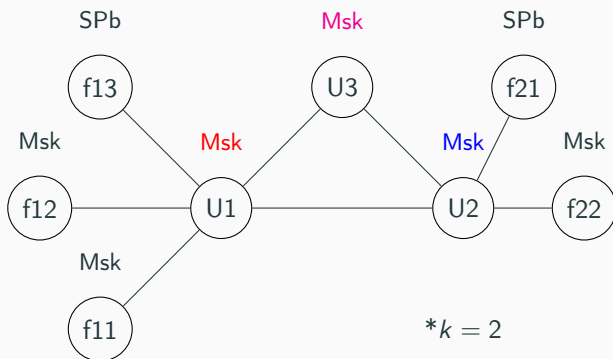
Социальный граф после первой итерации подхода на основе распространения меток

Label propagation для определения местоположения



Социальный граф после второй (и третьей тоже) итерации подхода на основе распространения меток

Label propagation для определения местоположения

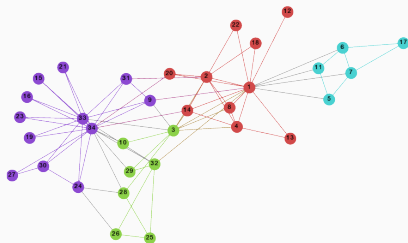


Социальный граф после четвертой итерации подхода на основе распространения меток

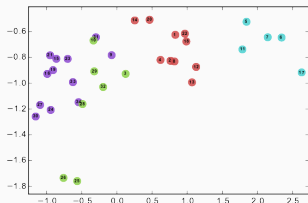
- SVM
 - Rout и др. (2013): City size, City population bins, Triads
- PGM
 - Chen, Liu и Zou (2016): Tie strength (Random Walk with Restart)
 - Jia и др. (2016): Influential friends (Sequential Random Walk with Restart)

Наш подход

Graph node embeddings



(a) Input: Karate Graph

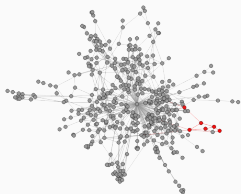


(b) Output: Representation

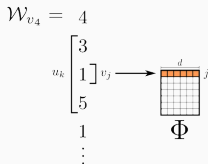
Вложение Zachary's Karate Club в плоскость

Graph node embeddings: DeepWalk

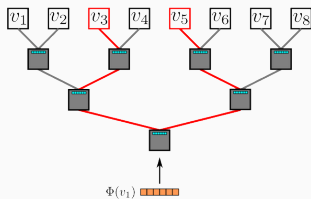
Perozzi, Al-Rfou и Skiena (2014): **DeepWalk**



(a) Random walk generation.



(b) Representation mapping.



(c) Hierarchical Softmax.

Overview of DeepWalk

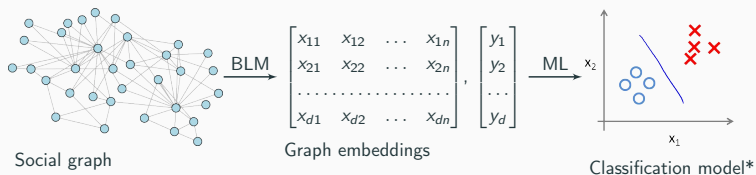
Ivanov и Bartunov (2015): **Bilinear Link Model**

$$\vec{In} \in \mathbb{R}^D, \vec{Out} \in \mathbb{R}^D \quad (3)$$

Вероятность ребра между вершинами $v \rightarrow u$

$$p(v|u) = \frac{\exp(\vec{In}_u^T \vec{Out}_v)}{\sum_{w \in V} \exp(\vec{In}_u^T \vec{Out}_w)} \quad (4)$$

Подход на основе graph node embeddings



* Мы использовали нейронную сеть - многослойный перцептрон

- Размер датасета: 99 055 808 пользователей
- Обучающая выборка: 15 003 815 ($\sim 15\%$) выбранных случайным образом пользователей
- Тестовые выборки:
 - (a) 84 051 993 ($\sim 85\%$) - пользователи не вошедшие в обучающую выборку
 - (b) 4 588 844 ($\sim 4.6\%$) - пользователи не вошедшие в обучающую выборку и оставившие хотя бы одно сообщение в топ-миллионе активных групп
- Количество классов: 277

(а) Пользователи вне зависимости от наличия сообщений

Метрика		Подход	
		LP	GE
Accuracy		0.463	0.516
Macro	Precision	0.522	0.308
	Recall	0.150	0.167
	F1-score	0.192	0.203
Weighted	Precision	0.760	0.729
	Recall	0.463	0.516
	F1-score	0.555	0.595

Обозначения: LP - подход на основе распространения меток, GE - подход на основе векторного представления

(b) Только пользователи, оставившие хотя бы одно сообщение

Метрика		Подход	
		LP	GE
Accuracy		0.786	0.839
Macro	Precision	0.374	0.302
	Recall	0.241	0.268
	F1-score	0.255	0.271
Weighted	Precision	0.819	0.823
	Recall	0.786	0.840
	F1-score	0.788	0.827

Обозначения: LP - подход на основе распространения меток, GE - подход на основе векторного представления

The End

