

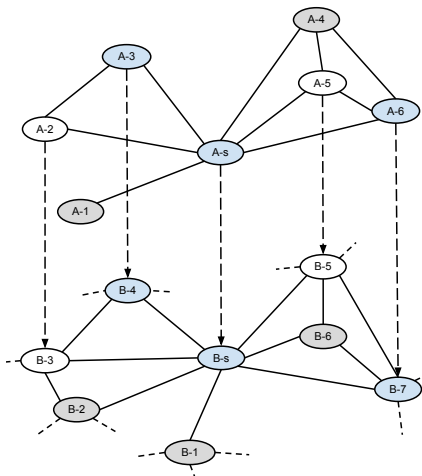
Joint Link-Attribute User Identity Resolution In Online Social Networks

Сергей Бартунов, Антон Коршунов

ИСП РАН

22 февраля 2012 г.

Постановка задачи



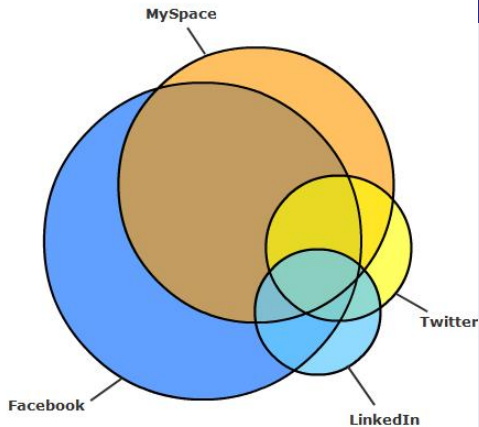
Данные

- Два социальных графа $\langle A, B \rangle$:
 - Поля профилей (имя, адрес, день рождения и т.д.)
 - Социальные связи (друзья, подписка, ...)

Задача

Найти как можно больше пар профилей $(v, u) \mid v \in A, u \in B$, принадлежащих одному человеку

SNS Usage Overlap



Source: Anderson Analytics 2009

Мотивация

- Многие используют социальные сети и имеют несколько аккаунтов
- Есть множество нишевых социальных сетей
- Социальная информация может сильно помочь во многих приложениях, нужен более полный социальный граф
- Объединение контактов

Профили в социальных сетях



Sergey Bartunov

@sbos

Беспечный спецзодок! Зайцы съедят меня, когда я стану травой
Moscow

Профили в социальных сетях

- Некоторые содержат множество данных (Facebook)
- Некоторые практически никакой (Twitter)



Sergey Bartunov •

Works at ISP RAS

Studied at Московский государственный университет имени Ломоносова

Upload a profile picture

Friends

All Friends (55)

Contact Information

To edit, click on highlighted profile field labels

Profile: Create username

Emails: sbos@sbos.in
sbos.net@gmail.com

Twitter: sbos

Skype: xsbosx

Phones: 8 9851090410
Add phone

Website: _____

Basic Information

To edit, click on profile field labels

Sex: Male

Birthday: November 12

Current City: Moscow, Russia

Hometown: Moscow, Russia

Family: _____

Relationship: In a relationship

Interested In: Women

Languages: Russian, English and Albanian

Political: БАТОЧКА

Attribute-based UIR

- Поля профилей сравниваются с помощью функций нечеткого сравнения строк
- Результаты взвешиваются, суммируются и сравниваются с пороговым значением

Недостатки

- Не всегда пользователи аккуратно заполняют поля профилей или держат их в актуальном состоянии
- Иногда люди придумывают ники, а не вводят реальные имена
- Одинаковые имена не всегда означают одного владельца
- Профили вообще не всегда доступны из-за приватности

Сравнение частично сопоставленных списков контактов

- Сначала профили сопоставляются по полям
- Затем в качестве дополнительной информации привлекается показатель близости сопоставленных друзей

Похожесть списков контактов

- $J(L_v, L_u) = \frac{|L_v \cap L_u|}{|L_v \cup L_u|}$
- $\cos(L_v, L_u) = \frac{|L_v \cap L_u|}{\sqrt{|L_v| |L_u|}}$
- $\text{dice}(L_v, L_u) = \frac{2|L_v \cap L_u|}{|L_v| + |L_u|}$

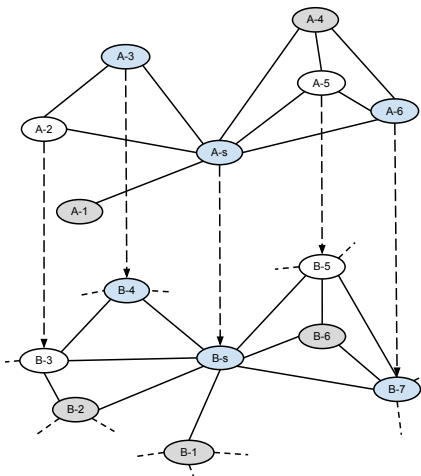
Сравнение частично сопоставленных списков контактов

- Сначала профили сопоставляются по полям
- Затем в качестве дополнительной информации привлекается показатель близости сопоставленных друзей

Недостатки

- Техника опирается на ненадежные поля профилей
- Такой показатель близости очевидно нестабилен

Обозначения



Локальная перспектива

- Центральный профиль
- Все профили, с ним связанные
- Списки контактов этих профилей

Задача

Найти оптимальное отображение

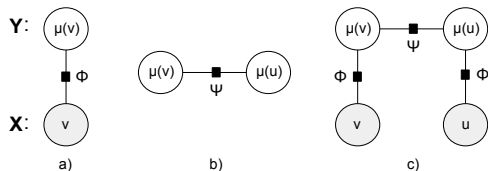
$$\mu : V(A) \rightarrow V(B) \cup N$$

N - нейтральная проекция

Основные положения

- Необходимо использовать, как информацию из профилей, так и социальные связи
- Задачи выбора проекций $\mu(v)$ и $\mu(u)$ для связанных вершин v и u взаимосвязаны
- Если v и u связаны в графе A , то их проекции в графе B должны быть близки друг к другу

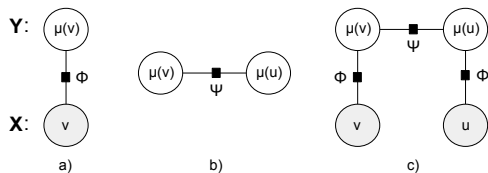
Joint Link-Attribute Model



Вероятностная модель

- Для каждой вершины v в графе A :
 - Наблюдаемая переменная x_v - профиль v
 - Скрытая переменная $y_v = \mu(v)$ - проекция профиля v в графе B
- Переменные x_v и y_v связаны фактором Φ
- Переменные y_v и y_u связаны фактором $\Psi \Leftrightarrow (v, u) \in A$

Joint Link-Attribute Model

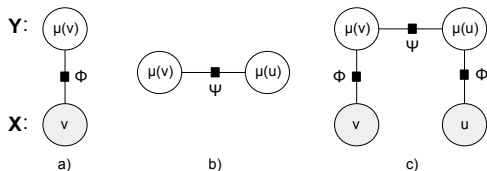


Энергия модели

$$E(\mathbf{Y}|\mathbf{X}) = \sum_{v \in V(A)} \Phi(y_v | x_v) + \sum_{(v, u) \in E(A)} \Psi(y_v, y_u)$$

- $\Phi(y_v | x_v) \sim \text{profile-distance}(v, \mu(v))$ - (не)похожесть профиля на свою проекцию
- $\Psi(y_v, y_u) \sim \text{network-distance}(\mu(v), \mu(u))$ - расстояние между проекциями в графе B

Joint Link-Attribute Model

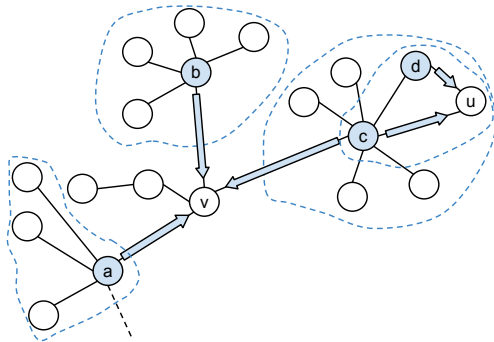


Joint nature

$$E(\mathbf{Y}|\mathbf{X}) = \sum_{v \in V(A)} \Phi(y_v | x_v) + \sum_{(v, u) \in E(A)} \Psi(y_v, y_u)$$

- a) Если не рассматривать социальные связи, то $\Psi \equiv 0$, и модель вырождается до тривиальной системы
- b) Если недоступна информация о полях профилей, то $\Phi \equiv 0$, и решается задача деанонимизации
- c) Вся информация доступна

Заранее известные проекции



Заранее известные проекции

- Для некоторых вершин проекции могут быть известны заранее
- Некоторые можно считать таковыми, если $\text{profile-distance}(v, \text{pr}(v)) \leq \Delta$

Похожесть профилей в Twitter и Facebook

Схема сравнения

| Facebook | Twitter | Функция сравнения |
|----------|-------------|---------------------|
| Name | Name | VMN |
| | Screen name | Screen Name measure |
| Website | URL | URL measure |

Функции сравнения

- Screen Name проверяет не совпалили ли явно ник и имя
- URL measure проверяет не указал ли явно пользователь второй профиль
- VMN позаимствована из Vosecky et. al., *User identification across multiple social networks*, 2009.

Похожесть профилей в Twitter и Facebook

Вектор похожести

- К каждой паре полей применяется соответствующая функция сравнения
- В результате получается *вектор похожести* $V(v, \mu(v))$

Функция (не)похожести

- Вектор $V(v, \mu(v))$ можно использовать как набор признаков для обучения бинарного классификатора
- $\text{profile-distance}(v, \mu(v)) = P(\text{разные люди} | V(v, \mu(v)))$

Похожесть профилей в Twitter и Facebook

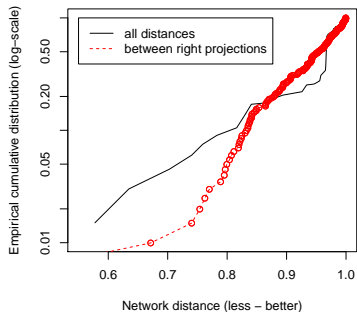
Сравнение классификаторов

| алгоритм | полнота | точность | F_1 |
|----------------------|--------------|--------------|-------------|
| Naive Bayes | 0.862 | 0.308 | 0.453 |
| C4.5 | 0.569 | 0.86 | 0.685 |
| C4.5 с MultiBoosting | 0.669 | 0.879 | 0.76 |

Выводы

- В целом, это работает
- Ни один классификатор не смог идеально „объяснить” принадлежность профиля

Очистка результатов



Плохие результаты

- Плохая связность (см. рисунок)
- Мало вершин с заранее известными проекциями
- ...

Вывод: результаты надо очищать

- Меньше - лучше (ближе)
- Нет разумного порогового значения

Очистка результатов - тривиальное решение

Взаимное проецирование

- Получить μ из A в B
- Получить ν из B в A
- Если $v \neq \nu(\mu(v))$, то $\mu(v) \leftarrow \mathbf{N}$

Свойства

- Просто и интуитивно
- Слишком грубая техника очистки результатов - не учитывает *причину* ошибки
- Требуется в ~ 2 раза больше времени

Очистка результатов - классификатор

Признаки

- 1 profile-distance(v , $pr(v)$)
- 2 Средняя графовая близость к проекциям смежных вершин
- 3 Доля заранее известных проекций среди смежных вершин
- 4 Взаимо-согласованность смежных вершин с заранее известными проекциями:

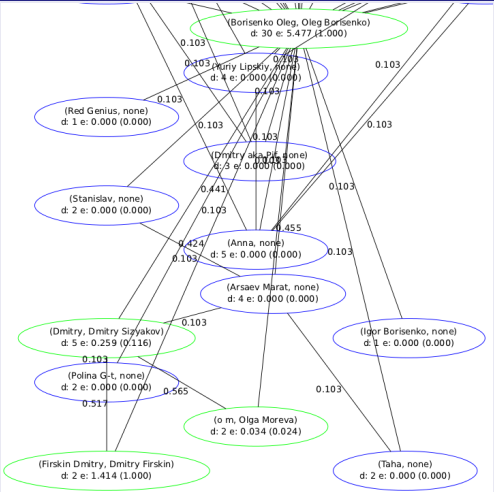
$$\frac{1}{n} \cdot \sum_v \frac{1}{n-1} \sum_{u \neq v} \text{network-distance}(pr(v), pr(u))$$

Сравнение классификаторов

| алгоритм | полнота | точность | F_1 |
|------------------------|--------------|--------------|--------------|
| Naive Bayes | 0.762 | 0.256 | 0.383 |
| Support Vector Machine | 0.662 | 0.935 | 0.775 |
| C4.5 | 0.715 | 0.939 | 0.812 |
| C4.5 с MultiBoosting | 0.844 | 0.902 | 0.872 |

Результаты идентификации

Пример



Экспериментальные данные

Статистика

| | Twitter | Facebook |
|------------------------------|---------|----------|
| Основная выборка | | |
| # центральных пользователей | | 16 |
| # профилей | 398 | 977 |
| # связей | 1 728 | 10 256 |
| # сопоставленных профилей | | 141 |
| # заранее известных проекций | | 71 |
| Дополнительная выборка | | |
| # центральных пользователей | | 17 |
| # профилей | 1 499 | 7 425 |
| # связей | 15 943 | 172 219 |
| # сопоставленных профилей | | 161 |

Максимальное парасочетание

- Каждому профилю из A нужно сопоставить не более одного профиля из B
- Сопоставленные профили должны быть как можно более похожи по полям профилей
- Значение функции похожести должно быть не ниже некоторого порога
- Порог максимизирует точность

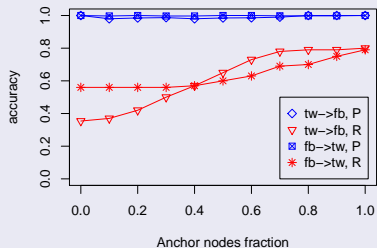
Базовые алгоритмы

- 1 Взвешенная сумма значений вектора $V(v, \mu(v))$
- 2 $1 - \text{profile-distance}(v, \mu(v))$

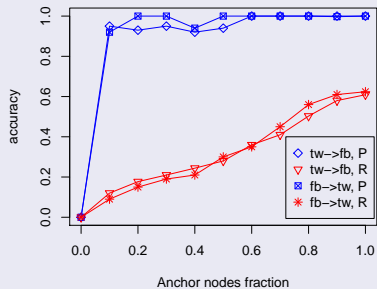
Результаты на основной выборке

| алгоритм | полн. | точн. | F_1 |
|-------------------------------------|------------|------------|-------------|
| безразличные к направлению проекции | | | |
| Базовый 1 (взвешенная сумма) | 0.45 | 0.94 | 0.61 |
| Базовый 2 (вероятностная похожесть) | 0.51 | 1.0 | 0.69 |
| JLA, взаимн. проекц., аноним. | 0.6 | 1.0 | 0.76 |
| JLA, взаимн. проекц. | 0.66 | 0.99 | 0.79 |
| Twitter → Facebook | | | |
| JLA, анонимн. ($\Phi \equiv 0$) | 0.62 | 1.0 | 0.77 |
| JLA | 0.79 | 1.0 | 0.89 |
| Facebook → Twitter | | | |
| JLA, анонимн. ($\Phi \equiv 0$) | 0.61 | 1.0 | 0.76 |
| JLA | 0.8 | 1.0 | 0.89 |

Идентификация



Деанонимизация



Повторная идентификация

Мотивация

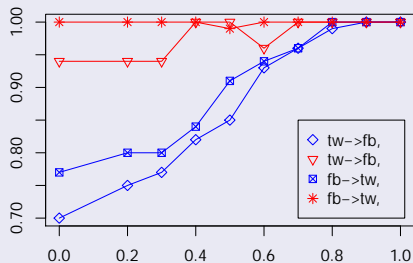
- Основная выборка мала
- Собрать хорошую выборку тяжело без помощи владельцев аккаунтов
- Нужно протестировать алгоритм автоматически на неразмеченной выборке

Постановка задачи

- При помощи второго базового алгоритма сопоставляются профили
- Фиксируется некоторая часть из них
- У всех возможных проекций для этой части профилей убирается вся информация, кроме связей
- Нужно найти проекции для этих профилей по связям

Повторная идентификация

Идентификация



Вывод

80% известных проекций
достаточно чтобы определить
оставшиеся 20%