



МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ ИМЕНИ М.В. ЛОМОНОСОВА
ФАКУЛЬТЕТ ВЫЧИСЛИТЕЛЬНОЙ МАТЕМАТИКИ И КИБЕРНЕТИКИ
КАФЕДРА МАТЕМАТИЧЕСКИХ МЕТОДОВ ПРОГНОЗИРОВАНИЯ

НИКИШИН ЕВГЕНИЙ СЕРГЕЕВИЧ

Методы выделения сообществ в социальных графах

КУРСОВАЯ РАБОТА

Научный руководитель:
д.ф-м.н., профессор
А.Г. Дьяконов

Москва, 2016

version 0.15

Содержание

1 Введение (неполное)	1
1.1 Модулярность	1
2 Разбиение на непересекающиеся сообщества	1
2.1 Edge Betweenness	1
2.2 Label Propagation	2
2.3 FastGreedy	3
2.4 WalkTrap	3
2.5 Infomap	4
2.6 Leading Eigenvector	5
2.7 MultiLevel	7
3 Разбиение на пересекающиеся сообщества	9
3.1 k-Clique Perlocation	9
3.2 BigCLAM	9
3.3 DEMON	10
3.4 CONGO	11
4 Эксперименты	12
4.1 Данные	12
4.2 Метрики качества	12
4.3 Результаты	14

1 Введение (неполное)

1.1 Модулярность

$$Q = \frac{1}{2m} \sum_{i,j} \left(A_{ij} - \frac{d_i d_j}{2m} \right) \delta(C_i, C_j)$$

2 Разбиение на непересекающиеся сообщества

2.1 Edge Betweenness

Для каждой пары вершин связного графа можно вычислить кратчайший путь, их соединяющий. Будем считать, что каждый такой путь имеет вес, равный $1/N$, где N — число возможных кратчайших путей между выбранной парой вершин. Если такие веса посчитать для всех пар вершин, то каждому ребру можно поставить в соответствие значение Edge betweenness — сумму весов путей, прошедших через это ребро.

Для ясности приведём следующую иллюстрацию:

В данном графе хочется выделить два сообщества: с вершинами 1-5 и 6-10. Граница же будет проходить через ребро, имеющее максимальный вес, 25. На этой идее и

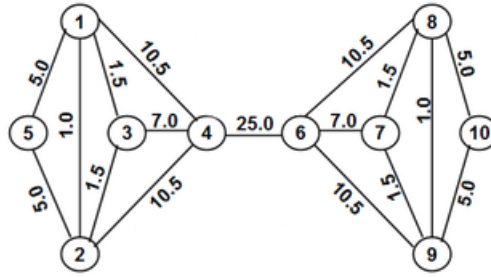


Рис. 1: Граф, для ребёр которого посчитаны значения Edge betweenness

основывается алгоритм: поэтапно удаляем ребра с наибольшим весом, а оставшиеся компоненты связности объявляем сообществами.

Собственно, сам алгоритм:

1. Инициализировать веса
2. Удалить ребро с наибольшим весом
3. Пересчитать веса для ребёр
4. Сообществами считаются все компоненты связности
5. Посчитать функционал модулярности (о нём будет сказано ранее)
6. Повторять с шага 2-6, пока есть рёбра

На каждой итерации процесса получается некое разбиение вершин. Последовательность таких разбиений, имеющая вид дерева, в листьях которого находятся сообщества с одной вершиной, а в корне — большое сообщество, содержащее все вершины, называется дендрограммой. Результатом работы алгоритма является ярус дендрограммы (т.е. разбиение), имеющий максимальную модулярность.

Из необходимости каждый раз пересчитывать веса следует главный минус: вычислительная сложность в худшем случае составляет $O(m^2n)$, где m — количество ребёр, n — количество вершин. Эксперименты показывают, что пересчитывать обычно приходится только веса для рёбер, которые были в одной компоненте связности, что несколько уменьшает сложность, однако зачастую этого оказывается недостаточно.

2.2 Label Propagation

Допустим, что большинство соседей какой-либо вершины принадлежат одному сообществу. Тогда, с высокой вероятностью, ему также будет принадлежать выбранная вершина. На этом предположении и строится алгоритм Label propagation: каждая вершина в графе определяется в то сообщество, которому принадлежит большинство его соседей. Если же таких сообществ несколько, то выбирается случайно одно из них. Пример:

В начальный момент времени всем вершинам ставится в соответствие отдельное сообщество. Затем происходят перераспределения сообществ. Из-за случайности

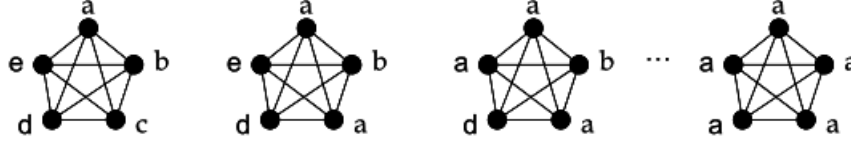


Рис. 2: Демонстрация работы алгоритма для полного графа

важно на каждой итерации изменять порядок обхода вершин. Алгоритм заканчивает работу, когда нечего изменять: все вершины относятся к тем сообществам, что и большинство их соседей. Авторы также советуют запускать несколько раз алгоритм и выбирать наилучшее из результирующих разбиений, либо пересекать их. Главное достоинство данного алгоритма, в противовес предыдущему, — почти линейная сложность. Однако на зашумленных графах зачастую происходит объединение всех вершин в одно сообщество.

2.3 FastGreedy

Алгоритм заключается в жадной оптимизации модулярности. Как и в прошлом методе, в каждой вершине графа инициализируется отдельное сообщество, а затем объединяются пары сообществ, приводящие к максимальному увеличению модулярности. При этом объединяются только инцидентные пары вершин, так как, в противном случае, модулярность не может увеличиться [во введении необходимо будет объяснить смысл модулярности, чтобы этот факт не вызывал вопросов].

Результатом работы алгоритма будет ярус дендрограммы, на котором модулярность максимальна.

Метод является вычислительно нетрудоёмким ($O(m \log n)$), легко применим к большим графам и, несмотря на жадность, зачастую неплохо справляется с задачей.

2.4 WalkTrap

Допустим, на вершинах графа задана такая метрика, что между двумя вершинами из разных сообществ расстояние велико, а из одного — мало. Тогда выделение сообществ можно рассматривать как задачу кластеризации вершин. Попытаемся ввести такую метрику, используя случайные блуждания. Объект может переместиться из вершины i в вершину j с вероятностью $P_{ij} = \frac{A_{ij}}{d_i}$, где A — матрица смежности, d_i — степень i . То есть на каждом шаге равновероятно выбирается "сосед" вершины i . Таким образом определяется матрица переходов P случайного блуждания. Она примечательна тем, что её степени являются вероятностями перехода из одной вершины в другую за соответствующее число шагов: вероятность перехода из i в j за t шагов равна $(P^t)_{ij}$. Также следует отметить, что $P = D^{-1}A$, где D — матрица со степенями вершин на диагонали. Используя этот аппарат можно ввести желаемую метрику на вершинах:

$$r_{ij} = \sqrt{\sum_{k=1}^n \frac{(P_{ik}^t - P_{jk}^t)^2}{d(k)}} = \|D^{-\frac{1}{2}} P_{i\bullet}^t - D^{-\frac{1}{2}} P_{j\bullet}^t\|,$$

где $P_{i\bullet}^t$ — вектор из вероятностей перехода за t шагов из вершины i во все другие. Вообще говоря, метрика зависит от t , авторы советуют брать $3 \leq t \leq 8$.

Естественным образом расстояние между вершинами обобщается на расстояние между сообществами:

$$r_{1C_2} = \|D^{-\frac{1}{2}}P_{C_1\bullet}^t - D^{-\frac{1}{2}}P_{C_2\bullet}^t\| = \sqrt{\sum_{k=1}^n \frac{(P_{C_1k}^t - P_{C_2k}^t)^2}{d(k)}},$$

где

$$P_{Cj}^t = \frac{1}{|C|} \sum_{i \in C} P_{ij}^t$$

Теперь, когда задана метрика, можно попытаться выделить кластеры в графе. Начальное разбиение — по одной вершине в каждом кластере $\mathcal{P}_1 = \{\{v\}, v \in V\}$. Также для всех пар инцидентных вершин считается расстояние. Далее для каждого k :

1. Выбрать C_1 и C_2 из \mathcal{P}_k согласно некоторому метрическому критерию.
2. Объединить два сообщества в новое $C_3 = C_1 \cup C_2$ и обновить разбиение $\mathcal{P}_{k+1} = (\mathcal{P}_k \setminus \{C_1, C_2\}) \cup C_3$.
3. Обновить расстояния между инцидентными сообществами.

После $n - 1$ шага получается дендрограмма разбиений, а $\mathcal{P}_n = \{V\}$. Таким образом, остался неясным только критерий выбора пар сообществ на шаге 1. Будем выбирать пару сообществ, минимизирующих приращение среднего квадратов расстояний между каждой вершиной и их сообществом при объединении сообществ. Т.е.

$$\Delta\sigma(C_1, C_2) = \frac{1}{n} \left(\sum_{i \in C_3} r_{iC_3}^2 - \sum_{i \in C_1} r_{iC_1}^2 - \sum_{i \in C_2} r_{iC_2}^2 \right) \rightarrow \min_{C_1, C_2}$$

Теперь осталось только получить результат, выбрав разбиение, на котором достигается максимума модулярность.

2.5 Infomap

В данном методе снова применяется подход, основанный на случайных блужданиях. У каждой вершины есть некоторая вероятность её посещения. С помощью кодов Хаффмана, в соответствии с этими вероятностями, можно закодировать путь блуждателя. Эта последовательность будет иметь некоторую длину. Однако, если использовать иерархическое кодирование (т.е. кодируем сообщество, затем кодируем вершины, попавшие в это сообщество. Коды вершин в разных группах могут совпадать), то можно сократить длину получившейся последовательности.

Как происходит иерархическое кодирование: при входе в сообщество записывается его уникальный код, затем записывается код вершины, в которую попали. Далее при переходах внутри сообщества пишутся только коды вершин. При выходе из сообщества пишется уникальный для него код выхода.

На этом и основывается метод Infomar: жадным способом минимизируется длина кода прогулки блуждателя.

Иллюстрация работы алгоритма (картинка уехала, но сейчас нет смысла править):

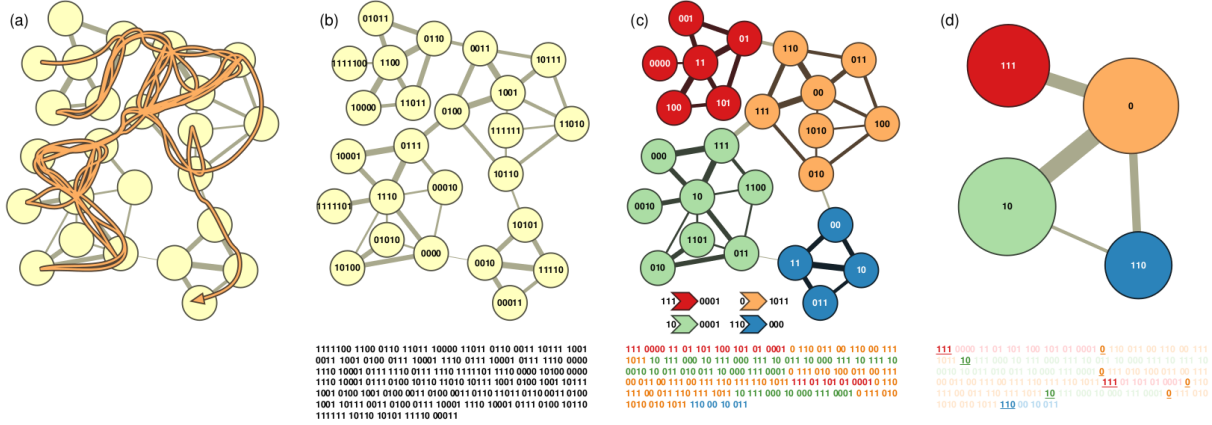


Рис. 3: На левом рисунке показан путь случайного блуждателя. На второй части изображены вершины с кодами Хаффмана, а ниже закодирован путь, изображенный на левом изображении. Далее показано кодирование с помощью иерархического метода. Ниже записаны коды сообществ и коды выхода из них. В самом низу закодирован путь. Длина кода уменьшилась. На последнем рисунке показаны сообщества и их коды

2.6 Leading Eigenvector

Для начала небольшой экскурс по спектральным методам выделения сообществ. Допустим, для простоты, что всего в графе 2 группы. Тогда предлагается, согласно неформальному определению сообществ, что количество ребёр между этими группами, также называемое *cut size*,

$$R = \frac{1}{2} \sum_{\substack{i,j \text{ в} \\ \text{разных} \\ \text{группах}}} A_{ij},$$

должно быть мало.

Чтобы получить более удобное представление вводится вектор индексов \mathbf{s}

$$s_i = \begin{cases} +1, & \text{если вершина } i \text{ принадлежит сообществу 1} \\ -1, & \text{если вершина } i \text{ принадлежит сообществу 2} \end{cases}$$

с n элементами. Тогда

$$\frac{1}{2}(1 - s_i s_j) = \begin{cases} 1, & \text{если вершины } i \text{ и } j \text{ принадлежат разным группам} \\ 0, & \text{если вершины } i \text{ и } j \text{ принадлежат одинаковым группам} \end{cases}$$

и величину *cut size* можно переписать в виде

$$R = \frac{1}{4} \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq n}} (1 - s_i s_j) A_{ij}$$

Используем следующую цепочку преобразований

$$\sum_{ij} A_{ij} = \sum_i d_i = \sum_i s_i^2 d_i = \sum_{ij} s_i s_j d_i \delta_{ij}$$

и перепишем *cut size* следующим образом:

$$R = \frac{1}{4} \sum_{ij} s_i s_j (d_i \delta_{ij} - A_{ij}) = \frac{1}{4} \mathbf{s}^T \mathbf{L} \mathbf{s}$$

где \mathbf{L} — матрица Лапласа с элементами

$$L_{ij} = \begin{cases} d_i, & \text{если } i = j \\ -1, & \text{если } i \neq j \text{ и между } i \text{ и } j \text{ есть ребро} \\ 0, & \text{иначе} \end{cases}$$

Далее надо отметить несколько замечательных свойств матрицы Лапласа:

1. Матрица симметричная, а, значит, собственные векторы образуют ортонормированный базис
2. Все собственные значения матрицы неотрицательны
3. Сумма по любой строке или по любому столбцу равна 0
4. Из свойств 2 и 3 следует, что всегда будет нулевое собственное значение и соответствующий ему собственный вектор $(1, 1, 1, \dots)/\sqrt{n}$

Можно пойти дальше и ещё упростить запись *cut size*: если разложить $\mathbf{s} = \sum_{i=1}^n a_i \mathbf{v}_i$, где $a_i = \mathbf{v}_i^T \mathbf{s}$, то

$$R = \sum_i a_i \mathbf{v}_i^T \mathbf{L} \sum_j a_j \mathbf{v}_j = \sum_{ij} a_i a_j \lambda_j \delta_{ij} = \sum_i a_i^2 \lambda_i$$

Таким образом, минимизацию R можно рассматривать как выбор a_i^2 , минимизирующих сумму. Как было отмечено, всегда существует собственный вектор из единиц. Если положить $\mathbf{s} = (1, 1, 1, \dots)$, то R становится равным нулю, что соответствует объединению всех вершин в одно сообщество. Такой тривиальный случай нас не интересует, поэтому рассматривается собственный вектор, соответствующий второму минимальному собственному значению. То есть мы будем подбирать вектор \mathbf{s} наиболее близким к $\mathbf{v}^{(2)}$. Учитывая ограничение, что значения \mathbf{s} могут быть только ± 1 , искомое \mathbf{s} принимает вид

$$s_i = \begin{cases} +1, & v_i^{(2)} \geq 0 \\ -1, & v_i^{(2)} < 0 \end{cases}$$

Это и есть спектральный метод в простейшем виде.

Однако авторы отмечают, что хорошее разделение — не совсем то, через которое проходит наименьшее число вершин. Поэтому они предлагают разделять те места, где количество рёбер меньше, чем ожидалось, или, наоборот, объединять те вершины, у которых количество рёбер больше, чем ожидалось.

Q = (количество вершин внутри сообщества) — (ожидаемое количество вершин)

$$= \frac{1}{2m} \sum_{i,j} \left(A_{ij} - \frac{d_i d_j}{2m} \right) \delta(C_i, C_j) \rightarrow \max$$

Используя $\delta(C_i, C_j) = \frac{1}{2}(s_i s_j + 1)$ перепишем

$$Q = \frac{1}{4m} \sum_{i,j} \left(A_{ij} - \frac{d_i d_j}{2m} \right) (s_i s_j + 1) = \frac{1}{4m} \mathbf{s}^T \mathbf{B} \mathbf{s},$$

где \mathbf{B} , называемая матрицей модулярности, во многих смыслах похожа на матрицу Лапласа

$$B_{ij} = A_{ij} - \frac{d_i d_j}{2m}$$

И именно настраивая вектор \mathbf{s} на собственный вектор, соответствующий максимальному собственному значению, получается метод Leading Eigenvector.

$$s_i = \begin{cases} +1, & u_i^{(1)} \geq 0 \\ -1, & u_i^{(1)} < 0 \end{cases}$$

Напомним, что будет относить к сообществу 1 те вершины, у которых соответствующее значение вектора \mathbf{s} равно плюс единице, и к сообществу 2 иначе. Подобным образом граф разбивается на сообщества, пока увеличивается значение модулярности.

2.7 MultiLevel

Алгоритм основан на оптимизации модулярности. Как и в многих предыдущих методах, каждой вершине сначала ставится в соответствие по сообществу. Далее чередуются следующие этапы:

1. Первый этап

- Для каждой вершины перебираем её соседей
- Перемещаем в сообщество соседа, при котором модулярность увеличивается максимально
- Если перемещение в любое другое сообщество может только уменьшить модулярность, то вершина остаётся в своём сообществе
- Последовательно повторяем, пока какое-либо улучшение возможно

2. Второй этап

- Создать метаграф из сообществ-вершин. При этом рёбра будут иметь веса, равные сумме весов всех рёбер из одного сообщества в другое или внутри сообщества (т.е. будет взвешенная петля)

- Перейти на первый этап для нового графа

Алгоритм прекращает работу, когда на обоих этапах модулярность не поддаётся улучшению. Все исходные вершины, которые входят в финальную метавершину, принадлежат одному сообществу.

Несколько замечаний:

- На первом этапе вершина может рассматриваться несколько раз
- Порядок перебора не сильно влияет на точность, однако может существенно влиять на время работы алгоритма
- На практике оказывается достаточно 3-4 итераций

Для ясности приведём иллюстрацию общей схемы работы алгоритма

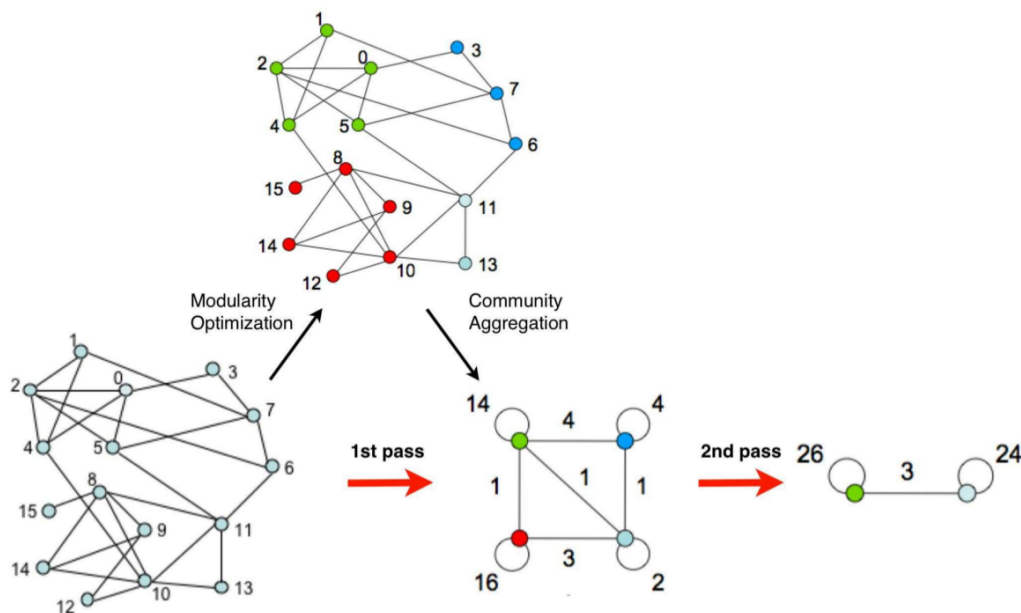


Рис. 4: Два прохода алгоритма. Для первого показаны оба этапа

3 Разбиение на пересекающиеся сообщества

3.1 k-Clique Perlocation

Clique perlocation method (CPM) основан на предположении, что сообщества состоят из пересекающихся полных подграфов. Алгоритм начинает работу с поиска всех клик размера k , после чего строится новый граф, вершинами которого являются найденные клики. Ребро образуется в случае, если пересечение вершин-клик состоит из $k - 1$ вершины исходного графа. Компоненты связности нового графа и будут определять найденные сообщества. Эксперименты показывают, что k хорошо брать в пределах от 3 до 6. Метод хорош своей интуитивностью, однако неприменим на графах с очень большим количеством вершин.

3.2 BigCLAM

Cluster Affiliation Model for Big Networks — вероятностная генеративная модель, сводящая задачу выделения сообществ к задаче неотрицательной матричной факторизации. Для начала немного изменим исходную постановку: теперь у нас будет двудольный граф, в одной доле которого находятся сообщества, а в другой — вершины, причём каждая вершина $u \in V$ не просто принадлежит сообществу $c \in C$, а принадлежит ему с каким-то неотрицательным весом F_{uc} (если не принадлежит, то вес равен нулю):

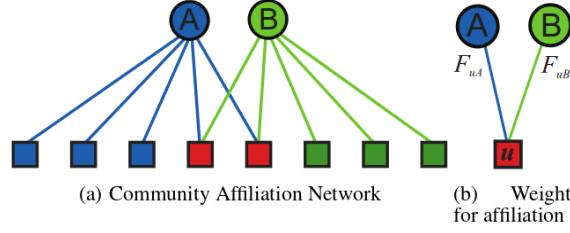


Рис. 5: Двудольный граф принадлежности. Рёбра с нулевыми весами не отображаются

Тогда, для заданной матрицы весов F предполагается, что каждое сообщество соединяет вершины u и v с вероятностью $1 - \exp(-F_{uc} \cdot F_{vc})$. Для всех сообществ вероятность ребра между u и v равна $1 - \exp(-\sum_c F_{uc} \cdot F_{vc})$ или, в краткой форме,

$$p(u, v) = 1 - \exp(-F_u \cdot F_v^T).$$

Теперь можно воспользоваться методом максимизации правдоподобия: для заданного графа $G(V, E)$ будем стараться найти K сообществ, при которых $\hat{F} \in \mathbb{R}^{N \times K}$ доставляет максимум правдоподобия:

$$\hat{F} = \arg \max_{F: F_{uc} \geq 0} l(F)$$

$$l(F) = \log P(G|F) = \sum_{(u,v) \in E} \log(1 - \exp(-F_u \cdot F_v^T)) - \sum_{(u,v) \notin E} F_u \cdot F_v^T$$

Эта задача оптимизации и сводится к неотрицательной факторизации матрицы смежности A графа G . Определять K предлагается по качеству на отложенной выборке.

В итоге получается матрица F степеней принадлежности вершин сообществам. Результирующая принадлежность (принадлежит/не принадлежит) определяется отсечением по порогу.

Осталось лишь описать поиск F . Будем использовать блочно-координатный градиентный метод. Предлагается обновлять каждое F_u при фиксированных остальных F_v , то есть обновлять принадлежности конкретной вершины при фиксированных принадлежностях других. Главная причина фиксации: задача становится выпуклой. То есть, для каждой вершины u решается вспомогательная задача:

$$\hat{F}_u = \arg \max_{F_{uc} \geq 0} l(F_u),$$

где

$$l(F_u) = \sum_{v \in \mathcal{N}(u)} \log(1 - \exp(-F_u \cdot F_v^T)) - \sum_{v \notin \mathcal{N}(u)} F_u \cdot F_v^T,$$

где $\mathcal{N}(u)$ — множество соседей вершины u . Именно благодаря суммированию только по соседям метод получается очень масштабируемым: реальные социальные графы больших размеров очень разреженные, то есть вершины имеют малое по сравнению с количеством вершин в графе количество соседей, вследствие чего каждое обновление имеет близкую к константе вычислительную сложность.

3.3 DEMON

Алгоритм Democratic Estimate of the Modular Organization of a Network является обобщением метода Label Propagation, описанного [ранее](#). Сначала для каждой вершины v строится эго-сеть: выбирается подграф, вершинами которого являются соседи v , а рёбрами — все рёбра между всеми соседями v . Далее, для данной эго-сети запускается Label Propagation, в результате работы которого получается некоторое разбиение $\mathcal{C}(v)$ на сообщества соседей v . После этого его необходимо объединить с итоговым покрытием \mathcal{C} , которое инициализируется пустым. Опишем, как происходит объединение. Два сообщества I и J объединяются в том и только в том случае, если не более ε процентов меньшего из них не содержится в большем из них. Например, для $\varepsilon = 0$ объединение будет происходить только, когда одно из сообществ полностью содержится в другом, а для $\varepsilon = 1$ объединение будет происходить всегда. Теперь можем описать сам алгоритм:

1. Инициализировать $\mathcal{C} = \emptyset$
2. Для вершины $v \in V$ построить эго-сеть и получить её разбиение $\mathcal{C}(v)$ с помощью Label Propagation
3. Для каждого из сообществ в $\mathcal{C}(v)$ и для каждого из сообществ в \mathcal{C} произвести объединение с заданным порогом.
4. Повторять шаги 2-4, пока есть нерассмотренные вершины.

3.4 CONGO

Cluster-Overlap Newman Girvan Optimized algorithm также является обобщением ранее описанного метода, а именно [Edge Betweenness](#). Автор вводит дополнительную операцию разбиения вершины для того, чтобы результатом было разбиение на пересекающиеся сообщества. Раньше каждому ребру ставилось в соответствие значение edge betweenness, с помощью которого происходило последовательное удаление рёбер (удалялось ребро с максимальным edge betweenness, после чего значения пересчитывались). Теперь дополнительно каждой вершине будем ставить в соответствие величину split betweenness. Представим, что вершину v заменили на $v1$ и $v2$. Тогда split betweenness вершины v будет равен количеству кратчайших путей, проходящих через виртуальное ребро между $v1$ и $v2$. При этом смежные с v рёбра делятся между $v1$ и $v2$ таким образом, чтобы величина split betweenness была максимальной для v . Пример:

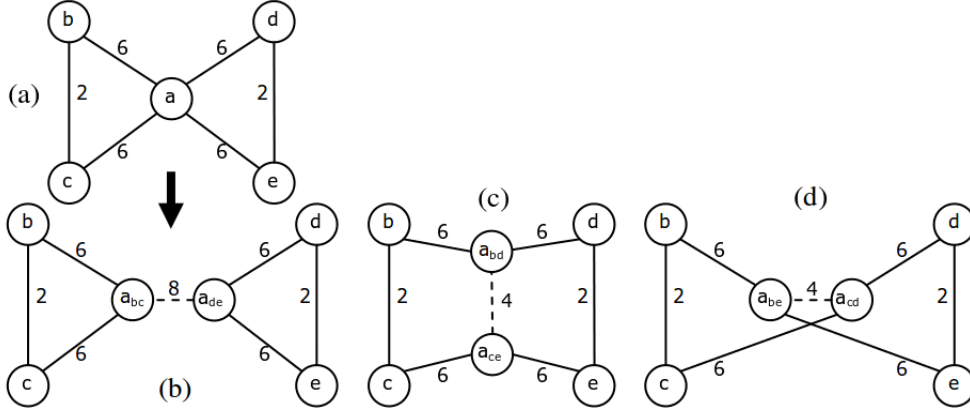


Рис. 6: (a) Граф. (b) Лучшее разбиение вершины а. (c), (d) Другие разбиения.

Далее последовательно повторяются удаления рёбер и разбиения вершин. Однако автор отмечает, что после каждой такой операции пересчитывать значения split и edge betweenness очень дорого, так как приходится "проходить" по всему графу. Поэтому предлагается пересчитывать значения betweenness только для путей, длина которых не превосходит некоторого h , являющегося параметром алгоритма. Итоговый алгоритм:

1. Посчитать все edge betweenness для рёбер и split betweenness для вершин
2. Найти ребро или вершину с максимальным значением betweenness
3. Удалить выбранное ребро или разбить выбранную вершину
4. Пересчитать величины betweenness в h -окрестности данного ребра или данной вершины
5. Повторять шаги 2-5, пока есть рёбра

4 Эксперименты

4.1 Данные

Для начала будем пробовать методы на модельных данных. Создадим граф с l сообществами, каждое из которых имеет по g вершин. При этом рёбра будут генерироваться случайно: с вероятностью p_{in} ребро появляется между вершинами из одного сообщества, с вероятностью p_{out} — между вершинами из разных. В наших экспериментах возьмём $l = 4, g = 64, p_{in} = 0.5$, а величину p_{out} будем варьировать, постепенно зашумляя граф и смотря на зависимость качества методов от зашумлённости графа. На этих данных будем тестировать методы для обнаружения непересекающихся сообществ (ground truth), поэтому будут использоваться метрики при известных ответах, которые будут описаны ниже.

Реальные данные возьмём из соревнования Learning Social Circles in Networks с платформы kaggle, представляющие из себя 110 эго-сетей (граф, вершинами которого являются друзья пользователя, а рёбрами — связи между ними, в случае изначия). На этих сетях также будем тестировать методы для обнаружения пересекающихся сообществ, однако истинные разбиения неизвестны, поэтому будут использоваться соответствующие метрики.

Наконец, данными для пересекающихся сообществ будут также эго-сети пользователей Facebook, для всех вершин каждой из которых известны принадлежности сообществам.

4.2 Метрики качества

Normalized Mutual Information

Если два разбиения похожи, то требуется небольшое количество информации, чтобы восстановить одно из разбиений по другому. Эта идея и лежит в основе метрики NMI, которая является мерой непохожести разбиений (мерой похожести тогда будет являться $1 - \text{NMI}$).

Представим себе два разбиения $\{x_i\}$ и $\{y_i\}$, где i — номер вершины, а x_i и y_i — соответствующие им номера сообществ. Также представим, что метки x и y — значения случайных величин X и Y , имеющих совместное распределение $P(X = x, Y = y) = \frac{n_{xy}}{n}$, где n — общее количество вершин, а n_{xy} — количество вершин, которые в разбиениях $\{x_i\}$ и $\{y_i\}$ имеют метки x и y . Аналогично $P(X = x) = \frac{n_x}{n}$, $P(Y = y) = \frac{n_y}{n}$.

Напомним определение энтропии и условной энтропии распределения:

$$H(X) = - \sum_x P(x) \log P(x), \quad H(X|Y) = - \sum_{x,y} P(x,y) \log P(x|y)$$

Тогда взаимная информация будет определяться как их разность:

$$I(X, Y) = H(X) - H(X|Y),$$

а нормировка производится на сумму отдельных энтропий:

$$I_{\text{norm}}(X, Y) = \frac{I(X, Y)}{H(X) + H(Y)}$$

Эта метрика будет использоваться как для непересекающихся, так и для пересекающихся сообществ с ground truth.

Split-Join Distance

Эта метрика является аналогом редакторского расстояния для разбиений: она измеряет минимальное количество операций, необходимых для перехода от одного разбиения к другому. Операциями могут быть:

- Добавить вершину к сообществу
- Удалить вершину из сообщества

- Создать сообщество с одной вершиной
- Удалить сообщество с одной вершиной

Эта метрика будет использоваться для непересекающихся сообществ с ground truth.

Modularity

Значения функционала модулярности, который многократно упоминался ранее, будем смотреть для эго-сетей с kaggle как меру правильности работы алгоритмов обнаружения непересекающихся сообществ для неразмеченных данных.

Omega Index

Omega Index измеряет количество согласованных пар вершин в двух покрытиях графов. Две вершины назовём согласованными, если они лежат в одинаковом количестве сообществ. То есть, Omega Index считает, сколько пар вершин принадлежат одновременно одному сообществу, двум сообществам и так далее.

Пусть K_1 и K_2 — количество сообществ в покрытиях C_1 и C_2 соответственно. Тогда

$$\omega(C_1, C_2) = \frac{\omega_u(C_1, C_2) - \omega_e(C_1, C_2)}{1 - \omega_e(C_1, C_2)},$$

где

$$\omega_u(C_1, C_2) = \frac{1}{M} \sum_{j=0}^{\max(K_1, K_2)} |t_j(C_1) \cap t_j(C_2)|,$$

$$\omega_e(C_1, C_2) = \frac{1}{M^2} \sum_{j=0}^{\max(K_1, K_2)} |t_j(C_1)| \cdot |t_j(C_2)|.$$

Здесь M равно $n(n-1)/2$ — количество пар вершин, а $t_j(C)$ — множество пар вершин, которые встречаются в покрытии C ровно j раз.

Omega Index равен единице, только в случае, если $\omega_u(C_1, C_2)$ равен единице, что означает точное совпадение C_1 и C_2 .

Эта метрика будет использоваться для пересекающихся сообществ с ground truth.

4.3 Результаты

Метод Edge Betweenness показал непозволительно высокую вычислительную сложность (более, чем в 10 раз дольше любого другого алгоритма) даже на простых модельных данных, поэтому из испытаний был исключён. Алгоритмы Label Propagation и Infomap таким недостатком не обладали, однако оказались крайне чувствительны даже к небольшому шуму. Лучшее всего с задачей справились методы MultiLevel и WalkTrap.

На реальных данных метод Label Propagation оказался не так плох, как на модельных, в отличие от Infomap, который снова не справился с задачей. В среднем, MultiLevel справился лучше других методов.

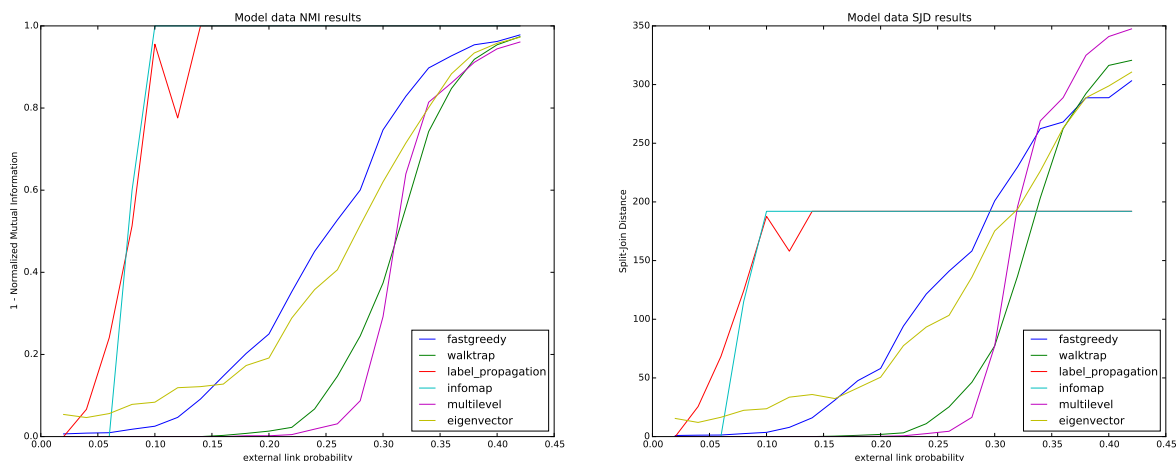


Рис. 7: Результаты работы методов обнаружения непересекающихся сообществ на модельных данных

Как видно, методы показали себя достаточно разнородно. Хочется отметить, что BigCLAM оказался действительно хорошо масштабируемым, практически не имел сложности по памяти и был очень быстр. Также неплох по всем этим параметрам оказался алгоритм DEMON. Весьма долго работал на любых данных CONGO, даже несмотря на название описывающей его статьи "A Fast Algorithm to Find Overlapping Communities in Networks". Clique Perlocation Method, в целом, показал себя неплохо на небольших сетях, однако при увеличении количества рёбер переставал справляться с задачей.

2do: разобраться с русской кодировкой и правильным оформлением списка литературы
2do: замечания в блокноте

Список литературы

- [1] Clauset, A. Finding community structure in very large networks / Aaron Clauset, M. E. J. Newman, Cristopher Moore // Physical Review E. — 2004. — <http://arxiv.org/abs/cond-mat/0408187>.
- [2] Girvan, M. Community structure in social and biological networks / Michelle Girvan, M. E. J. Newman // Proceedings of the National Academy of Sciences. — 2001. — <http://arxiv.org/abs/cond-mat/0112110>.
- [3] igraph library. — 2016. — <http://igraph.org/python/>.
- [4] Raghavan, U. N. Near linear time algorithm to detect community structures in large-scale networks / Usha Nandini Raghavan, Reka Albert, Soundar Kumara // Physical Review E. — 2007. — <http://arxiv.org/abs/0709.2938>.

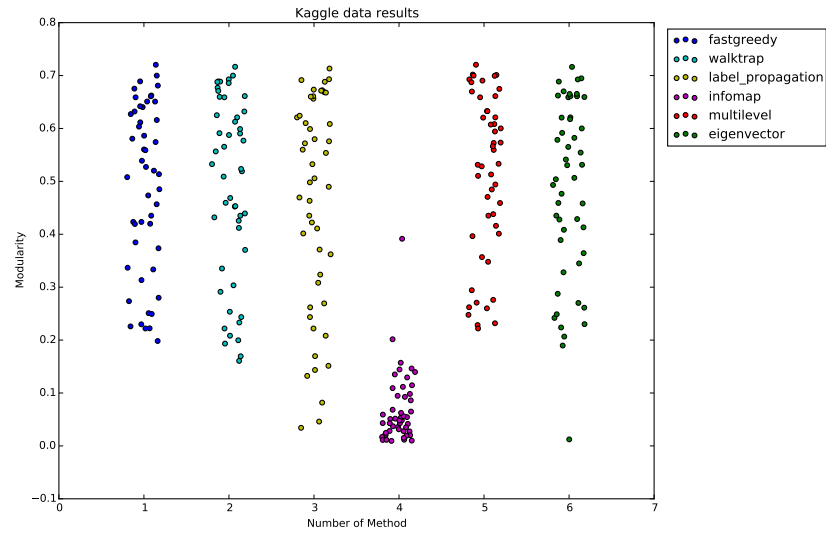


Рис. 8: Результаты работы методов обнаружения непересекающихся сообществ на реальных эго-сетях. Точкой обозначается результат на отдельной эго-сети. Ось абсцисс введена искусственно для удобства читателя

- [5] Slavnov, K. A. Social graph analysis. — 2015. — http://www.machinelearning.ru/wiki/images/6/60/2015_417_SlavnovKA.pdf.

Различные эго-сети		1–NMI			
# вершин	# рёбер	BigCLAM	CPM	CONGO	DEMON
59	146	0.2497	0.2975	0.1703	0.2315
66	270	0.3580	0.4305	0.3513	0.4305
159	1693	0.2898	0.3454	0.3278	0.4057
170	1656	0.1460	0.2730	0.1381	0.1585
227	3192	0.2342	0.2700	0.2366	0.2408
347	2519	0.0501	0.0476	0.0477	0.0489
547	4813	0.0120	0.0157	0.0241	0.0400
755	30025	0.1207	Memory	Time	0.1821
792	14024	0.2401	Memory	Time	0.3620
1045	26749	0.1382	Memory	Time	0.1245

Рис. 9: 1–Normalized Mutual Information для методов обнаружения пересекающихся сообществ. Memory — метод оказался слишком сложным по памяти. Time — метод не сошелся за приемлемое время

Различные эго-сети		Omega Index			
# вершин	# рёбер	BigCLAM	CPM	CONGO	DEMON
59	146	0.1058	0.1835	0.1299	0.0939
66	270	0.3307	0.3001	0.4413	0.3001
159	1693	0.3266	0.2622	0.3376	0.1521
170	1656	0.0521	0.0914	0.0749	0.0319
227	3192	0.0000	0.1793	0.0488	0.1817
347	2519	0.1257	0.2424	0.0619	0.1681
547	4813	0.0518	0.0080	0.0045	0.0043
755	30025	0.3917	Memory	Time	0.0000
792	14024	0.3378	Memory	Time	0.0181
1045	26749	0.1936	Memory	Time	0.0000

Рис. 10: Omega Index для методов обнаружения пересекающихся сообществ