

Кластеризация. Графовый подход, ЕМ и k means, иерархическая кластеризация

Сергей Лисицын

lisitsyn.s.o@gmail.com

6 апреля 2011 г.

Задача кластеризации

Clustering

- Имеется некоторое множество объектов
- Решением задачи является некоторое разбиение множества на непересекающиеся кластеры
- Кластеры – некоторые характерные группы объектов
- Количество кластеров для разбиения может быть задано изначально
- Эффективность разбиения оценивается функционалами внутриклассового и межклассового расстояния
- Задача ещё более некорректна чем классификация

Общая задача кластеризации

Пусть имеется некоторое множество объектов \mathcal{X} и конечное множество кластеров \mathcal{C} . Необходимо найти функциональную зависимость

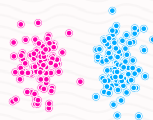
$$f: \mathcal{X} \rightarrow \mathcal{C},$$

причём (в отличие от классификации) ни одно значение f на \mathcal{X} априори не известно.

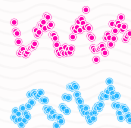
В задаче кластеризации функция f должна быть построена таким образом, чтобы выделять наиболее характерные группы объектов.

Свойства кластеров

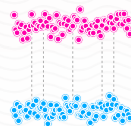
Компактность (compactness): кластеры компактны – объекты лежат внутри некоторого плотного множества



Связность (connectivity): объекты кластера можно выделить в связанные структуры



Пространственное разделение (spatial separation): кластеры отстоят друг от друга на некотором расстоянии



Функционалы качества кластеризации

- Межкластерное расстояние

$$\rho_c(c_1, c_2) = \sum_{f(x_1)=c_1, f(x_2)=c_2} \rho(x_1, x_2)$$

Эффективная кластеризация подразумевает максимизацию межкластерных расстояний: $\sum_{\substack{c_1 \neq c_2 \\ c_1, c_2 \in \mathcal{C}}} \rho_c(c_1, c_2) \rightarrow \max$

- Внутрикластерное расстояние

$$\rho_{\text{inner}}(c) = \sum_{f(x_1)=f(x_2)=c} \rho(x_1, x_2)$$

Эффективная кластеризация подразумевает минимизацию внутрикластерных расстояний: $\sum_{c \in \mathcal{C}} \rho_{\text{inner}}(c) \rightarrow \min$

- В некоторых случаях минимизируется отношение межкластерного и внутрикластерного расстояний

Цели кластеризации

- Выделение наиболее характерных групп объектов
- Выяснение структуры множества объектов
- Сокращение объема выборки объектов
- Обнаружение «необычных» объектов (novelty detection)

Графовый подход к кластеризации

- Множество объектов в метрическом пространстве представимо в виде неориентированного взвешенного графа, где вес ребра это значение метрики между его объектами-вершинами
- Кластеры в таком графе – непересекающиеся подграфы исходного
- Графы сами по себе достаточно удобны для кластеризации
- Представление объектов в таком случае знать не обязательно, достаточно расстояний между ними, вычисленных даже косвенно
- Кроме того, такие графы-кластеры удобно визуализировать

Алгоритм выделения связанных компонент

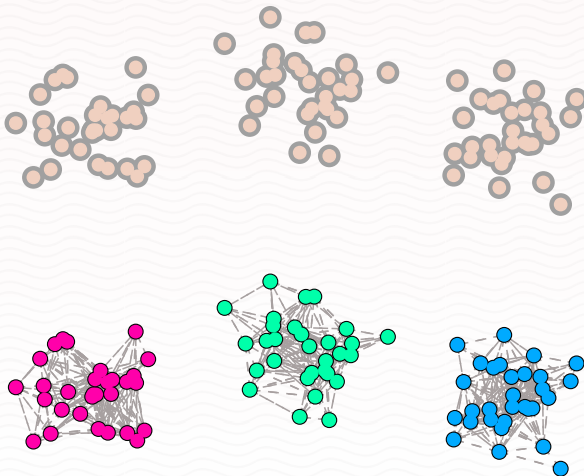
- Кластеры в пространстве можно представлять как группы связанных объектов некоторого графа – связанные компоненты
- Связная компонента – подграф, в котором все вершины взаимно достижимы
- Удаляя (или не создавая) ребра графа с весом, большим, чем некоторый параметр R_{max} , можно разбить граф на некоторое количество связанных компонент
- В распавшемся графе будут выделены некоторые характерные структуры с минимальным внутрикластерным расстоянием

Описание алгоритма выделения связных компонент

Вход: множество объектов \mathcal{X} , параметр R_{max}

Выход: множество пар (объект-кластер)

1. Создать граф без ребёр с узлами, соответствующими элементам множества \mathcal{X}
2. Соединить все узлы графа, соответствующие элементам $x_i, x_j \in \mathcal{X}, \forall i \neq j$, для которых $\rho(x_i, x_j) < R_{max}$
3. Возвратить соответствующие вершинам графа пары (элемент множества \mathcal{X} – номер связной компоненты)



```
1 class ConnectedComponentClusterer():
2     def clusterize(self, objs, R_max):
3         cluster_graph = graph()
4         N = len(objs)
5         cluster_graph.add_nodes(range(N))
6         for i in range(N):
7             for j in range(N):
8                 if not i==j and metric(objs[i], objs[j]) < R_max:
9                     try: cluster_graph.add_edge((i, j))
10                        except: pass
11         component_idx = connected_components(cluster_graph)
12         clusters = defaultdict(list)
13         for i in component_idx.keys():
14             clusters[component_idx[i]].append(objs[i])
```

Алгоритм построения покрывающего дерева

Spanning tree clustering

- Минимальное покрывающее дерево – древовидная структура с минимальным суммарным расстоянием между соседними объектами
- При удалении из построенного дерева $K - 1$ самых «длинных» ребёр maximизируется межкластерное расстояние
- Минимум суммарного расстояния между объектами минимизирует внутрикластерное расстояние
- Количество кластеров в таком случае задаётся заранее, алгоритм ведёт себя предсказуемое

Описание алгоритма кластеризации с помощью минимального покрывающего дерева

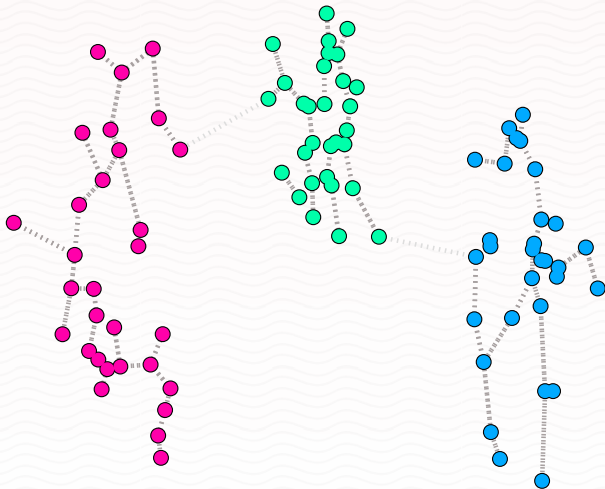
Вход: множество объектов \mathcal{X} , параметр K – количество кластеров

Выход: множество пар (объект-кластер)

1. Построить минимальное покрывающее дерево с узлами $0, 1, \dots$, для которого

$$\sum_{i,j} \rho(x_i, x_j) \rightarrow \min$$

2. Удалить $K - 1$ ребёр (i, j) максимальной метрики $\rho(x_i, x_j)$ построенного покрывающего дерева
3. Возвратить соответствующие вершинам графа пары (элемент множества \mathcal{X} – номер связной компоненты)

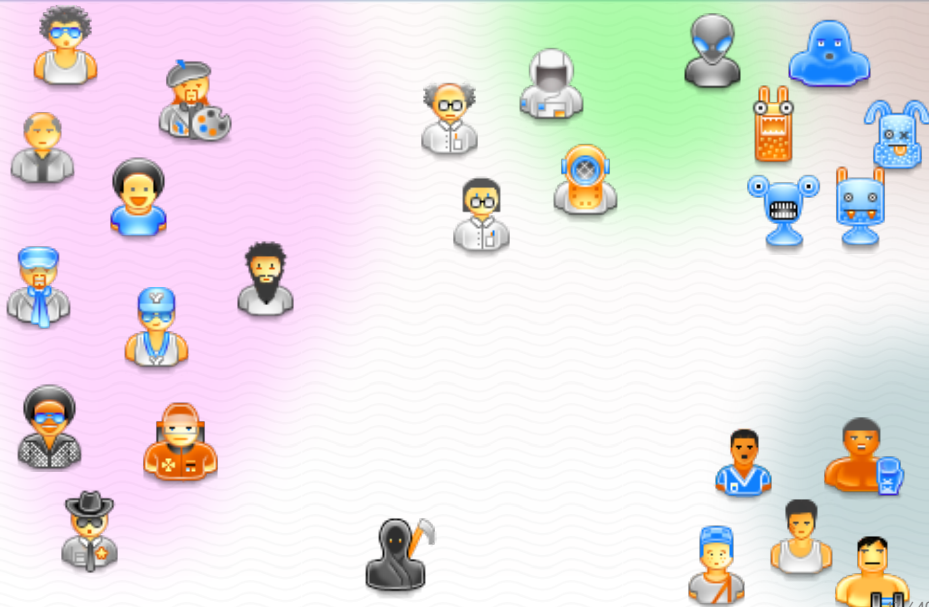


```
1 class SpanningTreeClusterer():
2     def clusterize(self, objs, metric, K):
3         N = len(objs)
4         c_graph = graph()
5         c_graph.add_nodes(range(N))
6         for i in range(N):
7             for j in range(i, N):
8                 c_graph.add_edge((i, j), wt=metric(objs[i], objs[j]))
9         stree, c_graph = minimal_spanning_tree(c_graph), graph()
10        c_graph.add_nodes(range(N))
11        del stree[stree.keys()[0]]
12        for i in stree.keys():
13            c_graph.add_edge((i, stree[i]),
14                             wt=metric(objs[i], objs[stree[i]]))
15        for _ in range(K-1):
16            max_edge = max(c_graph.edges(),
17                           key=lambda x: metric(objs[x[0]], objs[x[1]]))
18            c_graph.del_edge((max_edge))
19        component_idx = connected_components(c_graph)
20        clusters = defaultdict(list)
21        for i in component_idx.keys():
22            clusters[component_idx[i]].append(objs[i])
23        return clusters
```

Пример социальной кластеризации

- Объекты пространства – пользователи социальной сети (электронной почты, ICQ, etc)
- Никакой информации о самих объектах не имеется, можно определить только расстояния между объектами
- Метрика между двумя пользователями может определяться как обратное нормированное количество сообщений между ними
- Результатом кластеризации будут группы людей, наиболее тесносвязанные друг с другом (можно на это надеяться)

Пример социальной кластеризации



Эффективность алгоритмов

- Для обоих алгоритмов необходимо найти все взаимные метрики за $O(n^2)$, выделение кластеров выполняется поиском в ширину или глубину $O(|V| + |E|)$ (однако для построения покрывающего дерева существуют особые алгоритмы)
- Оба алгоритма неустойчивы к близкорасположенным кластерам – нет «баланса» между уменьшением функционалов качества, минимизируется только внутрикластерное расстояние
- Алгоритм выделения связных компонент плохо обусловлен по своему параметру R_{max} : небольшие изменения R_{max} могут повлечь резкое увеличение или уменьшение количества кластеров
- Алгоритм кластеризации с помощью покрывающего дерева выделяет заданное количество кластеров

Смеси распределений

- Выборка с группами может рассматриваться как группа нескольких многомерных распределений, такое предположение задаёт общую плотность распределения

$$p(x) = \sum_c P(c)p_c(x),$$

где $p_c(x)$ – функция плотности распределения кластера c , а $P(c)$ – априорная вероятность класса

- Кластер объектов – компонента этой смеси распределений с плотностью f_c
- Кластеризация сводится к нахождению компонент смеси распределений по эмпирическим оценкам плотности, а принадлежность конкретных объектов определяется максимумом $P(c)p_c(x)$

ЕМ-алгоритм

Expectation-Maximization

- Описан в 1977 году Демпстером, Лэйрдом и Рабином и предназначен для оценок максимального правдоподобия в статистических моделях
- В случае кластеризации ЕМ-алгоритм позволяет найти компоненты смеси распределений
- По полученным компонентам смеси распределений классификация производится достаточно легко
- Для кластеризации, как правило используют гауссианы:

$$p_c(x) = \mathcal{N}(x, \mu_c, \mathbf{R}_c)$$

- Алгоритм итеративен и выполняется до тех пор, пока всем объектам не будут назначены «стабильные» кластеры, не изменяющиеся в течение следующих итераций

Вход: множество объектов \mathcal{X} , количество кластеров k

Выход: множество пар (объект-кластер)

1. Задать начальные приближения $P(c) = 1/k$, за μ_c принять случайные объекты выборки, рассчитать $diag \mathbf{R}_c = \frac{\sum_{x \in \mathcal{X}} g_c(x)(x - \mu_c)^2}{|\mathcal{X}|P(c)}$
2. До тех пор, пока кластеры не стабилизируются:

2.1 E-шаг

Вычисление апостериорных вероятностей

$$g_c = \frac{P(c)\mathcal{N}(x, \mu_c, \mathbf{R}_c)}{\sum_k P(k)\mathcal{N}(x, \mu_k, \mathbf{R}_k)}$$

2.2 M-шаг

Пересчёт параметров гауссианов

$$P(c) = \frac{\sum_x g_c(x)}{|\mathcal{X}|}, \quad \mu_c = \frac{\sum_x g_c(x)x}{|\mathcal{X}|P(c)}, \quad diag \mathbf{R}_c = \frac{\sum_{x \in \mathcal{X}} g_c(x)(x - \mu_c)^2}{|\mathcal{X}|P(c)}$$

2.3 Для каждого объекта выбрать кластер $\arg \max_c g_c$

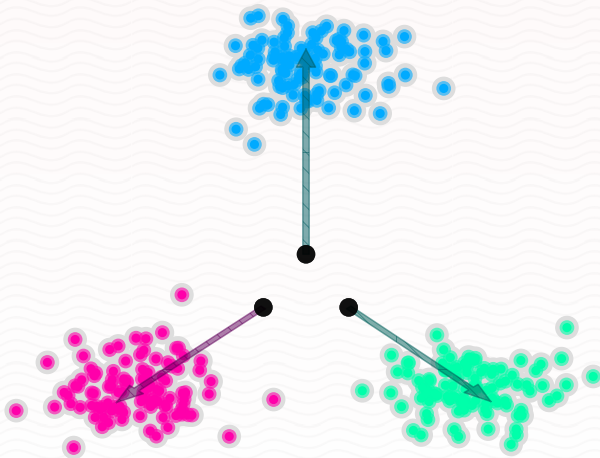
Эффективность

- Алгоритм EM в самом общем виде не очень практичен
- В случае плохого выбора скрытых переменных и неправильной оценки количества кластеров алгоритм плохо сходится
- Алгоритм очень трудоёмок в вычислениях
- Итеративная природа алгоритма требует определять либо количество итераций, либо критерий останова

Алгоритм кластеризации k средних

k means clustering

- Является упрощенным вариантом EM-алгоритма
- Идея алгоритма состоит в «стабилизации» центров кластеров на M-шаге и жёсткой привязке объектов к ближайшему центру на E-шаге
- Алгоритм итеративен: выбирается либо количество итераций, либо условие сходимости
- Выделяет только заданное количество компонент, но существуют модификации для переменного параметра k



Описание алгоритма k средних

Вход: множество объектов \mathcal{X} , количество центроидов k , количество итераций или условие останова алгоритма

Выход: множество пар (объект-кластер)

1. Занести в множество центроидов \mathbf{C} случайные объекты из \mathcal{X} , не изменяя \mathcal{X}
2. До тех пор, пока центроиды не стабилизировались (или не выполнилось заданное количество итераций):
 - 2.1 Для каждого объекта $x \in \mathcal{X}$: задать объекту x ближайший к нему центроид c : $f(x) = c, \rho(c, x) \rightarrow \min_c$
 - 2.2 Для каждого центроида $c \in \mathbf{C}$: пересчитать положение

$$c = \frac{\sum_{x, f(x)=c} x}{\sum_{x, f(x)=c} 1}$$

3. Вернуть пары (объект из \mathcal{X} – номер ближайшего к нему центроида)

```
1 class KMeansClusterer():
2     def clusterize(self, objects, k, max_iter=10):
3         cs = numpy.array(random.sample(objects, k))
4         for iter in range(max_iter):
5             new_cs = numpy.zeros([k, 2])
6             counts = numpy.zeros(len(cs))
7             for x in objects:
8                 nearest_c = min(cs, key=lambda c: metric(c, x))
9                 for (index, centroid) in enumerate(cs):
10                     if numpy.all(centroid == nearest_c):
11                         new_cs[index] += x
12                         counts[index] += 1
13         cs = numpy.array(map(lambda x, y: x/y, new_cs, counts))
14         clusters = [[c] for c in cs]
15         for x in objects:
16             nearest_c = min(cs, key=lambda c: metric(c, x))
17             for cluster in clusters:
18                 if numpy.all(cluster[0] == nearest_c):
19                     cluster.append(x)
20         for cluster in clusters:
21             cluster.pop(0)
22         return clusters
```

Алгоритм k средних

k means

- k means – один из самых популярных алгоритмов кластеризации и находит применение во многих задачах
- Сходимость сильно зависит от начального выбора центроидов
- Композиция алгоритма с каким-либо более простым алгоритмом может помочь достичь лучшей сходимости
- Каждая итерация линейна по размеру выборки – $O(n)$
- Существует множество различных модификаций для лучшего выбора центроидов и лучшей сходимости

Иерархический подход

- Из кластеров может быть составлена иерархия: в вершине иерархии множество всех объектов, ниже группы объектов, в самом низу
- На любом уровне иерархии можно получить нужную степень разбиения на кластеры
- Метрика между объектами вполне однозначна, но метрику между кластерами следует определить каким-либо образом

Межкластерные расстояния

- Расстояние между ближайшими элементами (single linkage)

$$\rho(c_1, c_2) = \min_{x_1 \in c_1, x_2 \in c_2} \rho(x_1, x_2)$$

- Расстояние между самыми дальними элементами (complete linkage)

$$\rho(c_1, c_2) = \max_{x_1 \in c_1, x_2 \in c_2} \rho(x_1, x_2)$$

- Среднее групповое расстояние (group linkage)

$$\rho(c_1, c_2) = \frac{\sum_{x_1 \in c_1, x_2 \in c_2} \rho(x_1, x_2)}{|c_1| \cdot |c_2|}$$

Агломеративный подход к кластеризации

Agglomerative clustering

- Наиболее простой метод построения иерархии классов
- Изначально каждый объект содержится в отдельном кластере
- Последовательным объединением двух кластеров с минимальной взаимной метрикой
- С каждой итерацией количество кластеров уменьшается, процесс останавливается при достижении нужного их количества

Агломеративная кластеризация

Вход: множество объектов \mathcal{X} , параметр k количества кластеров

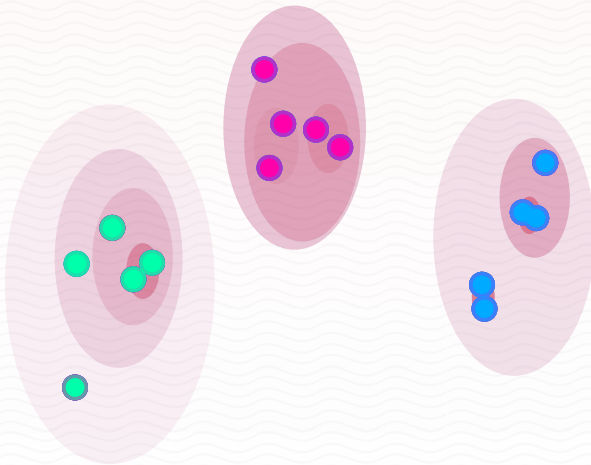
Выход: множество пар (объект-кластер)

- Создать множество C , состоящее из одноэлементных кластеров, элементов множества \mathcal{X}
- Пока количество кластеров больше, чем k
 - Выбрать два наиболее близких по метрике кластера из C :

$$(S_1, S_2) = \arg \min_{S_1, S_2 \in C, S_1 \neq S_2} \rho(S_1, S_2)$$

- Заменить в C кластеры S_1, S_2 кластером $S_1 \cup S_2$, то есть объединить их в один, исключая исходные
- Возвратить пары (объект из \mathcal{X} – номер кластера, в который попал этот объект)

```
1 class AgglomerativeClusterer():
2     def clusterize(self, objects, till, metric):
3         clusters = [[x] for x in objects]
4         iteration = 0
5         while len(clusters) > till:
6             n_idx = (0, 1)
7             n_dist = metric(clusters[0], clusters[1])
8             for i, cluster_lhs in enumerate(clusters):
9                 for j, cluster_rhs in enumerate(clusters):
10                     m = metric(cluster_lhs, cluster_rhs)
11                     if n_dist > m and not i == j:
12                         n_dist = m
13                         n_idx = (i, j)
14             agglomerated = clusters[n_idx[0]] + clusters[n_idx[1]]
15             clusters = \
16             [clusters[i] for i in range(len(clusters)) \
17              if i != n_idx[0] and i != n_idx[1]]
18             clusters.append(agglomerated)
19             iteration += 1
20         return clusters
```

Кластеризация дивизимным анализом

Divisive clustering

- Более сложный по сравнению с агломеративной кластеризацией метод
- Изначально все объекты содержатся в одном кластере
- Кластер(ы) некоторым образом разбиваются дихотомически
- С каждой итерацией количество кластеров увеличивается, процесс останавливается при достижении нужного их количества

Дивизимная кластеризация

Вход: множество объектов \mathcal{X} , параметр k количества кластеров

Выход: множество пар (объект-кластер)

- Все объекты из \mathcal{X} отнести к одному классу из \mathcal{C}
- Пока количество кластеров меньше, чем k
 - Разбить кластер с наибольшим внутрикластерным расстоянием на два с помощью какого-либо алгоритма кластеризации
- Возвратить пары (объект из \mathcal{X} – номер кластера, в который попал этот объект)

Источники

1. Воронцов, К.В. «Лекции по алгоритмам кластеризации и многомерного шкалирования»
2. Moore A. "K-means and Hierarchical Clustering"
3. Singh A. "Spectral Clustering"
4. Николенко С.И. «Алгоритмы кластеризации»
5. Ng A. "The k-means clustering algorithm"