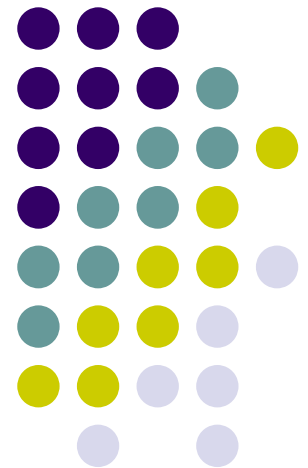
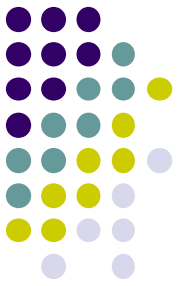


Кластер-анализ

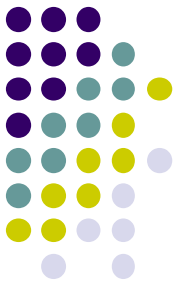
Корлякова М.О.
2018



Обучение

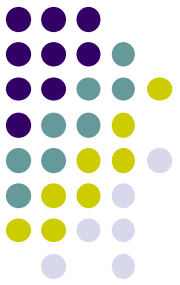


- Без учителя (выделение классов)
- С учителем (отнесение к классу)



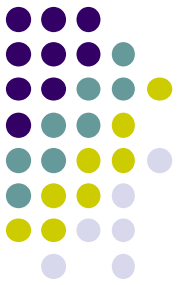
Задача обобщения

- n – число координат
- D_n – пространство в координатах $\langle x_1, x_2, \dots, x_n \rangle$
- x_i – номинальный, дискретный упорядоченный или непрерывный.
- X – объект из D_n
- $X = (a_1, a_2, \dots, a_n)$, a_i – значение x_i для X



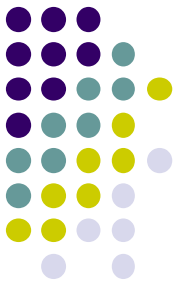
Задача обобщения

- O – множество объектов известных в D_n
- V – общее множество объектов в D_n



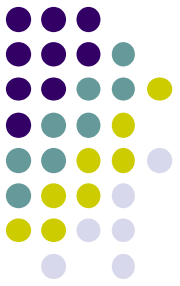
Задача обобщения

- Необходимо в D_n определить группы C_i , которые сформированы на основании близости по мере $m(X_j, C_i)$.
- $P(X_j, C_i) = 1$ для $m(X_j, C_i) < g$
- $P(X_j, C_i) = 0$ для $m(X_j, C_i) > g$



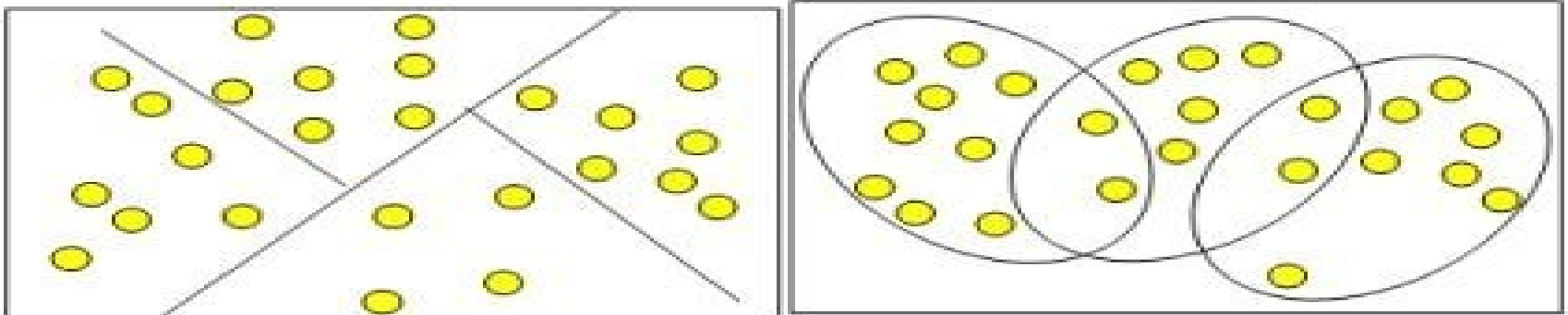
Кластеризация

- $T = \{(X_i)\}$
- Цель: Найти классы

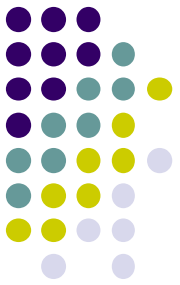


Типы кластерного анализа

- Одноуровневый
 - Фиксация числа кластеров
 - Фиксация размеров кластера
- Иерархический
- Кластеры без пересечения и с пересечением

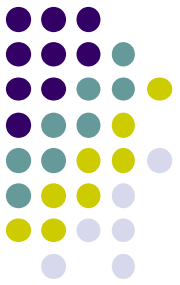


Расстояния между объектами



- Метрики : Минковский
- Меры: Хемминг
- И МНОГО ДРУГИХ МЕТОДОВ!!!!

Расстояние между множествами

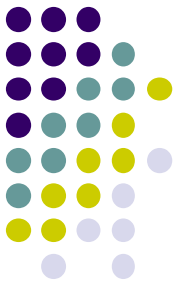


- Ближний сосед
- Средний
- Дальний сосед

- Метрика Хаусдорфа

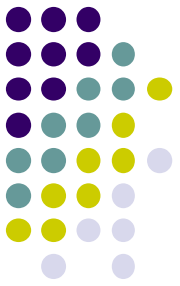
- И МНОГО ДРУГИХ МЕТОДОВ!!!!

Методы одноуровневого анализа



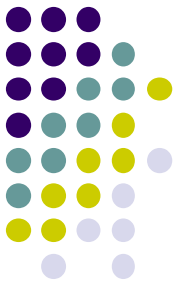
- Фиксация числа кластеров (k - средних)
 - Минимизация различий в примерах кластера
- Фиксация размера - FOREL
 - Поиски сгущений в пространстве
- Комбинации

Порядок решения задачи кластеризации



- Определить D_n – пространство признаков и $T=\{(X_i)\}$ – примеры
- Выбрать способ вычисления расстояния между объектами
- Определить отношение эквивалентности
- Определить тип алгоритма кластеризации.

Алгоритм кластеризации по максимальному расстоянию

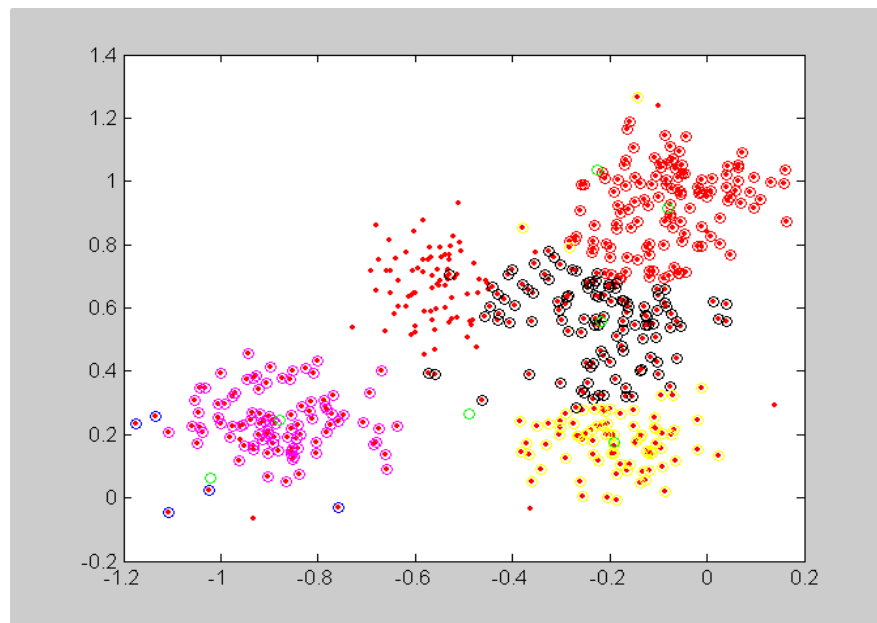
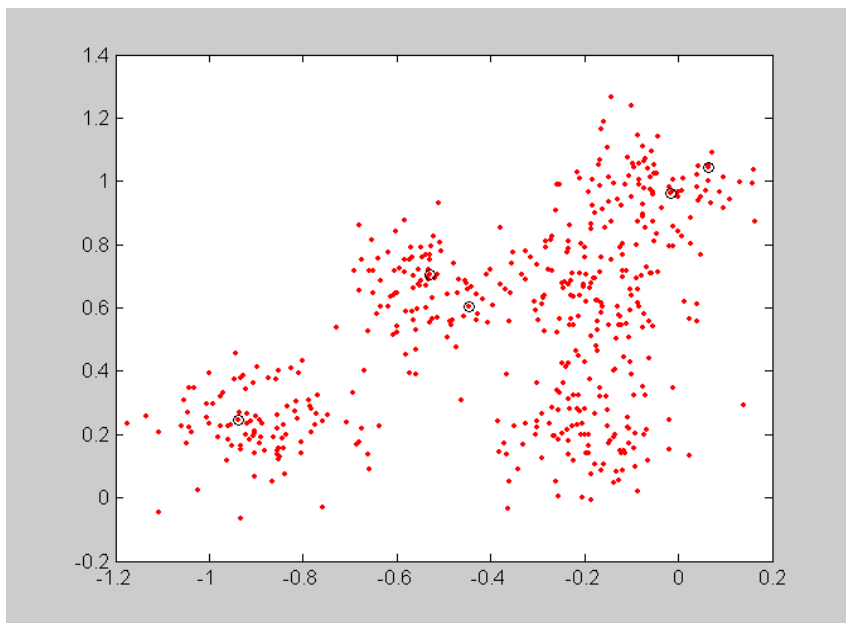


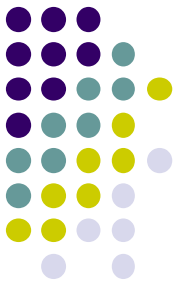
- Фиксирует размер кластера
 1. Зафиксировать диаметр кластера d_{\max}
 2. Установить центр кластера в свободный (вне других кластеров) объект X_i из обучающей выборки.
 3. Добавить к кластеру объект X_j такой, что для всех X_i из текущего кластера.
 4. Продолжать процесс расширения кластера пока не исчерпана выборка примеров.
 5. Вычислить центр кластера.
 6. Если остались свободные примеры, то перейти к процедуре формирования кластеров (п.2.), иначе остановить процесс.

Пример



Исходное множество Результат

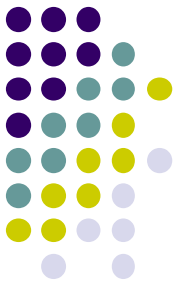




Алгоритм к-средних

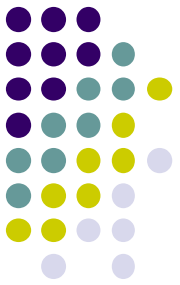
- Фиксирует число классов
 1. Номер итерации $s=0$
 2. Связать с каждым кластером K_j объект X_i из обучающей выборки (случайно).
 3. если число кластеров меньше N , то перейти к процедуре формирования кластеров (п.4.).
 4. Вычислить расстояния от всех объектов до всех центров кластеров.
 5. присоединить объект X_i к кластеру C_k , если $C_k = \min_{j=1..M} d(X_i, C_j)$
 6. повторить для всех объектов выборки.
 7. вычислить новое положение центров кластеров
$$\mathbf{centr}_b = \frac{1}{N_b} \sum_{i=1}^{N_b} \mathbf{x}_i$$
 8. где N_b – число примеров множества C_b .
 9. Повторять от п.4. пока кластер смещается более чем на ε (задано пользователем), иначе остановить процесс.

Достоинства алгоритма k-средних:



- простота использования;
- быстрота использования;
- понятность и прозрачность алгоритма.

Недостатки алгоритма k-средних:

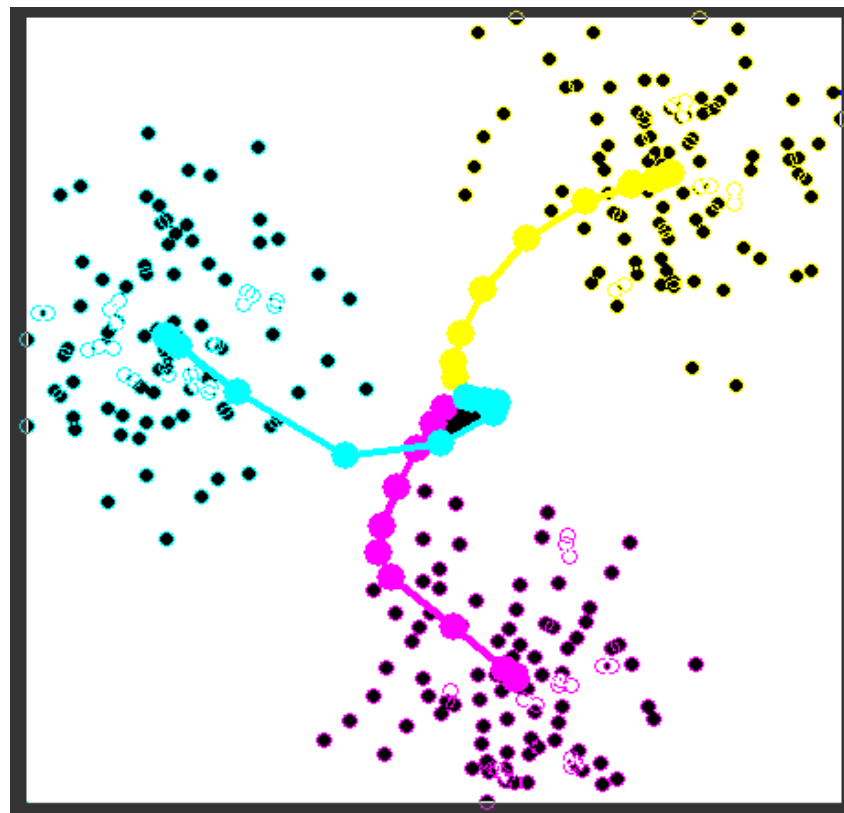
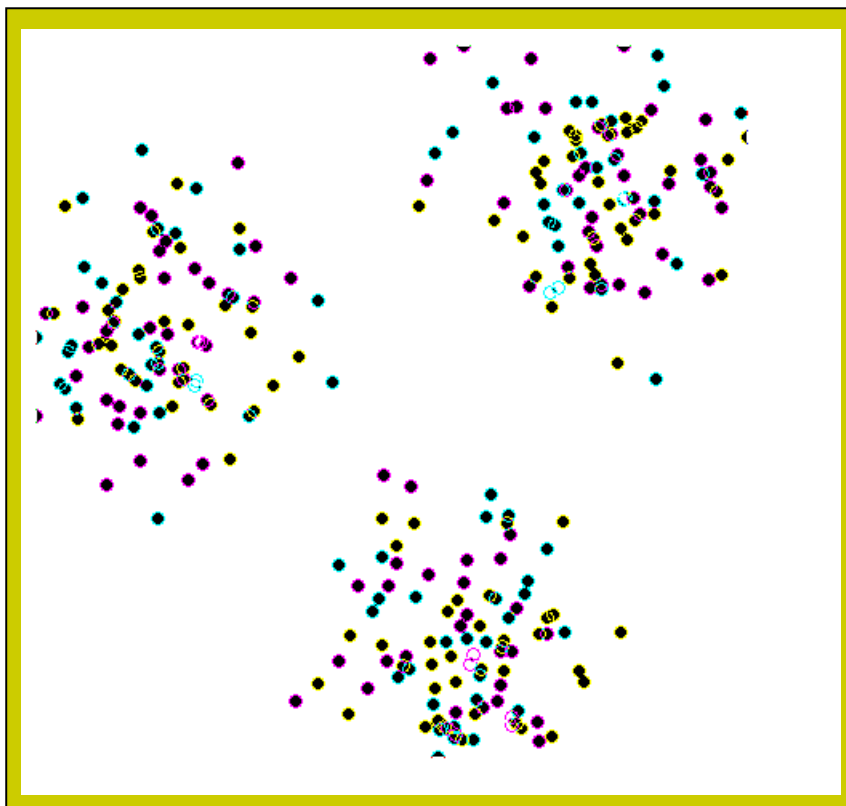


- Чувствителен к выбросам, которые могут искажать среднее;
- Может медленно работать на больших базах данных.

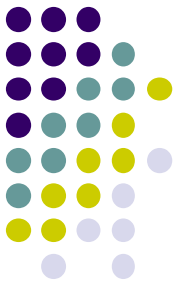
Пример



Исходное множество Результат



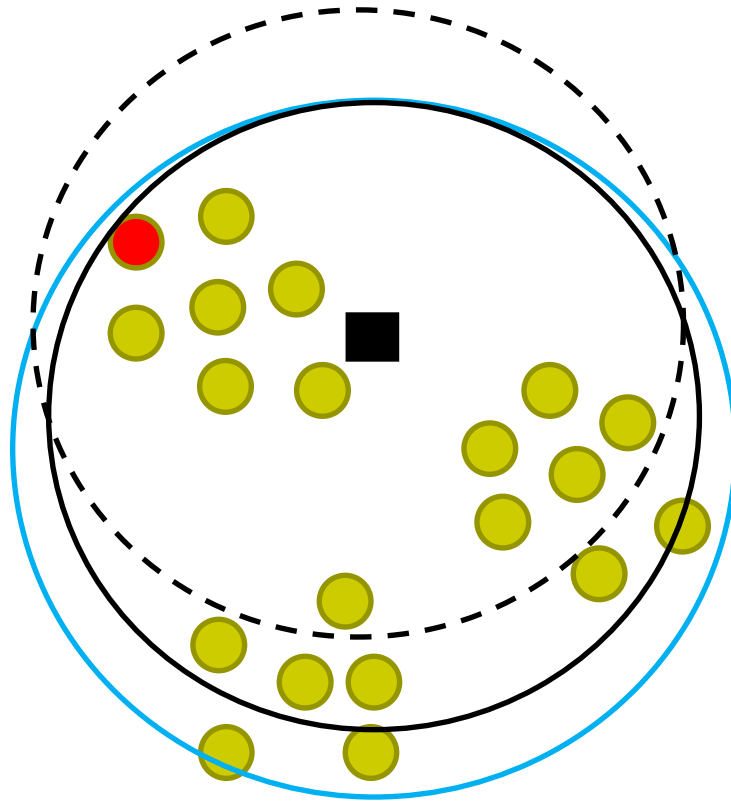
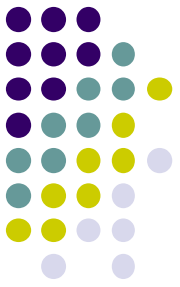
FOREL



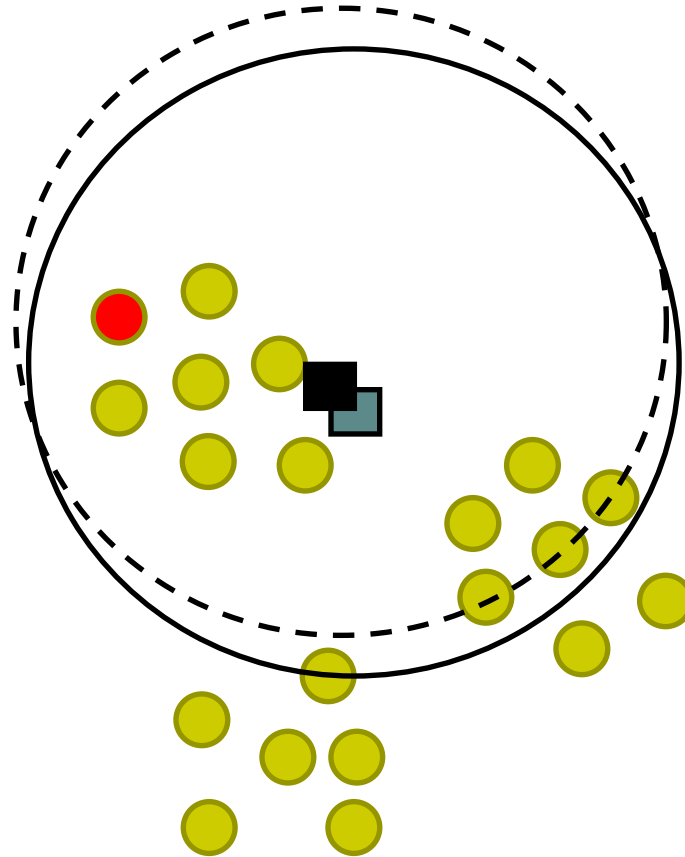
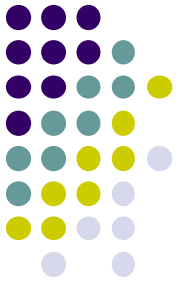
1. $T = \{X_i\}, i=1, N$
2. $R_0 = \max(\text{dist}(X_i, X_j))$ X_i, X_j из T – радиус кластера
3. $C_k = X_i$ – центр кластера
4. X_j из C_k , если $\text{dist}(C_k, X_j) < R_0 * 0.9$
5. Новый $C_k^* = 1/|C_k| \sum X_j$
6. Если $|C_k^* - C_k| > a$, то к 4, иначе 7
7. Точки C_k – исключаем из T
8. Если $|T| > 0$, то $k=k+1$ и к 1, иначе конец.

Критерий $\min \sum \text{dist}(C_k, X_i), X_i \text{ из } C_k$

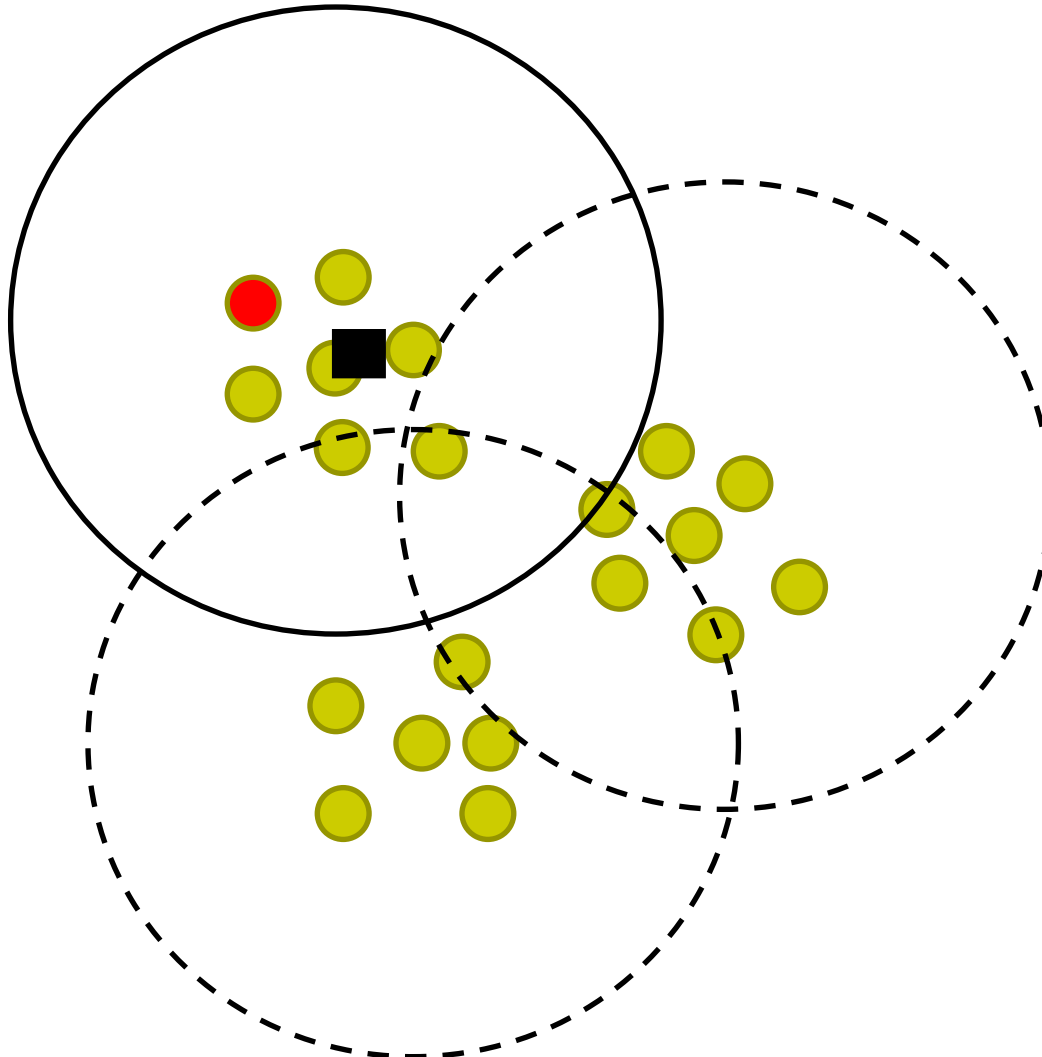
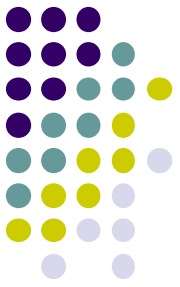
FOREL

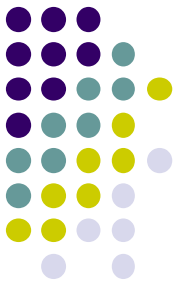


FOREL



FOREL



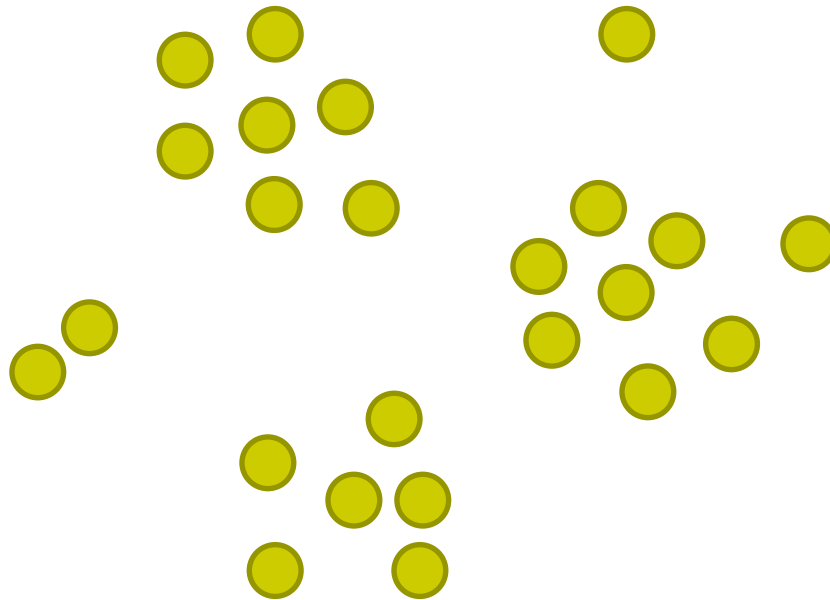
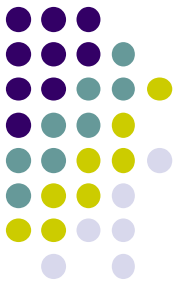


KOLAPS

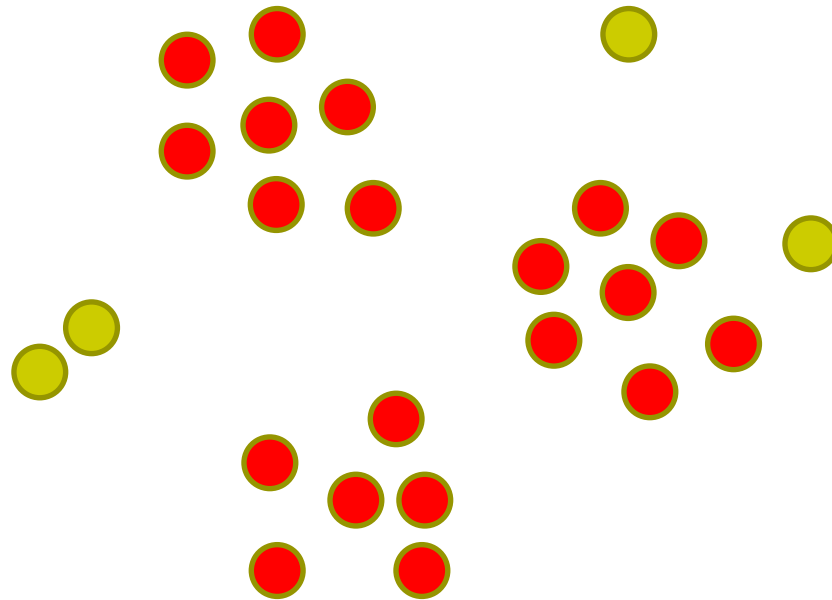
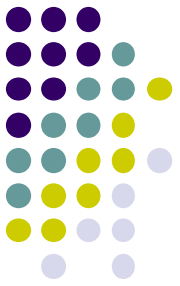
- $T=\{X_i\}, i=1, N$, d - мощность кластера
- Проводим генерацию по FOREL – $\{C_k\}, k=1, m$
- Находим все $|C_k|>d$ и заносим их в список L – кандидатов
- Кластеры кандидаты L_m :
 - Уменьшаем радиус от R до R_{min} с шагом dR
 - Если число теряемых точек за шаг увеличилось, то остановить сжатие

Критерий $\max \sum |L_m|$

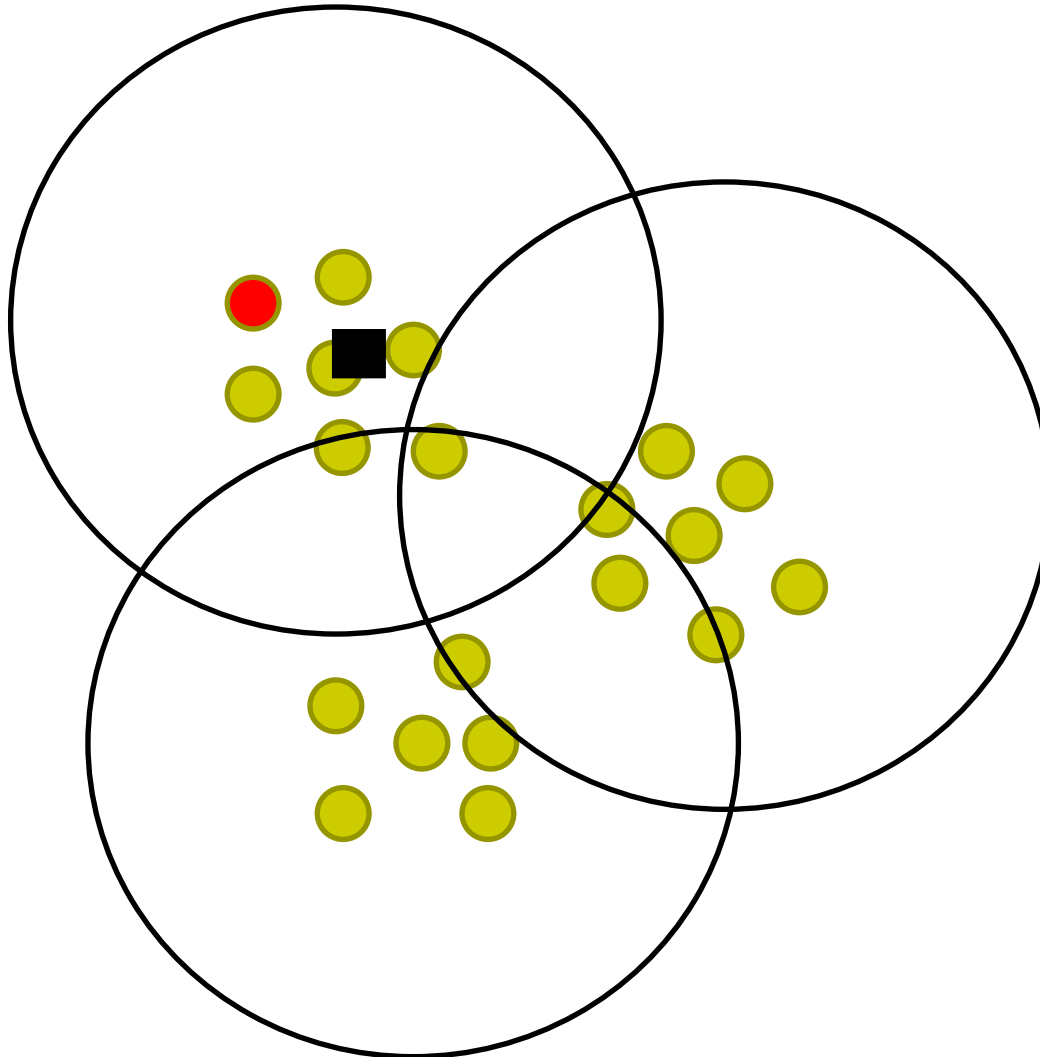
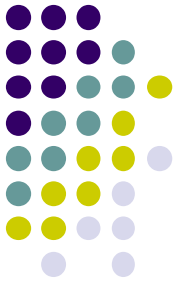
KOLAPS



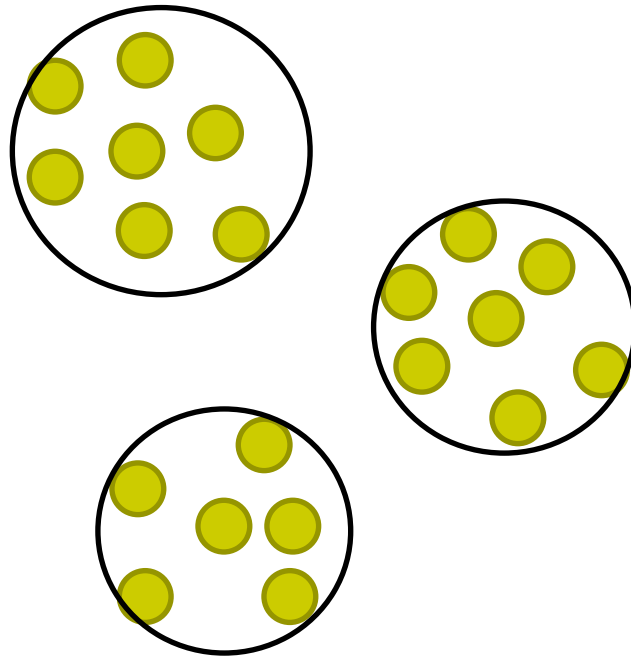
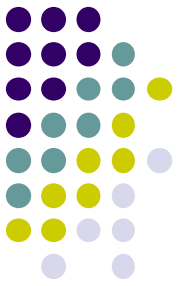
KOLAPS

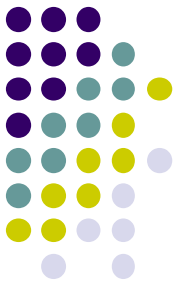


KOLAPS



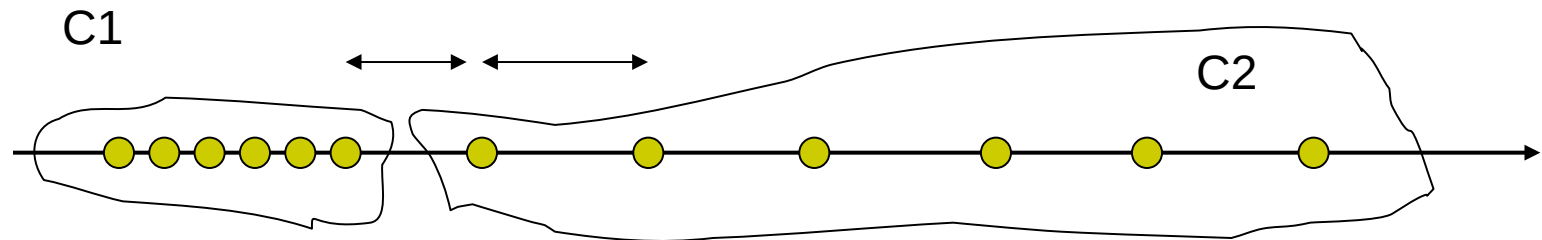
KOLAPS

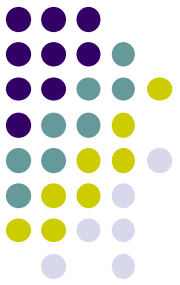




Гипотеза компактности

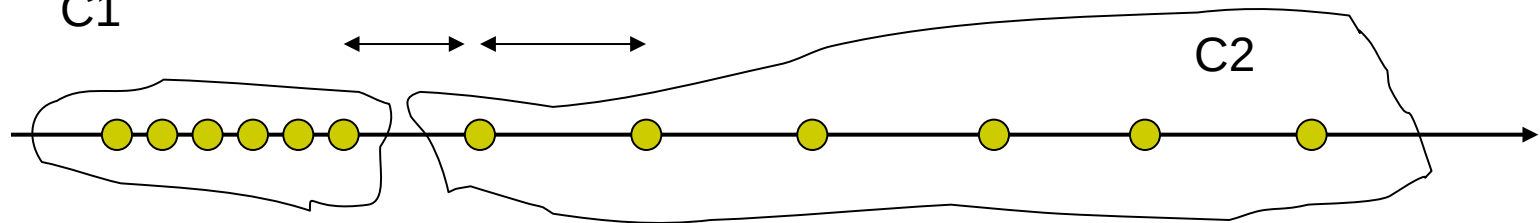
- Гипотеза λ -компактности
- Расстояние мало, но есть неоднородность.
- G – полный граф для $T=\{X_i\}$
- $A(a,b)$ – расстояние от точки a к b – длина ребра
- $D=\max(A(a,b))$





Гипотеза компактности

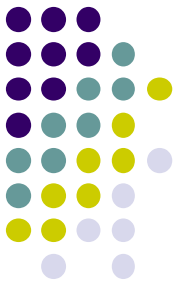
- $d = A(a, b) / D$
- Для ребер смежных с (a, b) определим $B_{\min} = \min(A(a, b))$
- $r^* = A(a, b) / B_{\min}$
- $R_{\max} = \max(r^*)$
- $R = r^* / R_{\max}$
- $\lambda = f(R, d) = R^2 d$



λ -KRAB



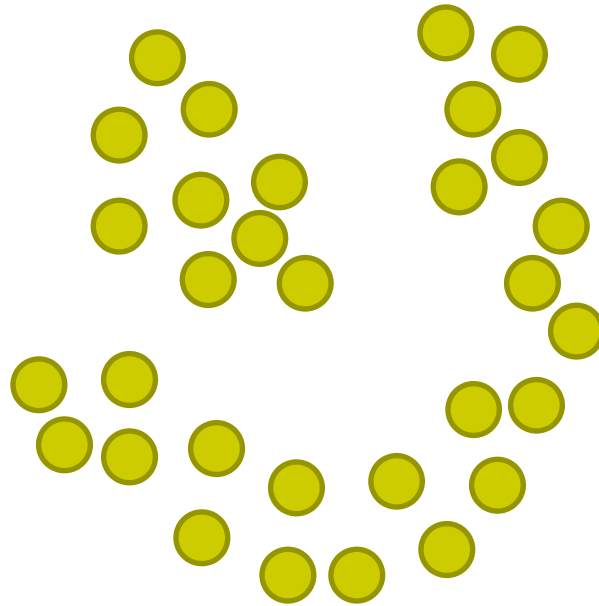
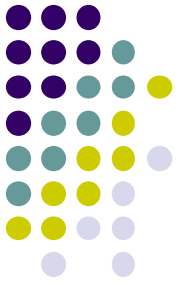
- Критерий равномерности таксонов (число объектов в кластерах приблизительно равно)
- $H = \sum_{i=1,k} (m_i/M)$, $i=1,k$, m_i – число объектов в кластере, M – общее число объектов
- Критерий алгоритма $\max (H^q R^s d^v)$ - q, s, v – параметры модели – степень влияния H, R и d ,
- Экспериментально - $\max (H^4 R^2 d)$



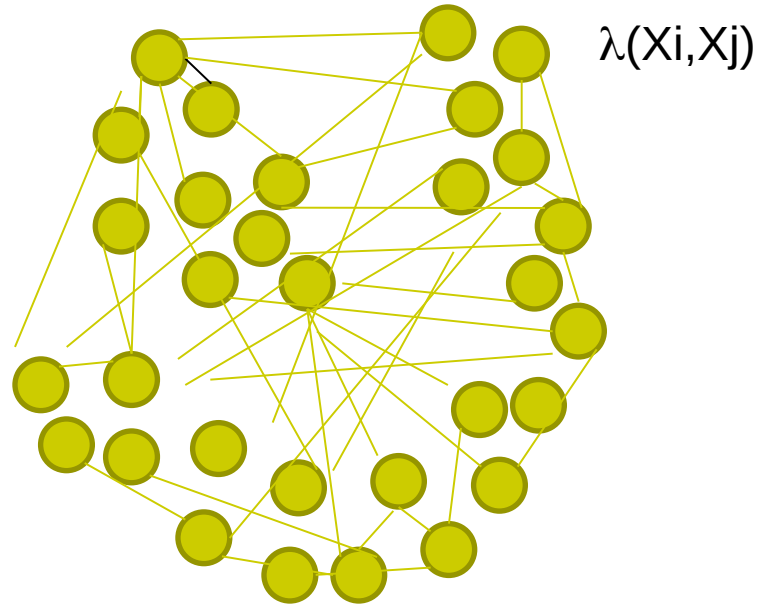
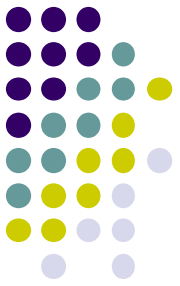
λ -KRAV – 2 кластера

1. Для всех $X_i, X_j \min(\lambda(X_i, X_j))$
2. Строим ребро $a(X_i, X_j)$ графа G .
3. Исключаем X_i, X_j из числа свободных
4. Если остались точки вне графа, то к 1, иначе к 5
5. G – кратчайший незамкнутый путь
6. Рассмотрели графы $G_1(a), G_2(a)$ с разрывом по ребру a и определили для каждого случая $f(a) = H(a)^4 R(a)^2 d(a)$ по 2-м кластерам.
7. Делаем разрыв в $a = \arg \max f(a)$

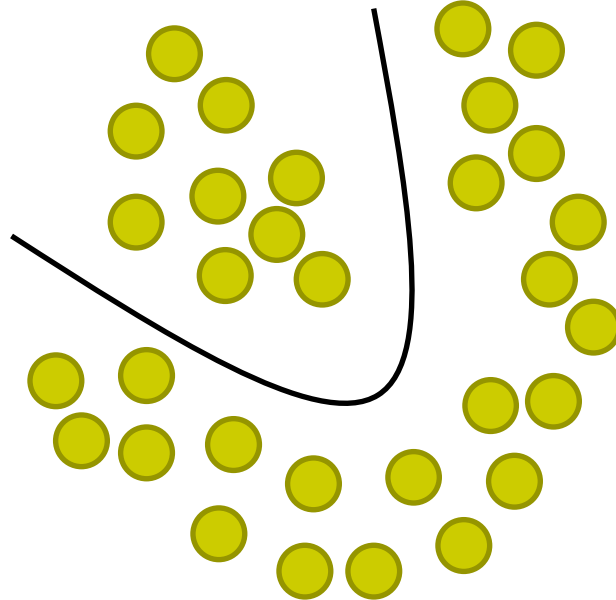
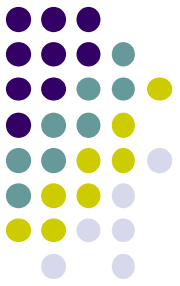
λ -KRAB



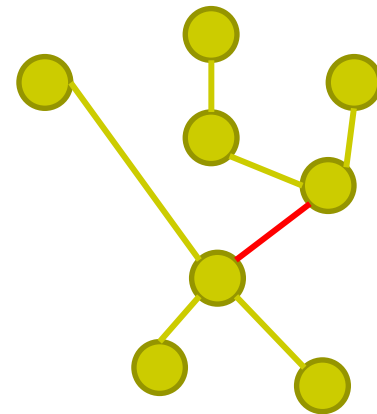
λ -KRAB



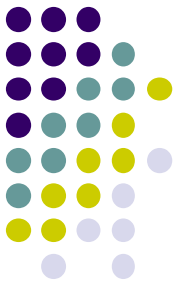
λ -KRAB



$\max f(a)$



Иерархическая кластеризация



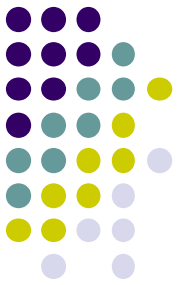
- Отношение эквивалентности

$$R_{\xi}(X_i, X_j) = \begin{cases} 0, \rho(X_i, X_j) < \xi \\ 1, \rho(X_i, X_j) \geq \xi \end{cases}$$

$\rho(X_i, X_j)$ - *мера близости*

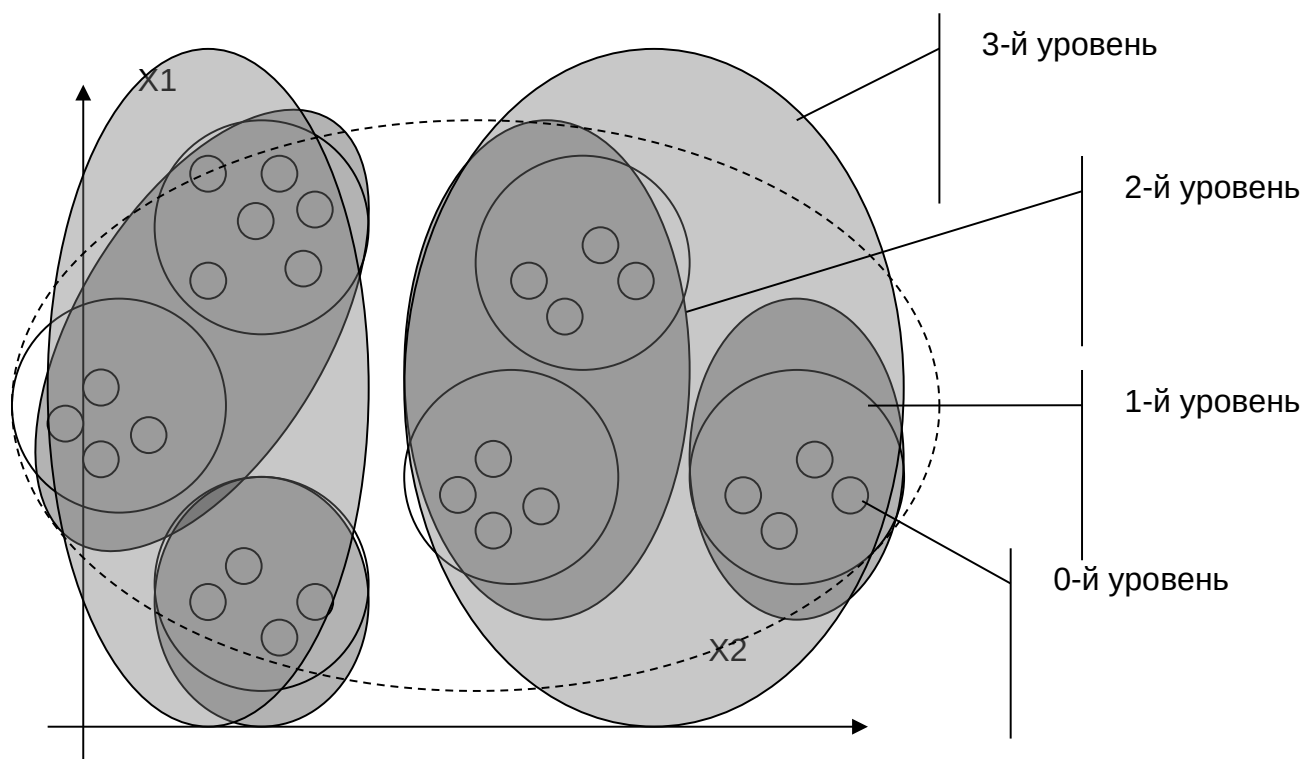
- $\xi=0$.
- Когда X_i, X_j *идентичны*, то $R_{\xi}(X_i, X_j) = 0$
- Работает эффективно для небольших объемов данных

Иерархическая кластеризация



- Каждому кластеру соответствует один уникальный объект.
- Для формирования следующего уровня – вычисляем центр кластеров текущего уровня и рассматривать эти центры в качестве входной информации следующего уровня. Для перехода на этот уровень следует увеличить порог $\xi + \delta$
- Продолжать процесс пока не останется один кластер.

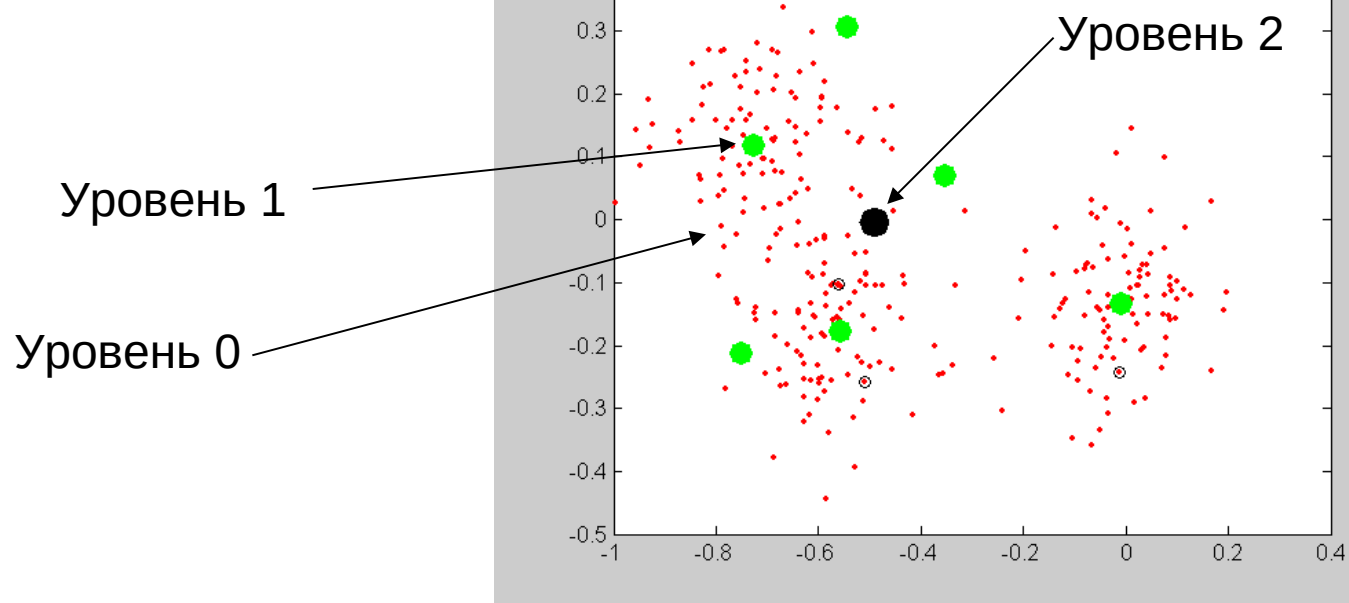
Иерархическая кластеризация



Пример



$N=3$;
 $n=100$;



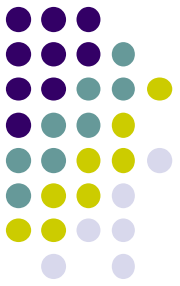
$P1 =$

-0.0110	-0.3528	-0.5584	-0.7277	-0.5441	-0.7518
-0.1335	0.0694	-0.1763	0.1191	0.3077	-0.2120

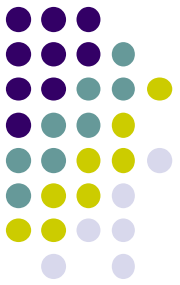
$P2 =$

-0.4910
-0.0043

Смысловые цели кластеризации

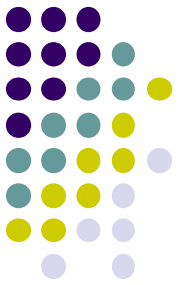


- Минимизировать изменчивость внутри кластеров,
- Максимизировать изменчивость между кластерами.



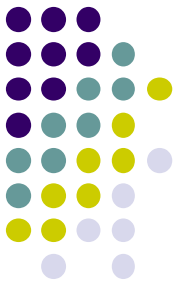
Новые методы

- Обработка сверхбольших объемов БД
- Требования:
 - Масштабируемость.
 - Работа в рамках оперативной памяти.
- методы кластеризации:
 - BIRCH,
 - CURE,
 - CHAMELEON,
 - ROCK



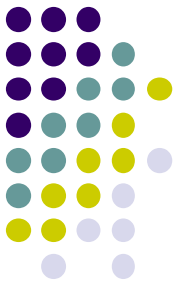
BIRCH

- Balanced Iterative Reducing and Clustering using Hierarchies - Тянь Зангом
- Алгоритм:
 - формируется предварительный набор кластеров.
 - к выявленным кластерам применяются другие алгоритмы кластеризации - пригодные для работы в оперативной памяти.



WaveCluster

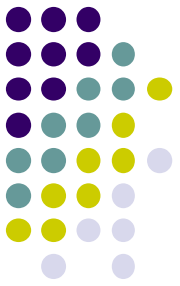
- Основан на волновых преобразованиях.
- Алгоритм:
 - данные обобщаются путем наложения на пространство данных многомерной решетки.
 - анализируются не отдельные точки, а обобщенные характеристики точек, попавших в одну ячейку решетки.
 - На последующих шагах для определения кластеров алгоритм применяет волновое преобразование к обобщенным данным.



особенности WaveCluster:

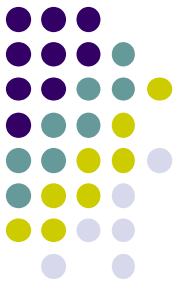
- сложность реализации;
- может обнаруживать кластеры произвольных форм;
- не чувствителен к шумам;
- применим только к данным низкой размерности.

Анализ результатов кластеризации.



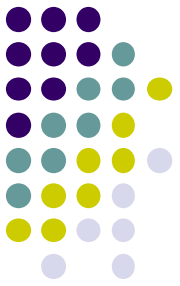
- не является ли полученное разбиение на кластеры случайным;
- является ли разбиение надежным и стабильным на подвыборках данных;
- существует ли взаимосвязь между результатами кластеризации и переменными, которые не участвовали в процессе кластеризации;
- можно ли интерпретировать полученные результаты кластеризации.

Процедуры проверки качества кластеризации:



- анализ результатов кластеризации, полученных на определенных выборках;
- кросс-проверка;
- проведение кластеризации при изменении порядка наблюдений в наборе данных;
- проведение кластеризации при удалении некоторых наблюдений;
- проведение кластеризации на небольших выборках.

Использование нескольких методов



- Отсутствие подобия не будет означать некорректность результатов,
- Присутствие похожих групп считается признаком качественной кластеризации.

Как сделать кластер анализ быстрее

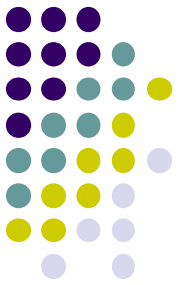


- Провести предобработку данных
 - Правильный выбор координат (оценка информативности)
 - Удаление выбросов (статистика и нормализация модели)
 - Редукция размерности
 - Факторный анализ (МЕТОД ГЛАВНЫХ КОМПОНЕНТ)
 - Многомерное шкалирование

Многомерное шкалирование



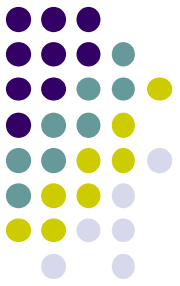
- целенаправленном преобразовании матриц сходства D , заранее сформированных на исходном множестве показателей
- Отображаем многомерные данные в пространство 2-х координат наибольшей значимости



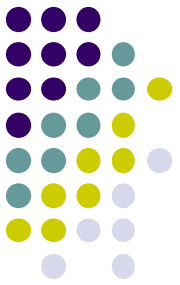
Этапы кластер анализа

- Использовать ли все наблюдения либо же исключить некоторые данные или выборки из набора данных.
- Выбор метрики и метода стандартизации исходных данных.
- Определение количества кластеров (для итеративного кластерного анализа).
- Определение метода кластеризации (правила объединения или связи)

Этапы кластер анализа

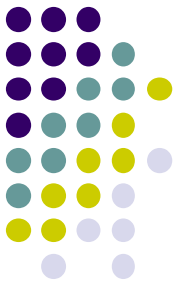


- Проведение кластеризации



Этапы кластер анализа

- Анализ результатов кластеризации.
 - не является ли полученное разбиение на кластеры случайным;
 - является ли разбиение надежным и стабильным на подвыборках данных;
 - существует ли взаимосвязь между результатами кластеризации и переменными, которые не участвовали в процессе кластеризации;
 - можно ли интерпретировать полученные результаты кластеризации



Этапы кластер анализа

- Оценка качества кластеризации и(возможно) возврат к предшествующим этапам
 - анализ результатов кластеризации, полученных на определенных выборках набора данных;
 - кросс-проверка;
 - проведение кластеризации при изменении порядка наблюдений в наборе данных;
 - проведение кластеризации при удалении некоторых наблюдений;
 - проведение кластеризации на небольших выборках.

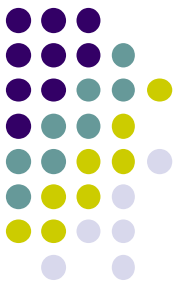
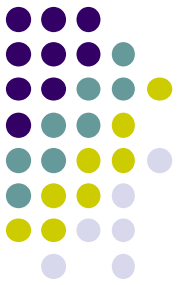


Таблица 5.2. Сравнение классификации и кластерзации

Характеристика	Классификация	Кластеризация
Контролируемость обучения	Контролируемое обучение	Неконтролируемое обучение
Стратегия	Обучение с учителем	Обучение без учителя
Наличие метки класса	Обучающее множество сопровождается меткой, указывающей класс, к которому относится наблюдение	Метки класса обучающего множества неизвестны
Основание для классификации	Новые данные классифицируются на основании обучающего множества	Дано множество данных с целью установления существования классов или кластеров данных



Этапы классификации

Кластеризация

изучение исходных данных на предмет наличия в них групп, классов и определение признаков, которые за это отвечают

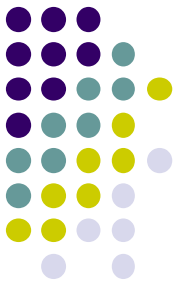


Построение модели

нахождение зависимости между значениями признаков объектов и принадлежность их к определенной группе

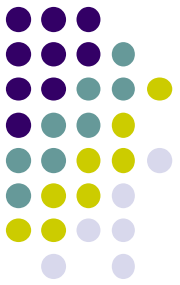


Классификация новых образцов
отождествление неизвестных образцов
с одним из известных классов



Возможные ситуации

1. **В начале ни одного класса не определено**
первым шагом в этом случае является
предварительный анализ данных на предмет
обнаружения потенциальных групп. В
зависимости от результата возможны
варианты:
 - **Имеется одна ярко выраженная группа**
 - **Имеется несколько ярко выраженных групп**Эти же варианты могут быть известны
априори



Возможные ситуации

2. Имеется одна ярко выраженная группа

В этом случае основная задача классификации найти и выделить типичную зависимость в данных для объектов, принадлежащих к одной группе и использовать ее для классификации новых объектов

3. Имеется несколько ярко выраженных групп

Необходимо использовать методы распознавания образов для выяснения принадлежности новых объектов к тому или иному классу. Задачу можно свести к предыдущей ситуации.