Обучение без учителя.

- набор данных, объекты $f(объект) = значение$ → дискр. ↘ непр.
- нет объектов, для которых знаем ответ.
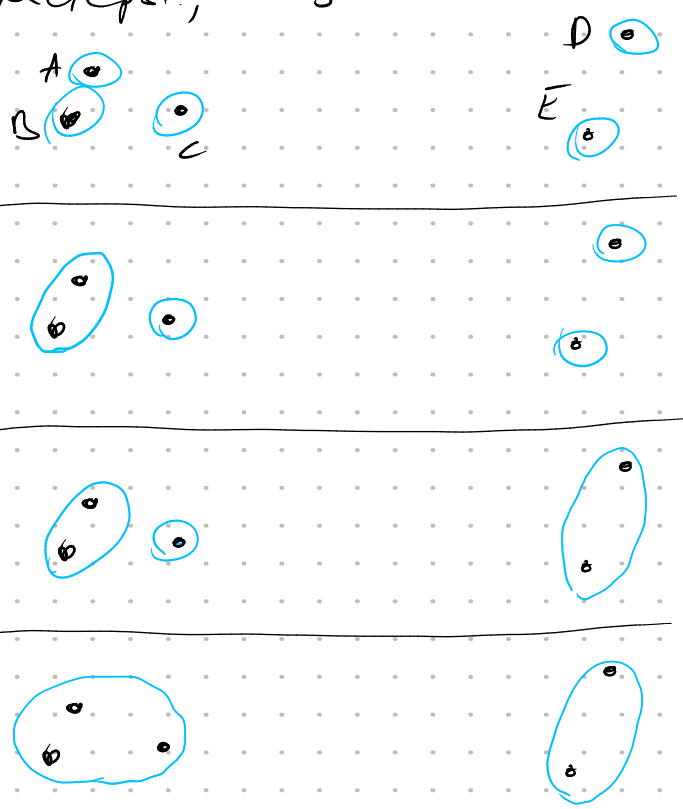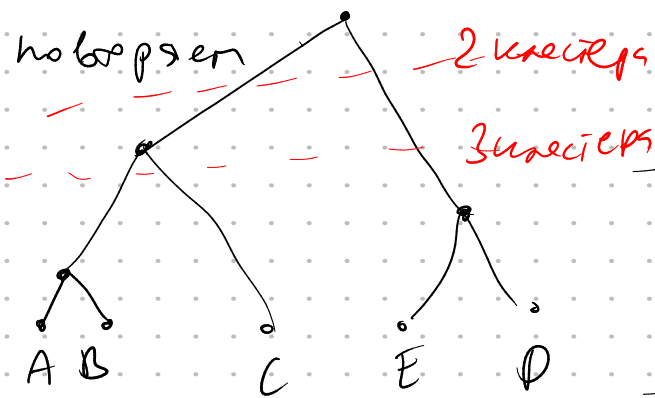
Если объекты = тексты, сами разбиваем на классы (темы)

Наглядно для объект = точки
кластеры - группы соседних точек.

# 1. Иерархическая кластеризация.
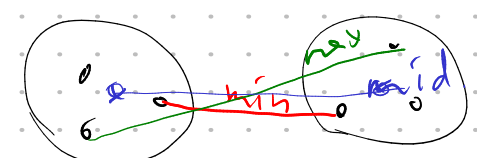
Начинаем с n точек, каждая точка = свой кластер
Находим самые близкие кластеры, объединяем
повторяем

— 2 кластера
— 3 кластера

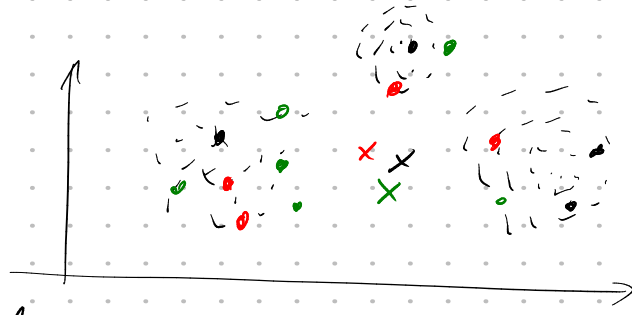когда остановили процесс,
тогда и получили оконч. -
точные кластеры

либо заранее знаете, сколько
хотите кластеров, либо в
процессе работы аккуратно
следите, не пора ли остановить.

Расстояние между кластерами

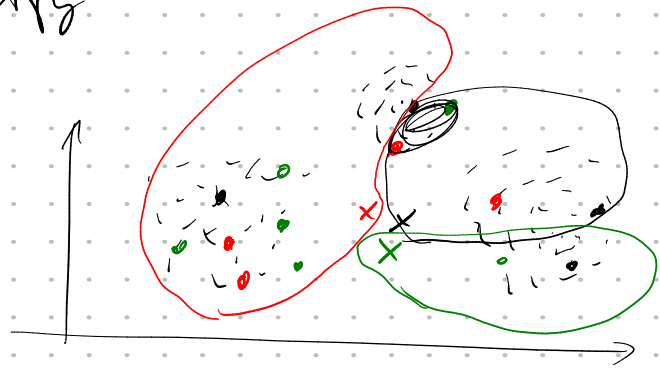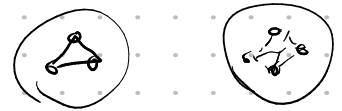Другие подходы - знаем, сколько кластеров хотим и сразу разбиваем на столько:

К-средних



сначала к случайных кластеров

нашли центры, и переназначаем кластер, исходя) точки и близ. центру

и т.д.



$I$ = Качество кластеризации: $\sum$ расстояний внутри кластеров $\rightarrow$ min



Еще есть спектральная кластеризация.

Но у нас же кластеризуют документы/предложения

$\uparrow$

для автоматич. реферирования

/ слова.

Нужно уметь превращать док/предл/слова в точки (в $\mathbb{R}^2$ $\mathbb{R}^{250}$ $\mathbb{R}^{100000}$).

Такие отображения будем называть: векторизация embedding

Для документов вариант:

Каждое значение в векторе — это одно слово.

сл1    сл2.    . . .        слN — все слова словаря.

text1
text2

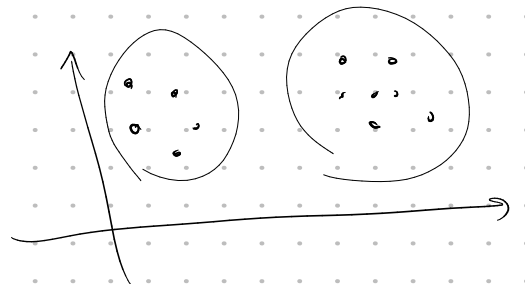1) на слов    10    5    0  0  0  0    1        — абсолютн.
2) частота    0.01  0.001  0  0  0  0    0.00.02    — относит.
3) tf-idf    = частота слова  ·  $\underbrace{\text{обр частот по słю}}$

$$\frac{1}{\log(doc\ freq + 1)}$$

term
frequency

. . .

Еще методы уменьшения размерности



text1
text2
⋮
textN

10000 значений



text1
text2
⋮
textN

200 значений

или даже

text1
text2
textN

2 значения

$x_1$  $x_2$

Примеры задач.

1. Кластеризация документов.   новый ⇄ темы
   tf-idf →

   + нарисовать их на плоскости.
   приведя к размерности 2 методом
   PCA (primary components analysis)



кластеризовать можно    1. в исх пространстве
                    ✓  2. в пр-ве после PCA → 100
                        3. — " — → 2

— — — — — — — — — — — —

2. Кластеризация предложений из текстов двух
   авторов. Совпадёт ли результат кластеризации
   с реальным разбиением на авторов?

3. Документ. Кластеризуем предложения.
   Реферат документа — это набор предложений по
   одному из кластера.



4. Смыслы слов.
   Толковый словарь          коса
                             — длинные волосы
                             — инструмент
                             — ...

идти
= перемещаете ногами
= происходит
= красиво. сочетается
То |

можно определять возможные смыслы любой типики.

Корпус, где используется слово:

применение слова



‖‖‖‖ ‖‖ ‖ слова ‿‿‿ ‿‿‿ ‿‿‿

контекст. (окно)

Все применения одного слова можно кластеризовать по контекстам. Кластер = смысл.

Как векторизовать контекст?    1) tf-idf.

2) word2vec

(позже)

какие слова использовать?

1) "идти" — многозначное слово с тонкими различными смыслами.

2) "Коса"   явные разные смыслы
   "замок"

3) "банан" + "шар"   ← как взять эти одно слово

50  11:20