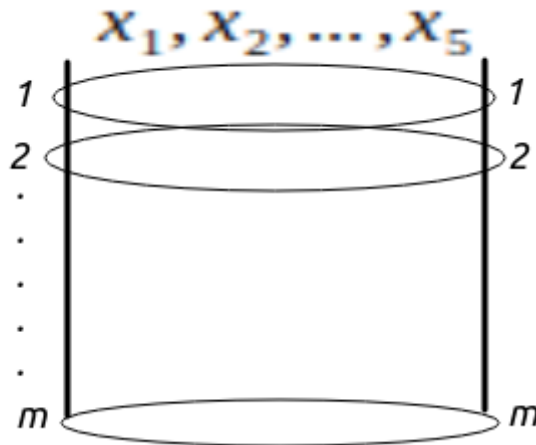


График параллельных координат

Пусть будет 5 признаков:

x_1, x_2, \dots, x_5

Каждый из них может состоять из сколько угодно элементов



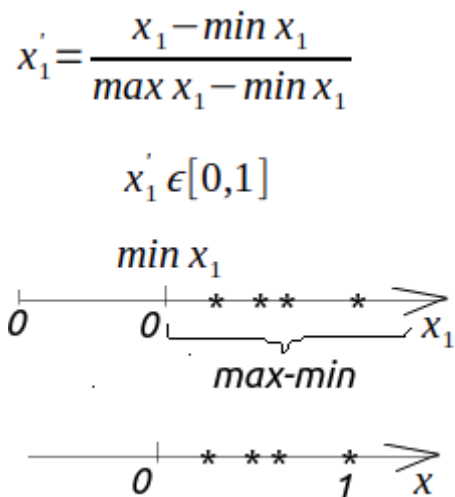
Тут у нас m измерений, чему оно равно нас не волнует

1. Отнормируем каждый признак на интервал $[0;1]$

2. Отнимим минимальное значение от каждого и поделим на размах

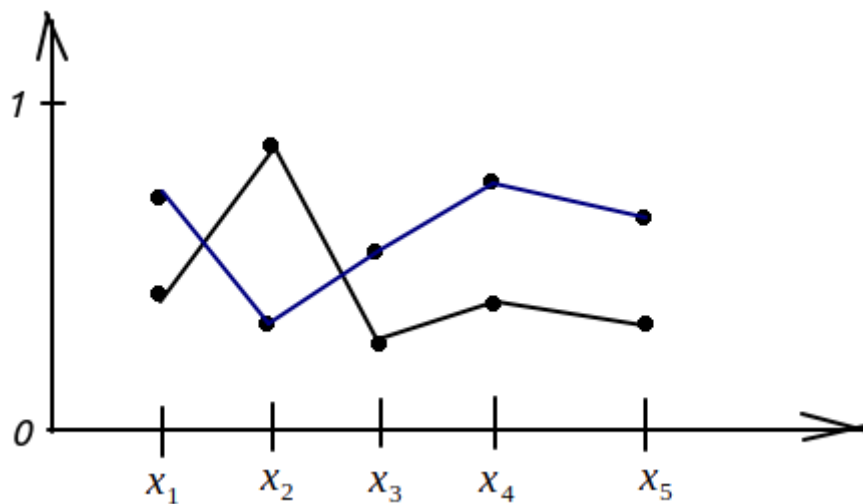
3. Первое, что мы делаем, это отнимаем минимумы. Минимум становится равным 0, а вот максимум = макс-мин

4. Когда поделим на размах все эти точки (то на прямой появится единица=макс)



Как строится график после нормировки?

Очень просто. По оси абсцисс - номер или имя признака, по оси ординат - нормированные значения от 0 до 1



И вот мы берем первый объект и у него $x_1=0.5$, $x_2=0.8$, $x_3=\text{масенький}$, x_4 и x_5 где-то равны. Соединяем эти точки прямыми.

Переходим ко второму объекту. Снова наносим все точки и соединяем их.

Делаем так, пока все объекты мы не описали.

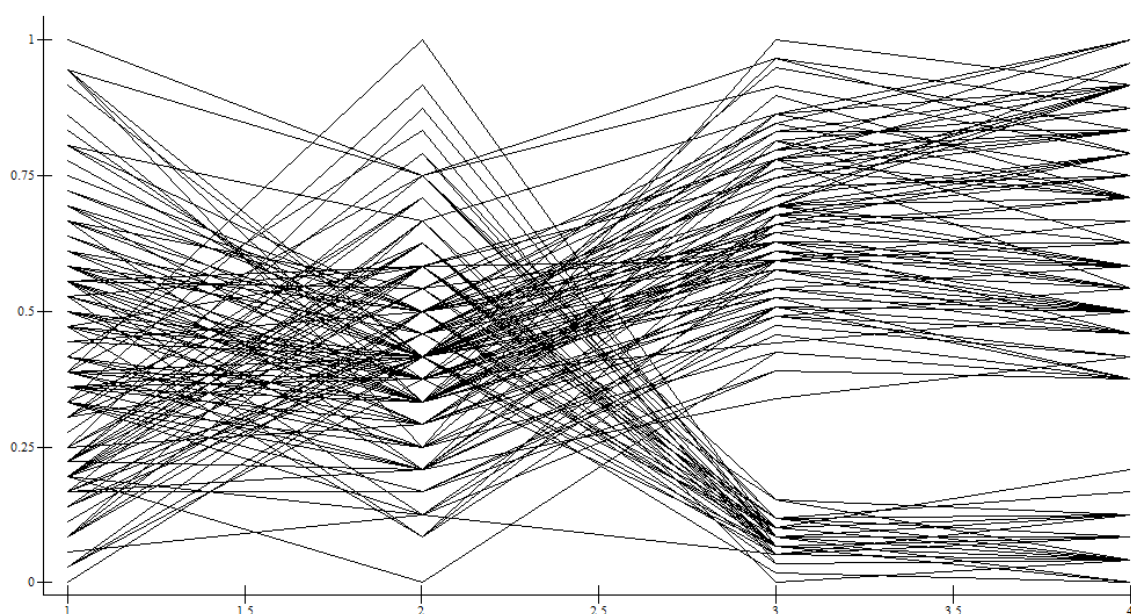
Реализация на APL:

Соединяем все ириски в матрицу `ir` и проверяем размерность. У нас 150 цветочков и 4 признака

```
pir←ir1,ir2,ir3
150 4
```

Получаем график ломаных

```
(150/1) parcoor ir
```



Что глядя на график можно сказать? Если наклонить голову вправо и представить как будет выглядеть гистограмма, она будет близка к равномерному распределению. От минимума до максимума мы видим, что все 150 цветочков перемешаны в кучу и нет никакой структуры соответственно. (Это по первому признаку)

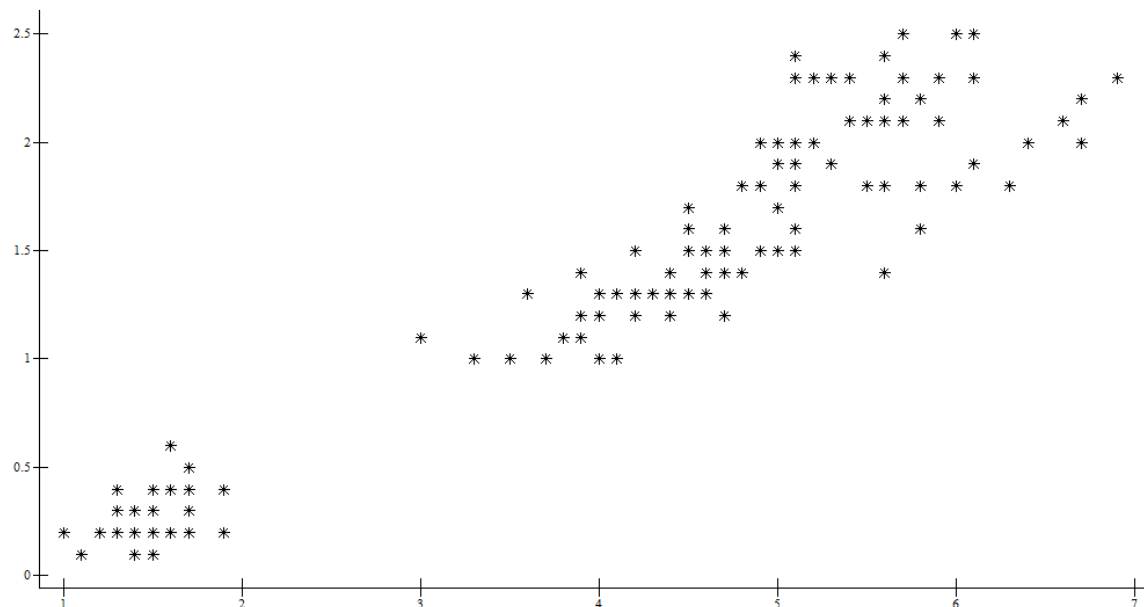
Глядя на второй признак, видим, что все цветки снова перемешаны.

А вот глядя на третий признак. ОПА!!! Видим две группы. Часть цветков имеют малое значение этого признака, а другая часть - большие значения этого признака.

Аналогично по признаку 4.

Выбираем интересные на наш взгляд признаки 3 и 4. Построим для них графичек.

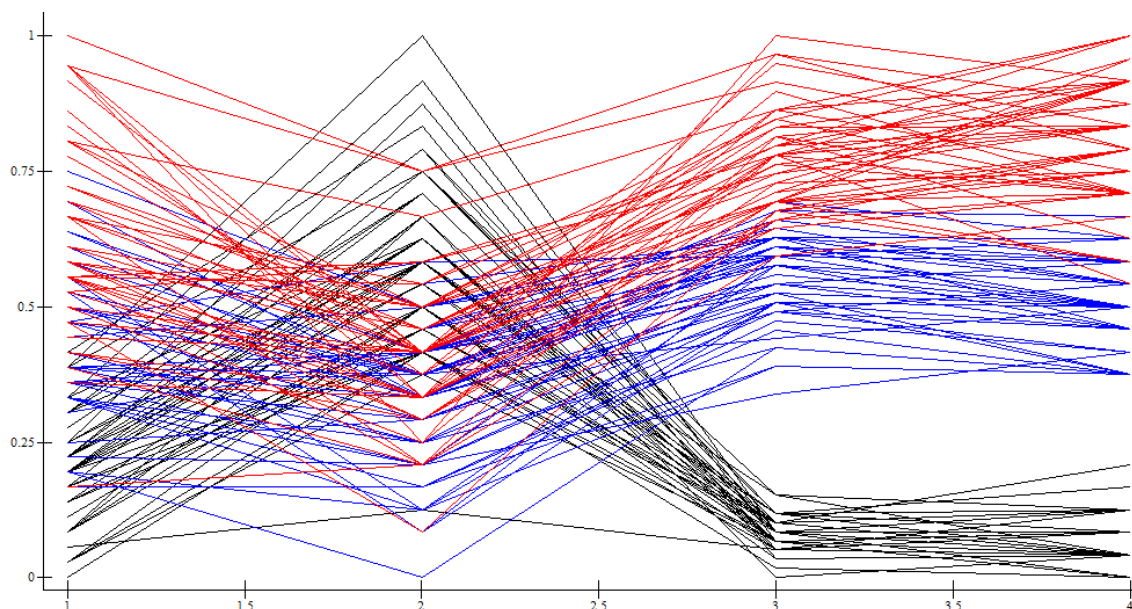
```
plot ir[,3 4]
```



Видим две группы, два кластера. Мы открыли то, чего не ожидали, что все цветки разделены на 2 группы. Мы не ожидали, что все цветки разобьются на две группы.

Теперь мы рассматриваем эти две кучки, Зовем Марию Ивановну и она нам говорит. Что одна группа принадлежит одному сорту, а во второй группе у нас смешаны два сорта. И Мария Ивановна нам их рисует. Но перед этим мы построим покрашенный график параллельных координат.

```
(50/13) parcoord ir
```



По первому признаку черный, красный, синий сорта перемешаны к чертовой матери.

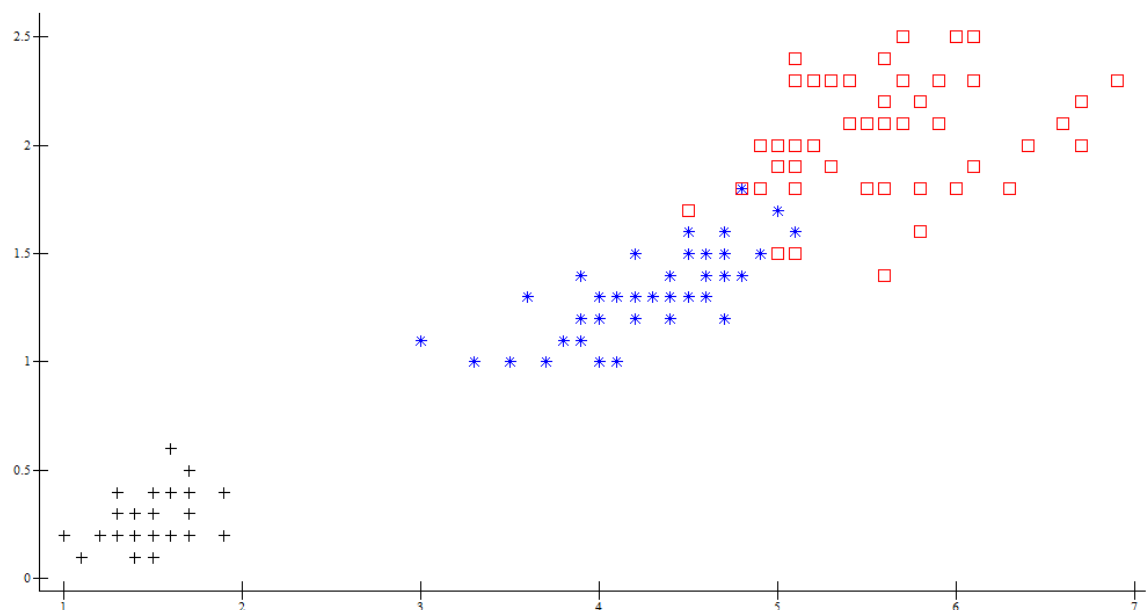
По второму признаку аналогично.

А по третьему и четвертому. Черненький (первый сорт) сильно отстоит от синего и красного сортов. Есть возможность поставить порог.

Мы видим, что не смотря на то, что чуть-чуть они перемешаны (4 красненьких попали к синеньким), т.е. мы не можем безошибочно разделить эти два класса, но можно сказать, что синенькие имеют меньшие значения признаков, чем красные.

Теперь можем постоить плоскость 3 и 4 признака и покрасить.

```
(50/13)plotc c[1]ir[,3 4]
```

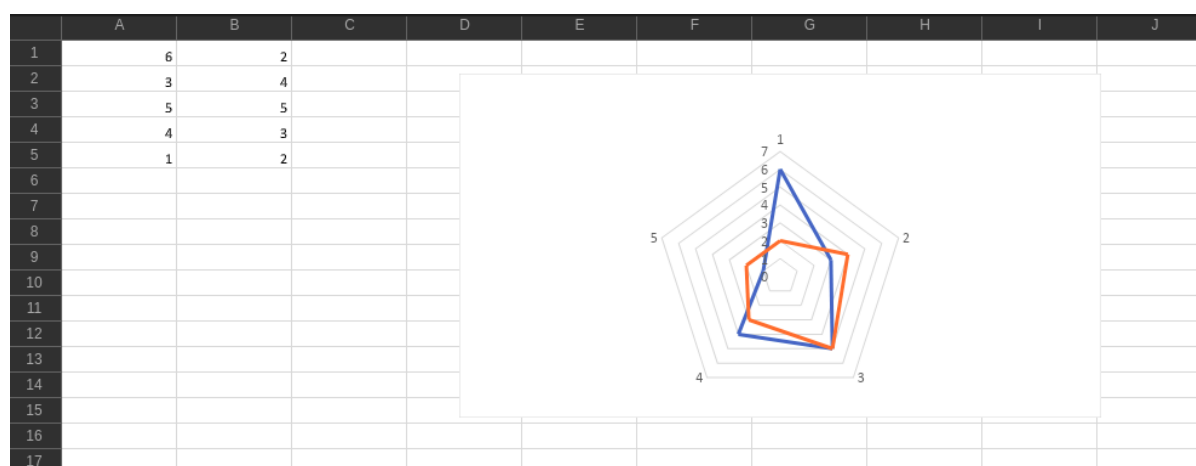


Этот образ отличается от первого, тем, что в принципе этому графику наплевать, сколько у нас признаков 4 или 400. Это очень мощный метод визуализации.

Еще один метод на котором мы подробно не будет останавливаться.

Radar plot (у него много названий)

Продемонстрируем в Экселе. Это тот же самый графичек, но который мы строим в полярных координатах.



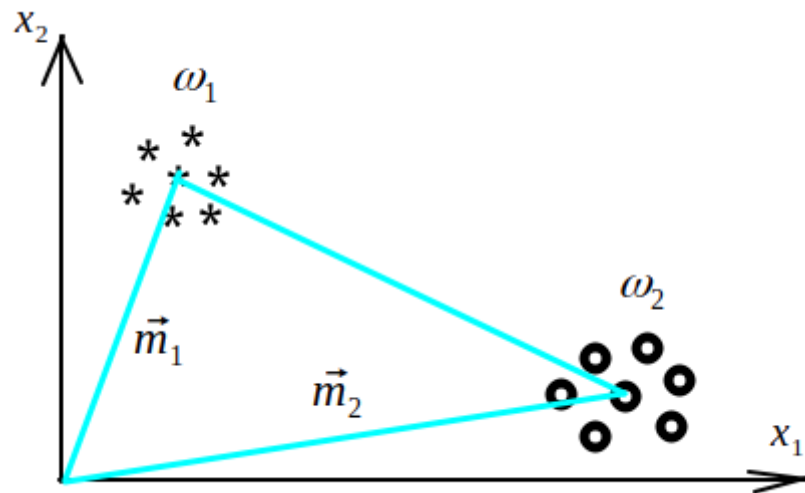
Видим, как у нас отличаются эти два столбца по своим значениям. И если их много, мы можем тоже выявить какие-то разные типы. Часть с острыми уголками, часть без острых уголков.

Вот это мы можем положить картинки и сказать: "Дети, разложите пожалуйста картинки на две кучки" и они прекрасно проведут этот кластерный анализ.

Все эти графики позволяют выявить какие-то интересные признаки и понизить размерность пространства, т.е. из большого числа X выбрать небольшое и в небольшой размерности мы можем замечательно смотреть.

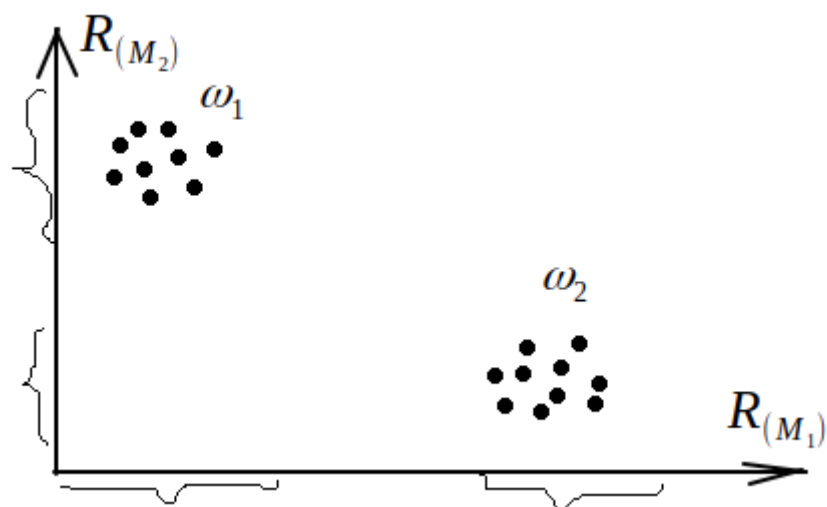
Теперь посмотрим простые два примера понижения размерности пространства:

1. **Расстояние до двух фиксированных точек.** это центры классов чаще всего фиксированных точек. Если известно, что есть два класса ω_1 и ω_2 , то две фиксированные точки. Так это будет выглядеть



M1-центр первого класса

M2-центр второго класса, строим следующий график



по оси абсцисс расстояние до R_1 , по ординат - до R_2

Сначала идут ω_1 , потом ω_2

У ω_1 расстояние до M_1 маленькое, а до M_2 - большое. Для ω_2 наоборот.

У нас изменятся деления по осям, а вид графика незначительно изменится.

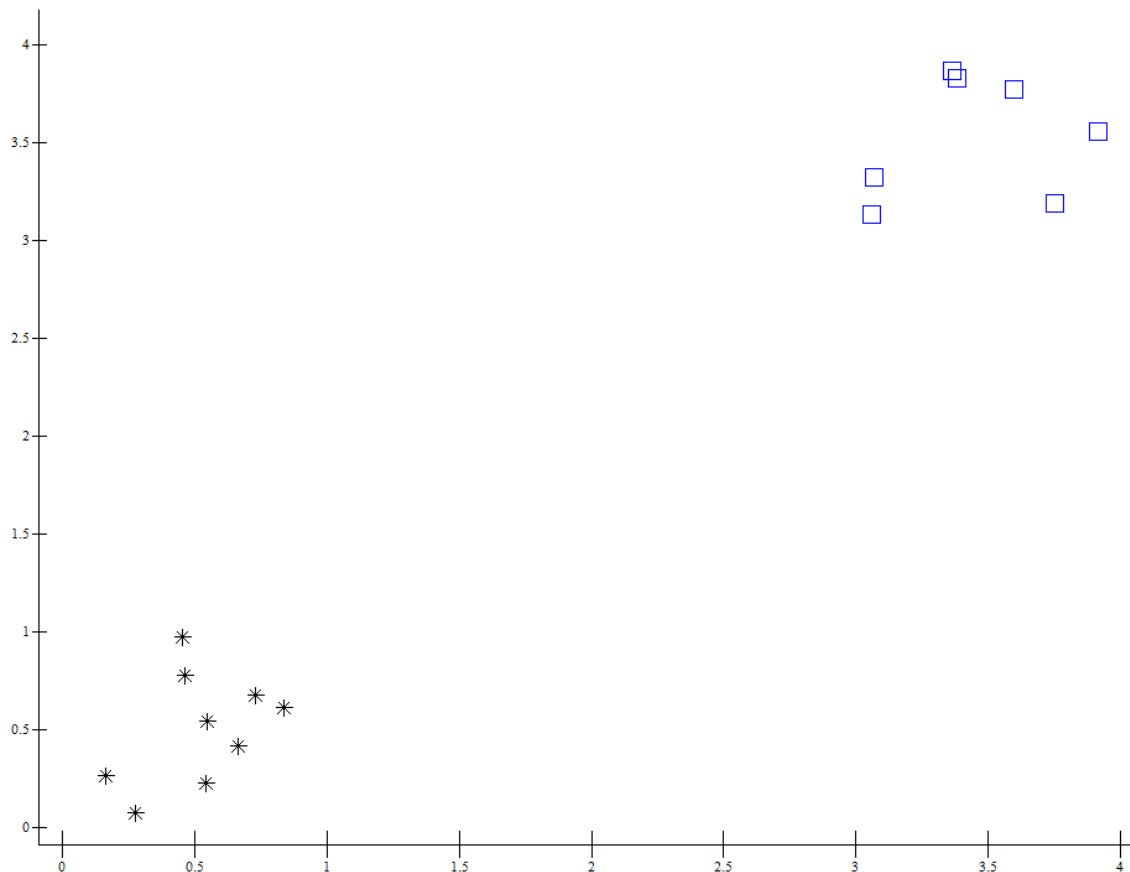
А зачем мы его строим? Его не надо было строить в двумерном пространстве потому что и так все видим, а вот в многомерном пространстве нам его надо построить.

Начнем с двумерного:

```
a←?9 2p0      Я А-первый класс(9 случайных точек от 0 до 1)
b←3+?7 2p0     Я В-второй класс(7 случайных точек от 0 до 1) и мы прибавим 3-
ку
```

Строим:

```
plot (c[1]a)(c[1]b)
```



Если вычислим M1. Напишем ф-цию ave и она будет суммировать матрицу по первому измерению и делить на число строк.

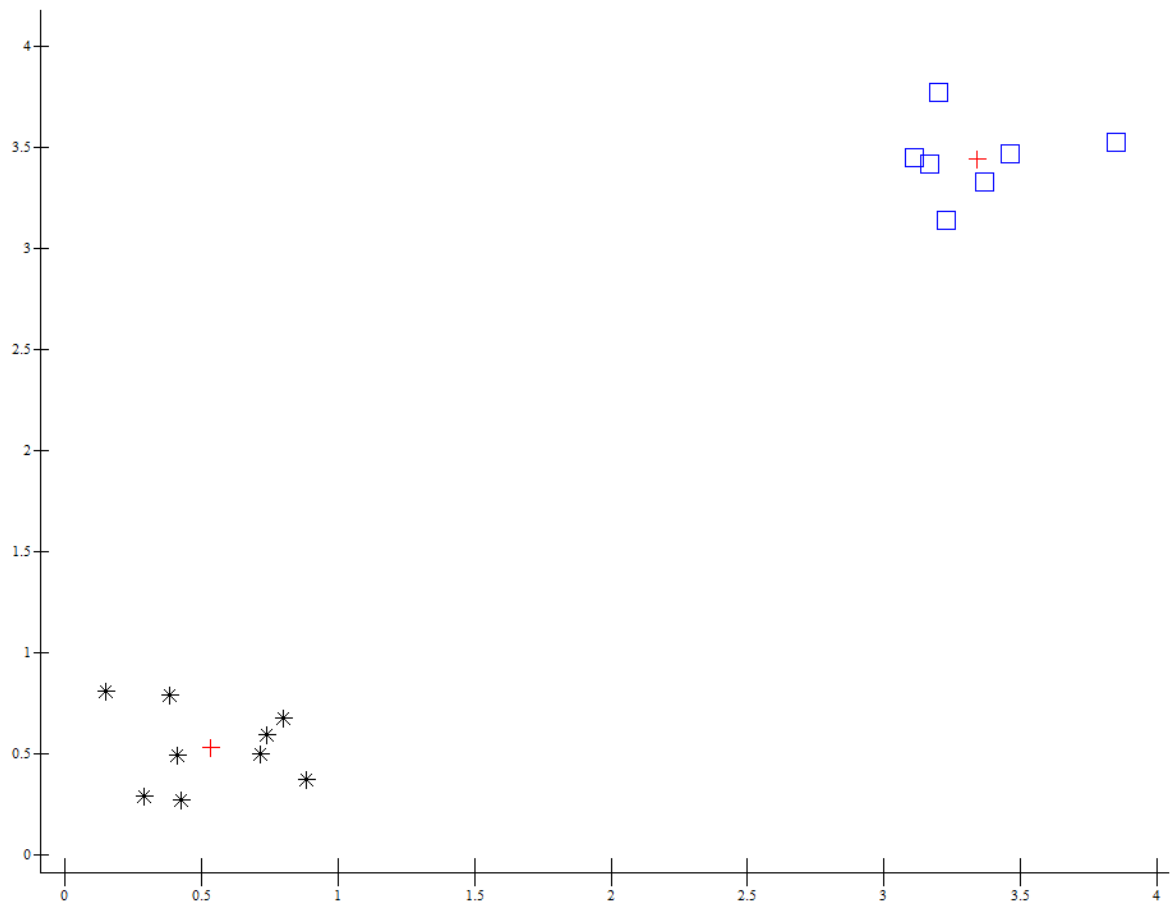
```
ave←{(+/w)÷#w}
```

Посчитали центры

```
m1←ave a
m2←ave b
```

Теперь мы можем их пометить на графике

```
color 'red'
marker m1  Я центр первого класса
marker m2  Я центр второго класса
```

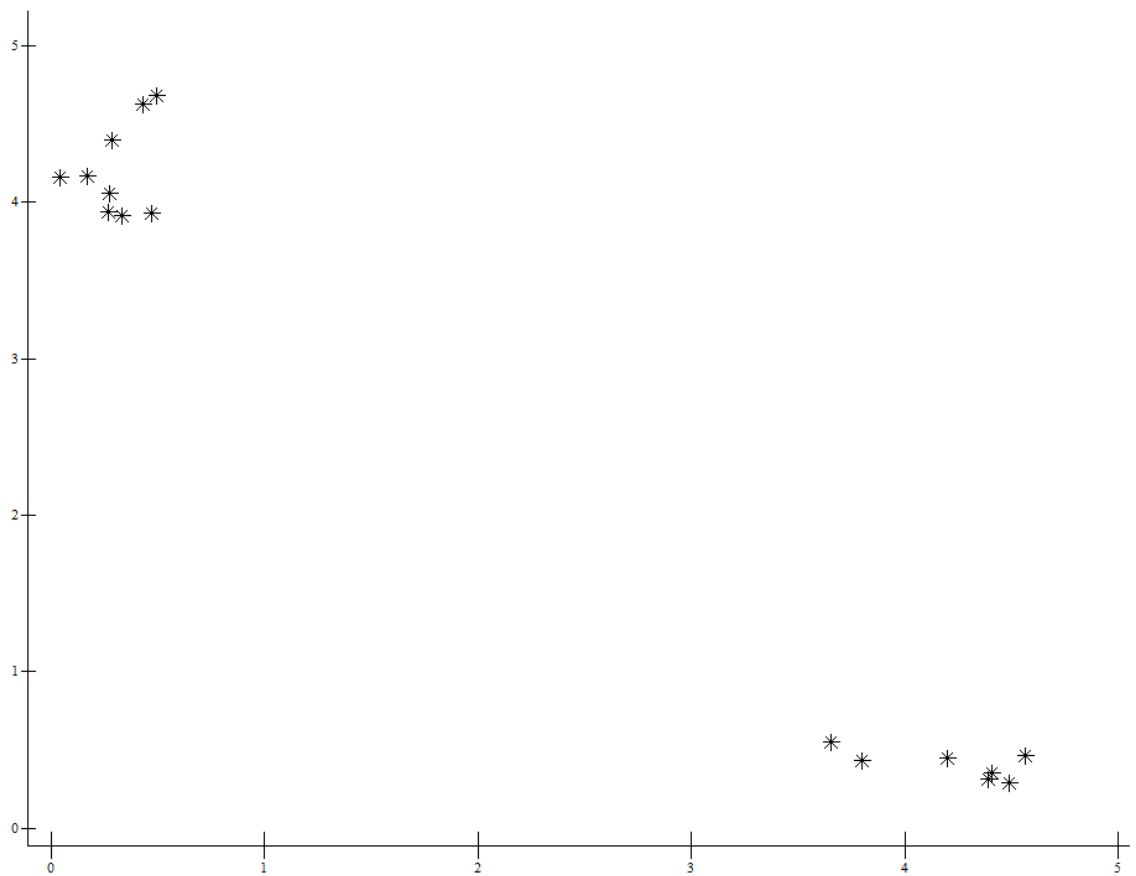


Теперь слепим a и b в одну матрицу, чтоб не мучиться и $R1$ - это расстояние от $M1$ до всех точек. $R2$ аналогично.

```
ab←a7b  
r1←(+/(ab-[2]m1)*2)*0.5  
r2←(+/(ab-[2]m2)*2)*0.5
```

Построим

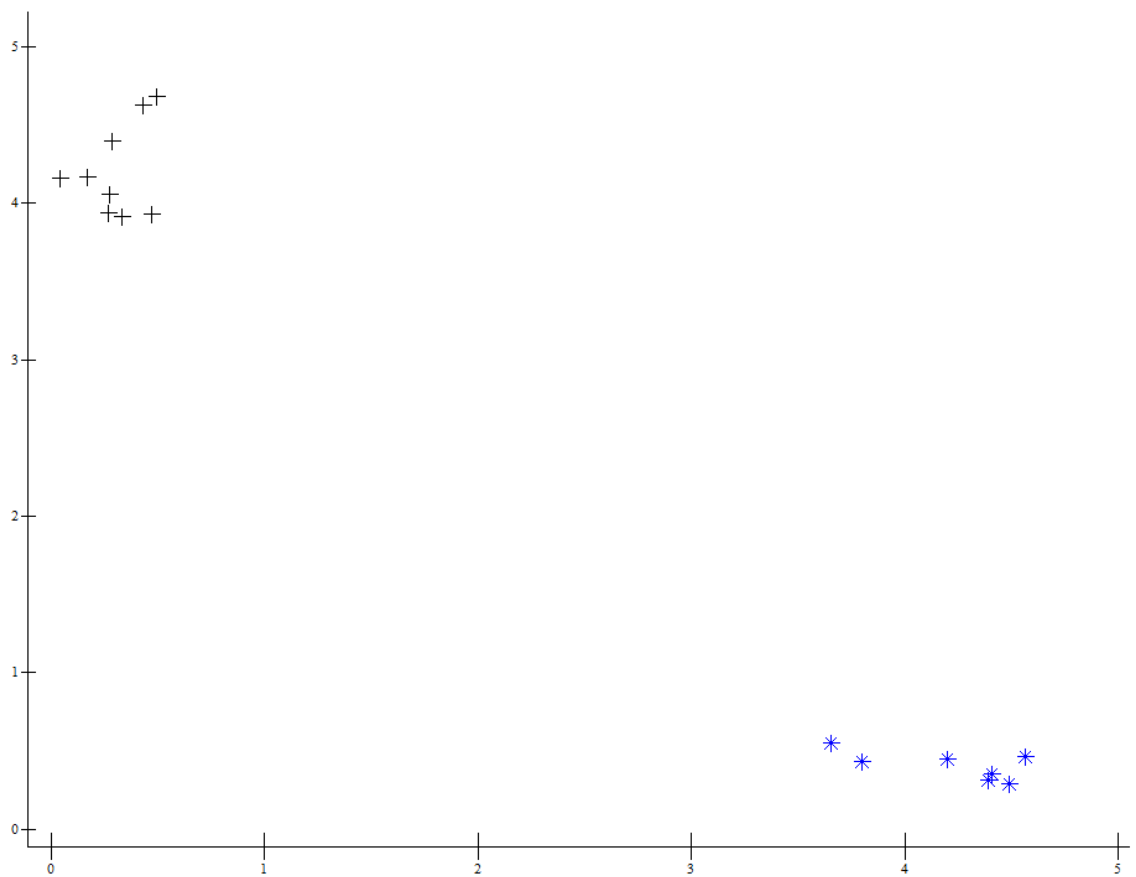
```
plot r1 r2
```



Распределение перевернулось, но ничего равным счетом не поменялось

Покрасим:

```
(9 7/1 2)plotc r1 r2
```



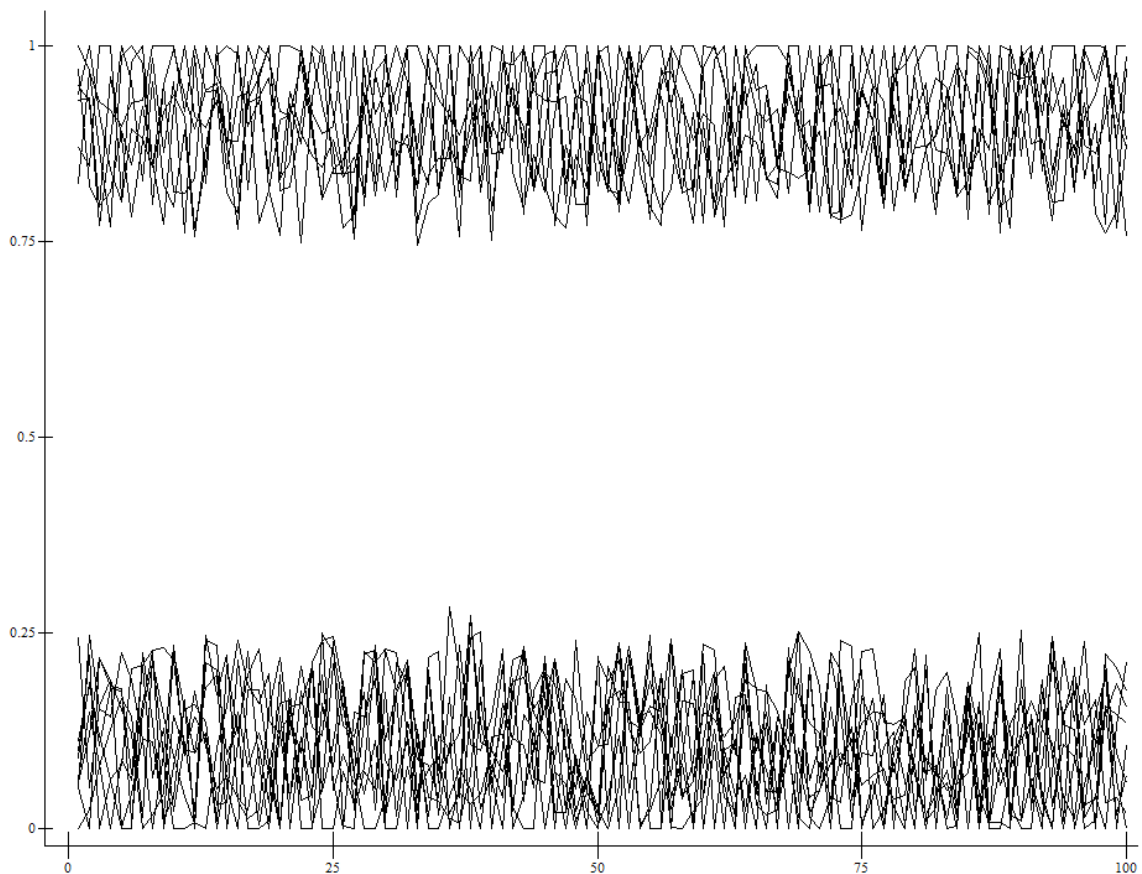
Когда у нас и так все было видно, то не имело смысла, что-то городить.

Тогда представим, что теперь у нас не двумерное, а **стомерное пространство**!

```
a←?9 100p0
b←3+?7 100p0
pa
9 100
pb
7 100
```

Никакого графика мы здесь не нарисуем, можем сделать параллельные координаты. Всего у нас 16 точек.

```
(16/1)parcoor a7b
```



Здесь трудно понять, что к чему относится.

Считаем опять M1 и M2 - центры:

```
m1←ave a
m2←ave b
pm1
100
pm2
100
```

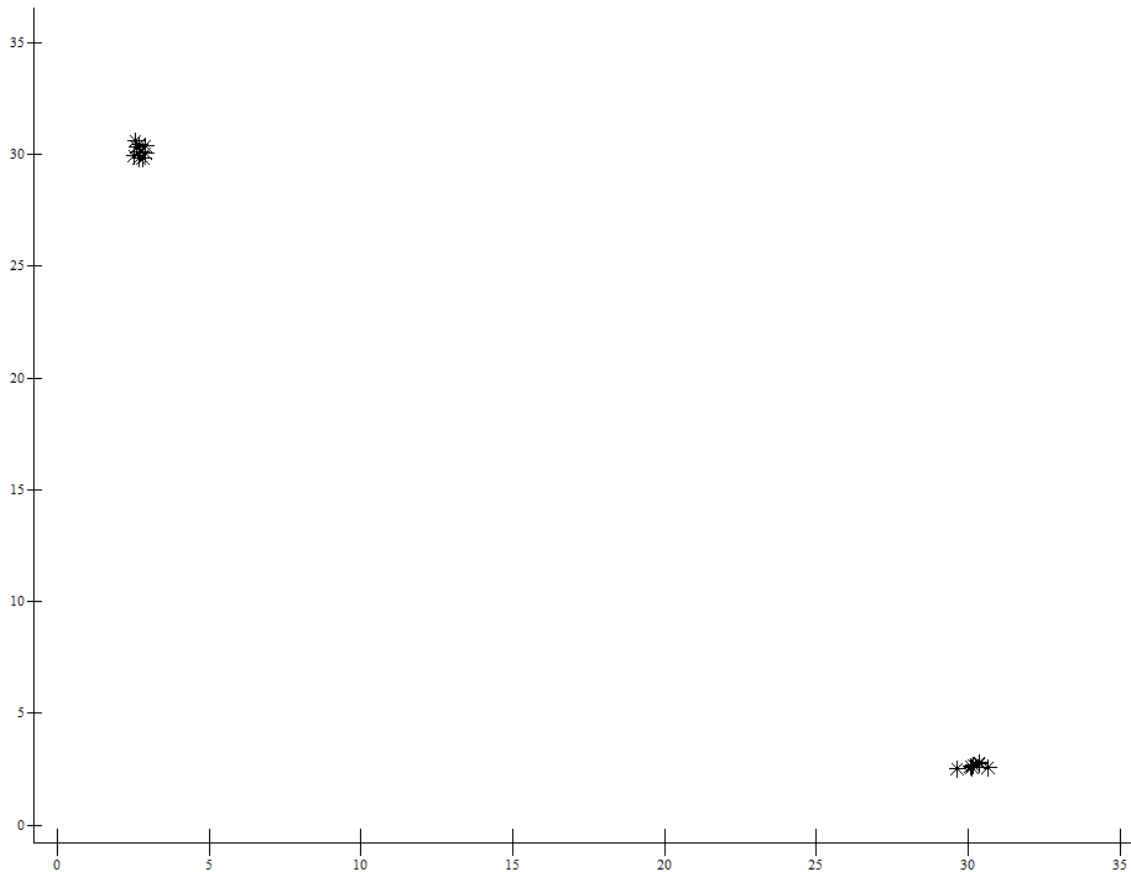
Они у нас теперь 100-мерные.

Делаем a и b. Считаем расстояние первого центра до всех точек

```
ab←a~b
r1←(+/(ab-[2]m1)*2)*0.5
r2←(+/(ab-[2]m2)*2)*0.5
```

И стоим графичек. Видим две группы на плоскости.

```
plot r1 r2
```

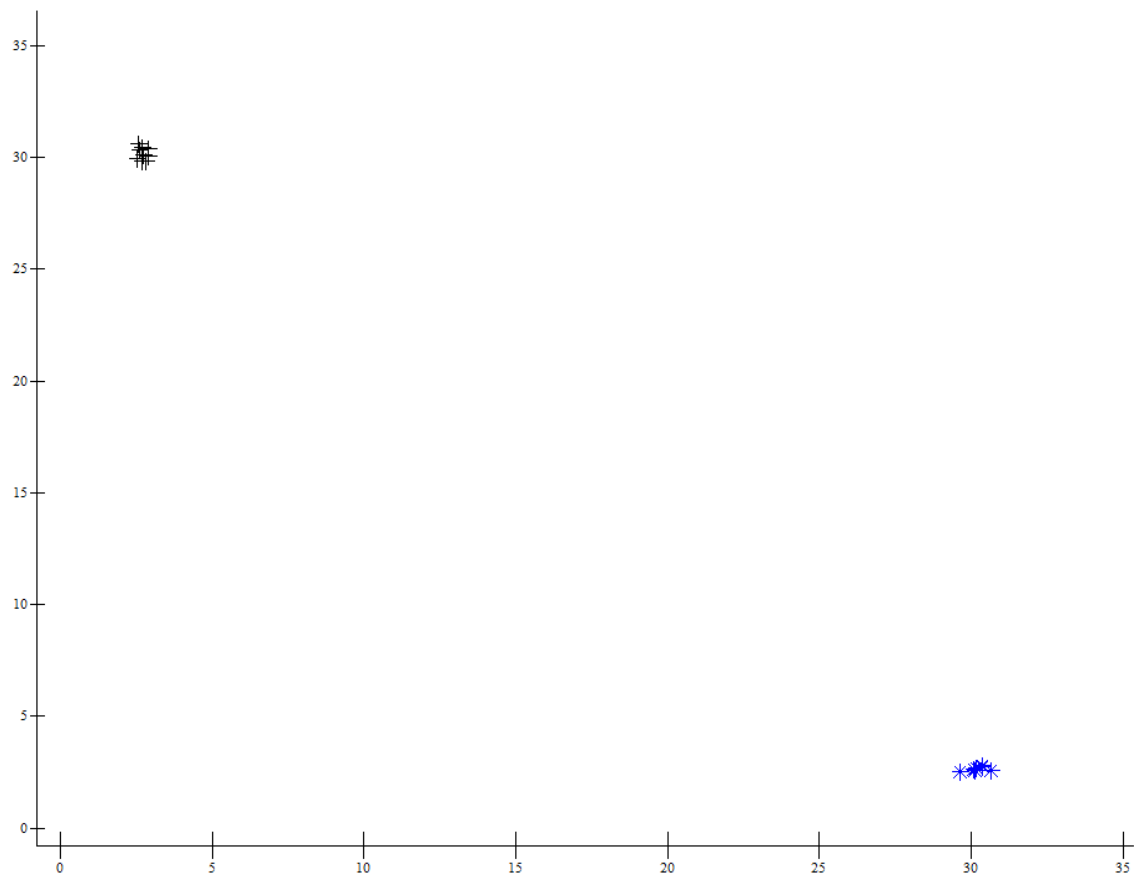


Можем взять 1 признак (кто меньше/больше 15)

```
+/r1<15
9
+/r1>15
7
```

Важно, что мы были в 100мерном пространстве и нифига не увидели, зная что первые точки относятся к одному классу, а другие ко второму. Посчитали расстояния до центров от всех точек.

```
(9 7/1 2)plotc r1 r2
```



Чтобы по 100 раз не набирать, сделаем ф-цию:

```
)ed m12proj

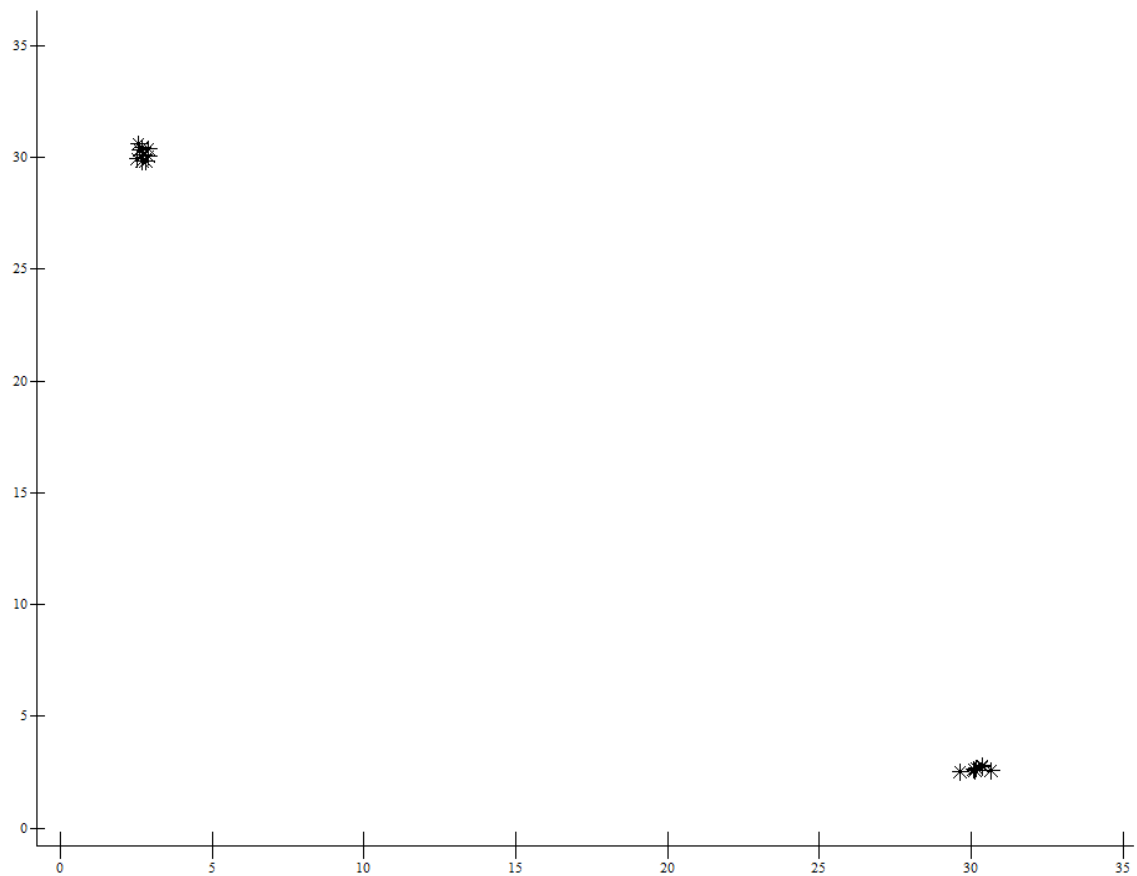
[0]   r←a m12proj b;ab;m1;m2;r1;r2  Я аргументы матрицы a и b. Результат -
вектор                                векторов расстояний.
[1]   m1←{(+÷ω)÷≠ω}a                Я считаем среднее по a
[2]   m2←{(+÷ω)÷≠ω}b                Я считаем среднее по b
[3]   ab←a⌣b                        Я делаем матрицу a и b, где у нас слепятся
обе
[4]   r1←(+/(ab-[2]m1)*2)*0.5        Я посчитаем расстояние
[5]   r2←(+/(ab-[2]m2)*2)*0.5
[6]   r←r1 r2                        Я в результат передадим r1 и r2

▽
```

Зелененьки-переменные глобальные(они нужны только на момент исполнения ф-ции), когда записываем в [0] они становятся локальными и стали беленькими.

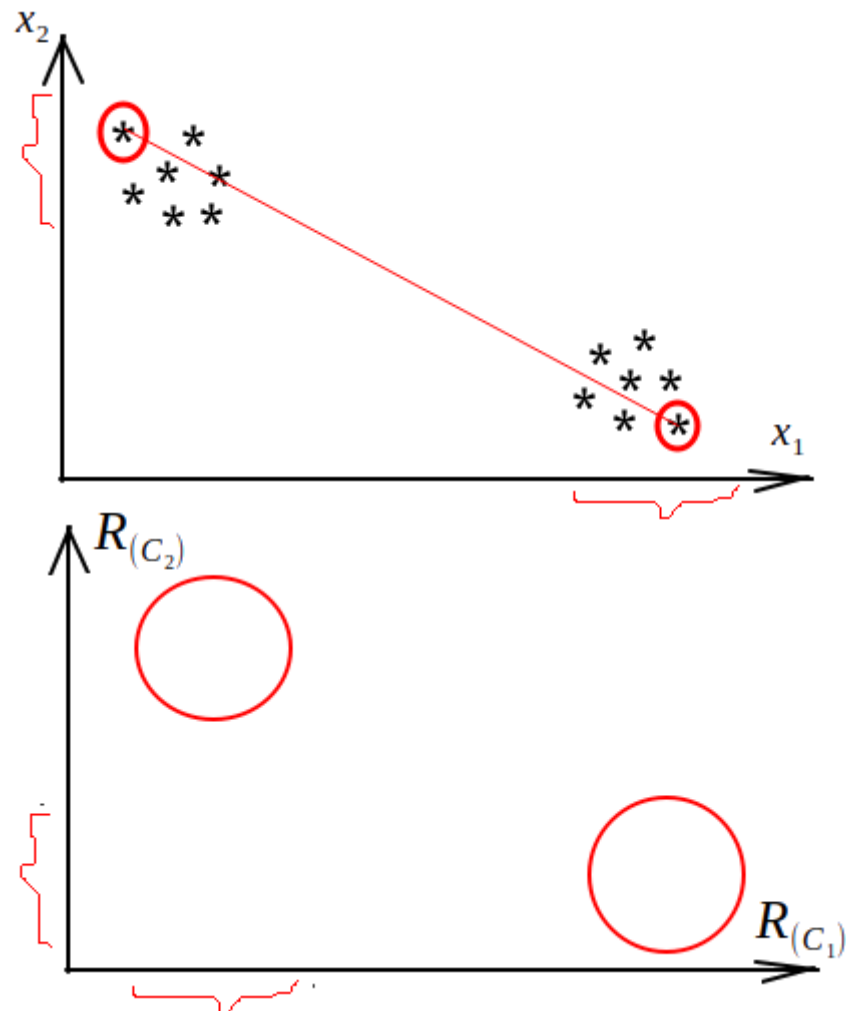
Построим:

```
plot a m12proj b
```



На выходе у нас будет вектор расстояний до первого класса и вектор расстояний до второго класса.

2. Если у нас тот случай, который очень часто встречается. Мы не знаем есть ли какие-то классы или их нет.



Какие две точки кажутся особенными? (зеленым обвели) Их особенность в том, что они наиболее удаленные друг от друга. Если посчитаем все возможные между всеми точками расстояния и найдем максимум, то он будет в этих двух точках.

Снова считаем расстояния. До C_1 - маленькие, а до C_2 - большие.

Нам наплевать в сколькомерном пространстве считать расстояния между точками

$$R_{xy} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Если $n=2$ - теорема Пифагора

Если $n=3$, берем ту же формулу, x и y - лежащий на столе и на полу шарик. Начало координат в углу комнаты. Один конец рулетки от шарика на столе, другой конец - шарик на полу и получим то же расстояние, что и по формуле.

Если $n=4$, то с рулеткой не залезем. И после 4-х нам уже наплевать.

Расстояние мы это посчитаем:

```
pab
16 100
dist←{+/(α-ω)*2)*.5}
```

есть точки

```

      1 2
1 2
      3 4
3 4

```

какое между ними расстояние? будет корень из 8(рисунок)

```

      8*.5
2.828427125

```

Имея *dist* ф-ции считаем расстояние между всеми *a* и *b*

Есть два вектора векторов, внешнее произведение - выполняет попарно(слева берется первый элемент и складывается со всеми эл-тами правого аргумента и т.д.)

```

      3 2 4 °.+ 5 3
8 6
7 5
9 7

```

Прелесть в том, что мы можем любой операнд давать (умножение, сложение, вычитание), если дадим два вектора векторов в нашу ф-цию *dist*-вычислятся расстояния между итым и житым и получаем результат в виде матрицы.

```

      3 2 4 °.× 5 3
15 9
10 6
20 12
      3 2 4 °.- 5 3
-2 0
-3 -1
-1 1
      (1 2)(3 4) °.dist (2 1)(4 2)
1.414213562 3
3.16227766 2.236067977

```

Превратим *a* и *b* в 16 векторов (каждый размерности 100)

```

      pab
16 100
      ab<-c[2]ab
      pab
16
      p`ab
100 100 100 100 100 100 100 100 100 100 100 100 100 100 100 100

```

Получили матрицу 16 на 16 парных расстояний, по диагонали матрицы 0

```

r<-ab °.dist ab

```

Сделаем вот такую *m*

```

m<-?3 3p9

```

Надо найти максимальный элемент в этой m (будет равно 9)

```
m=1/m
# логическая матрица, где 1=max
0 0 1
0 0 0
0 0 0
m
4 6 9
2 6 6
7 5 7
```

Координаты точки:

```
1m=1/m
1 3
```

В нашем случае r:

```
1r=1/r
3 15 15 3
```

Т.к. матрицы симметричные возьмем из них первые i=3, j=15

```
i j=1 1r=1/r
i
3
j
15
```

r1 - расстояния всех ab до первой точки:

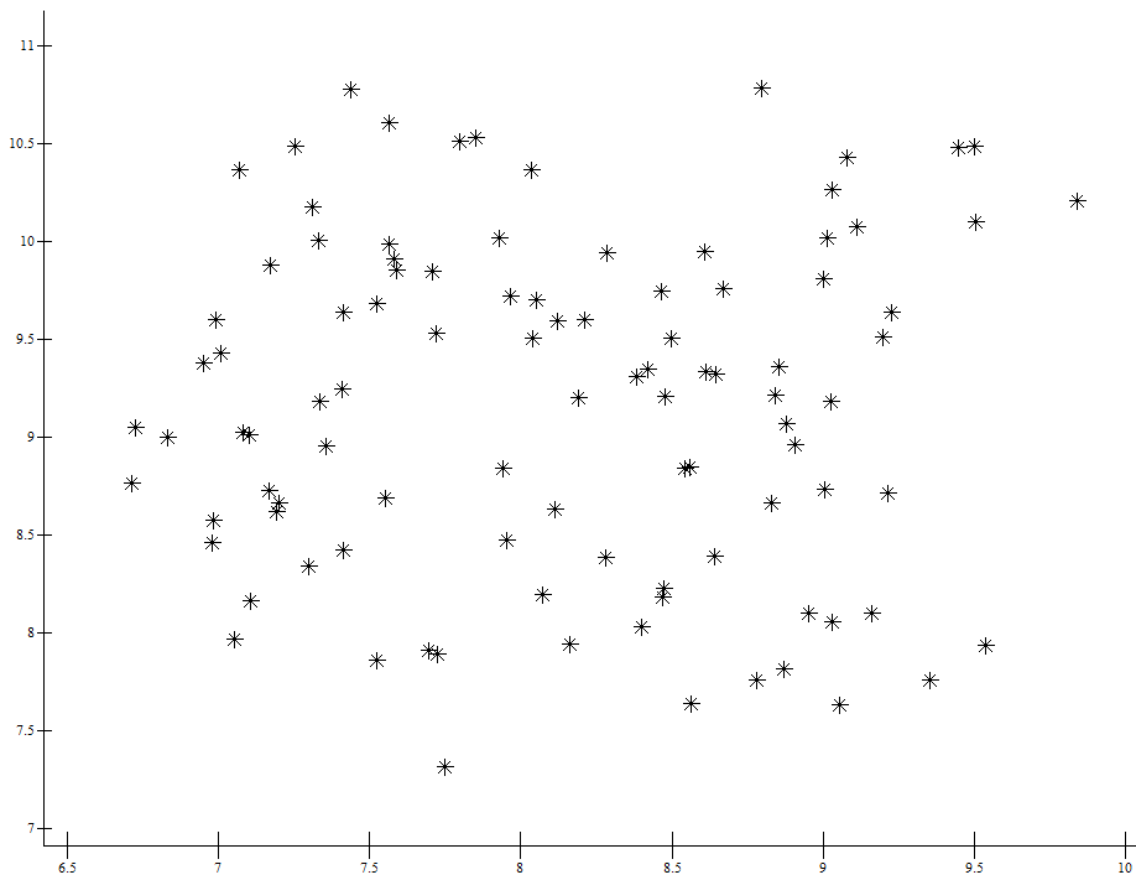
```
r1<-+(ab-ab[3])*2)*0.5
```

Так же считаем r2:

```
r2<-+(ab-ab[15])*2)*0.5
```

Строим:

```
plot ~r1 r2
```



Вот все расстояния от мервой до второй точки, т.к. опять мы понизили все наши расстояния.

Соберем все это хозяйство в ф-цию maxproj (расстояние через две максимальные равноудаленные точки) и проверим:

```

vr'maxproj'
∇ r←maxproj ab;i;j;r1;r2
[1] r←ab°. {+/(α-ω)*2}ab
[2] i j←↑1r=[/,r
[3] r1←(+/"(ab-ab[i])*2)*0.5
[4] r2←(+/"(ab-ab[j])*2)*0.5
[5] r←r1 r2
∇
r1 r2←maxproj ab
p`r1 r2
16 16
plot r1 r2
plot maxproj ab

```