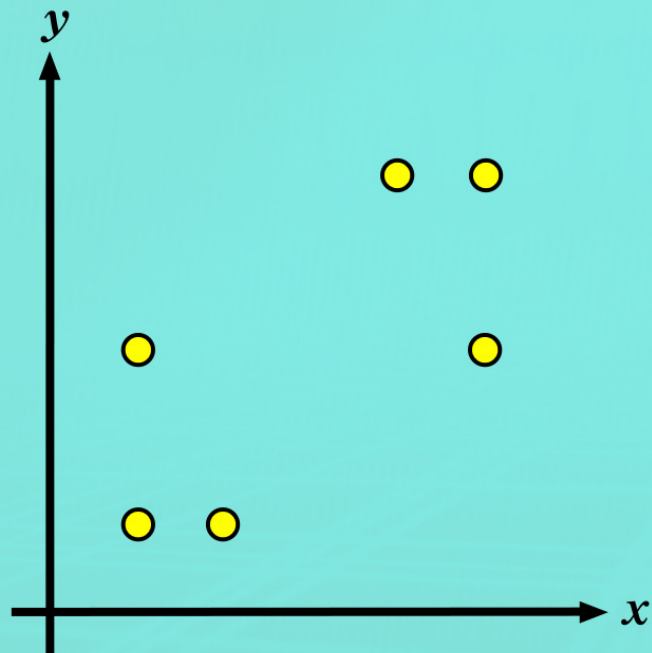


データの分析と

知識発見

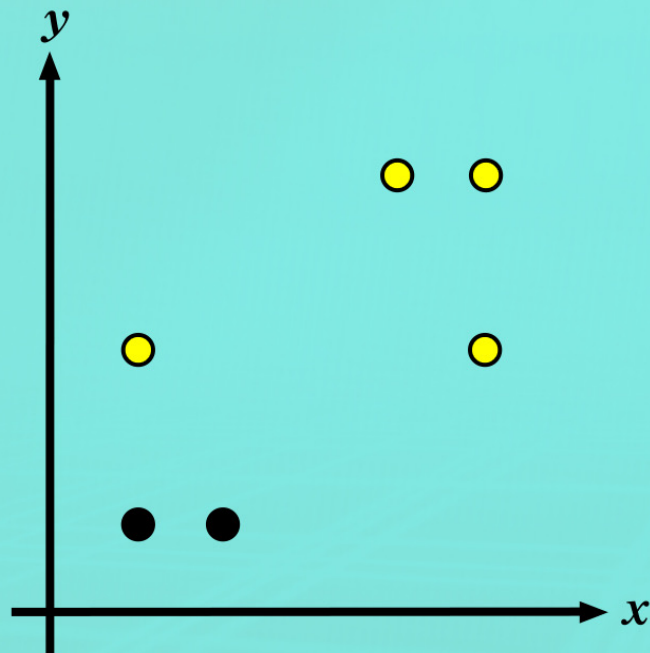
k 平均法

6個の点を2つのグループに分ける



k 平均法

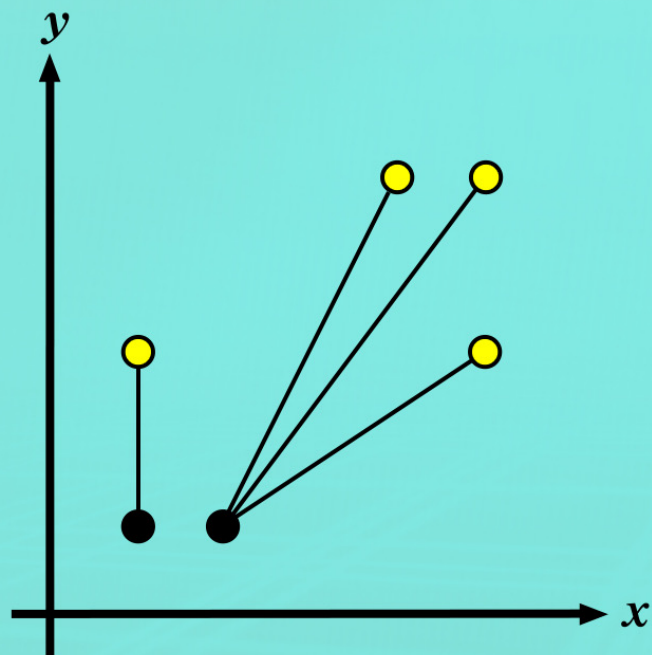
6個の点を2つのグループに分ける



- ◆ 2個の点を選び
その点を2つのグループの
代表点とする

k 平均法

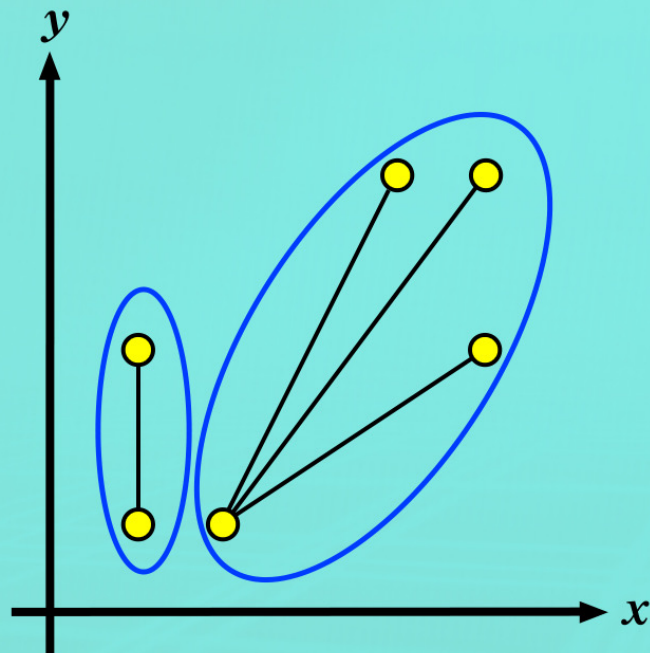
6個の点を2つのグループに分ける



◆ その2点との距離を計算する

k 平均法

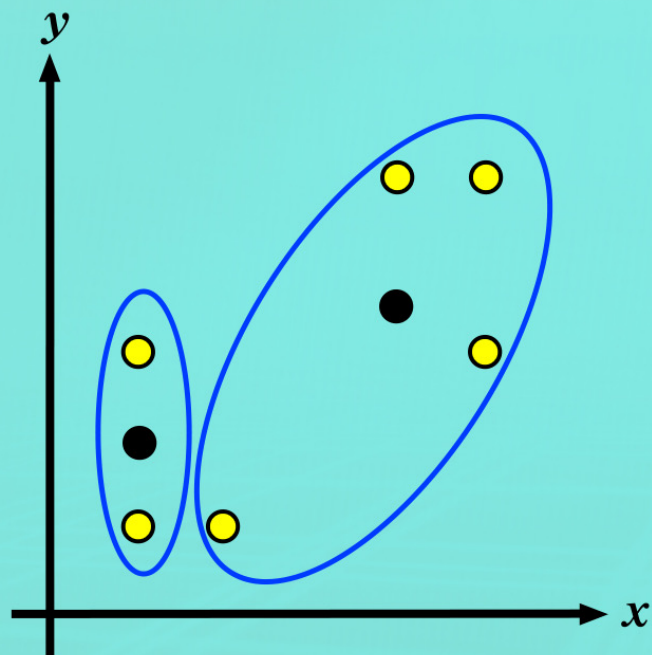
6個の点を2つのグループに分ける



- ◆ それぞれの点は
近い方の代表点のグループに
属することにする

k 平均法

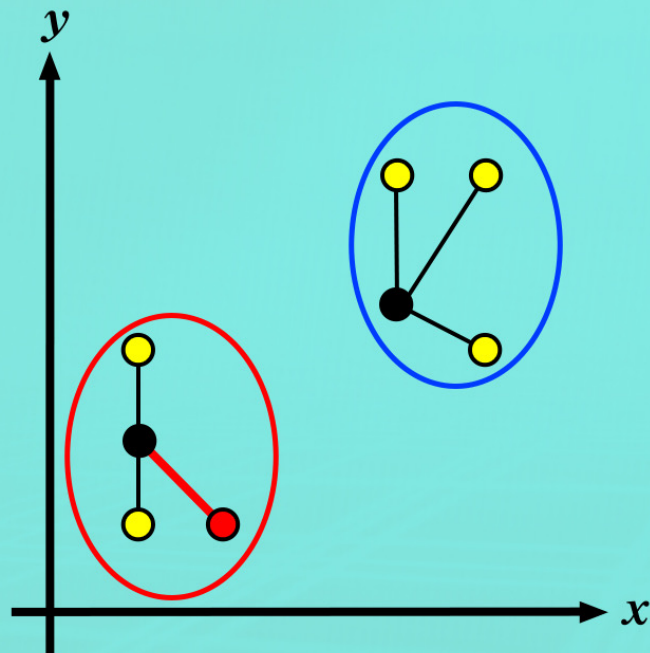
6個の点を2つのグループに分ける



◆ それぞれのグループの重心を
新たな代表の点とする

k 平均法

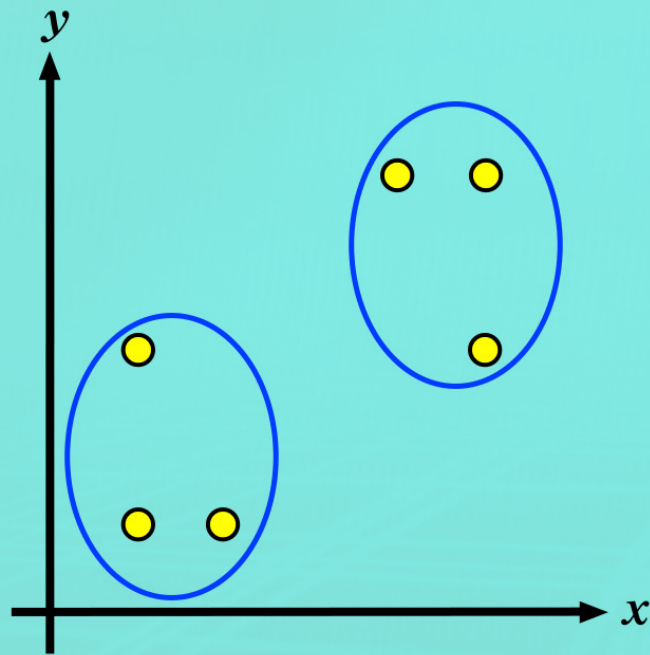
6個の点を2つのグループに分ける



- ◆ 先ほどと同様に
近い方の代表点と同じ
グループとして分類する

k 平均法

6個の点を2つのグループに分ける



- ◆ 代表となる点を選び直す
- ◆ 距離を計算する
- ◆ 近い方のグループの点とする

という作業を繰り返す

グループに変化がなければ終了とする

Rによる演習(1)

```
> library(tidyverse)
> set.seed(38)
> sigma <- 0.25
> ca <- tibble( x=rnorm(50,1,sigma), y=rnorm(50,1,sigma), cl_t = "A")
> cb <- tibble( x=rnorm(50,-1,sigma), y=rnorm(50,1,sigma), cl_t = "B")
> cc <- tibble( x=rnorm(50,-1,sigma), y=rnorm(50,-1,sigma), cl_t = "C")
> cd <- tibble( x=rnorm(50,1,sigma), y=rnorm(50,-1,sigma), cl_t = "D")
> t0 <- rbind(ca,cb,cc,cd)
> ggplot(t0,aes(x=x, y=y, color=cl_t) ) + geom_point()
> k0 <- t0 %>% select(x , y) %>% kmeans(centers = 4)
> k0
> t0 %<>% mutate(cl_k=factor(k0$cluster, label=c("A","B","C","D"), levels=c(4,2,1,3) ) )
> t0 %>% select(cl_t,cl_k) %>% table()
```

```
# (1, 1) を中心に
# (-1, 1) を中心に
# (-1, -1) を中心に
# ( 1, -1) を中心に
# 4つのグループ
```

```
# k平均法
```

```
# k平均法
```

```
# kmeans(訓練データの座標, 中心の数)
```

Rによる演習(2)

```
> min <- -2  
> max <- 2  
> N_test <- 100  
> int <- (max-min) / (N_test-1)  
> x <- rep(seq(min,max,int), each=N_test)  
> y <- rep(seq(min,max,int), N_test)  
> test <- tibble(x=x, y=y)
```

```
# rep は繰り返し, each とするとそれぞれ N_test個  
# eachがない場合には N_test 周繰り返す
```

```
> library(class)  
> cl_test <- t0 %>% select(x, y) %>% knn(test, cl=t0$cl_t, k=5)  
> ggplot(test,aes(x=x,y=y) )+ geom_point(aes(color=cl_test),size=0.5 )
```

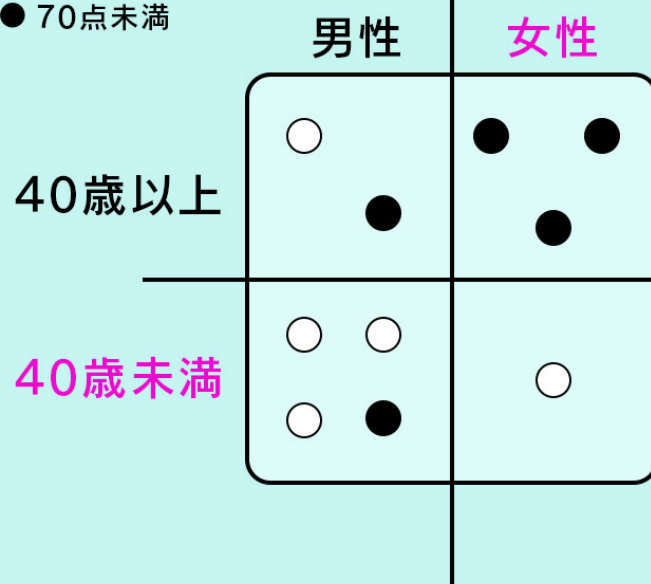
```
# k近傍法  
# cl_test の色で表示
```

```
# k近傍法  
# knn(訓練データの座標, 検証データの座標, 訓練データのラベル, 近傍の数)  
# 検証データの各座標はk個の近傍にある訓練データのラベルの多数決で決める
```

不純度とデータの分割

学生番号	年齢	性別	点数
1	40歳未満	男性	70点以上
2	40歳未満	女性	70点以上
3	40歳未満	男性	70点以上
4	40歳未満	男性	70点以上
5	40歳以上	男性	70点以上
6	40歳未満	男性	70点未満
7	40歳以上	女性	70点未満
8	40歳以上	女性	70点未満
9	40歳以上	男性	70点未満
10	40歳以上	女性	70点未満

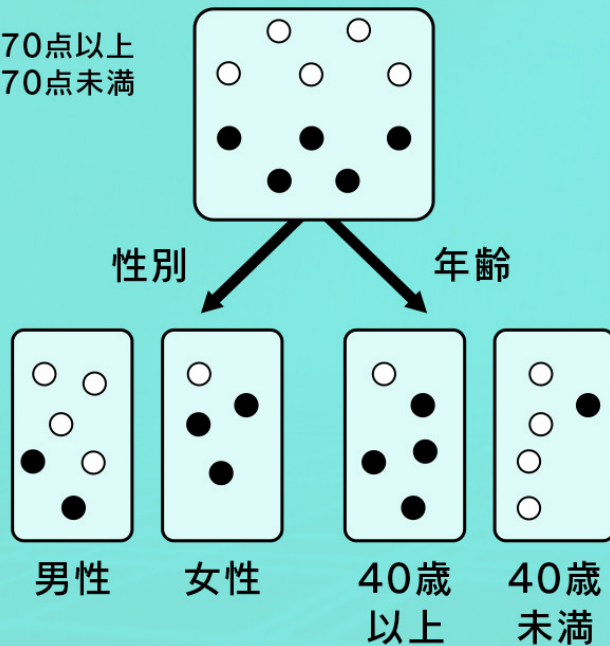
○ 70点以上
● 70点未満



不純度とデータの分割

集合の分割

○ 70点以上
● 70点未満



ジニ係数

不純度が低くなるように分割

$$I = 1 - p_A^2 - p_B^2 \\ = 2p_A(1 - p_A)$$

AとBという 2種類のボール:

1つ取り出して 元に戻して
もう1つ取り出すとする

不純度が低い

⇔ 別のものを取り出す割合が小さい

Rによる演習(3)

```
> library(rpart)
> rp_t <- rpart(data=t0, cl_t ~ x + y, method="class")
> library(rpart.plot)
> rpart.plot(rp_t)
> rp_test <- predict(rp_t, newdata=test, type="class")
```

分類木

predict で予測, type でclass を指定

```
> ggplot(test,aes(x=x,y=y)) + geom_point(aes(color=rp_test), size=0.5 )
```

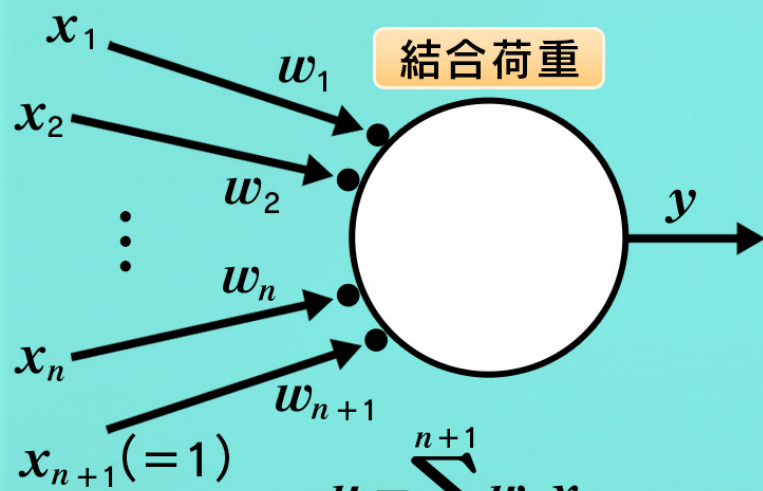
```
> library(nnet)
> t1 <- t0 %>% mutate(cl_t = as.factor(cl_t) )
> nn_t <- nnet(data=t1, cl_t ~ x+y, size=5)
> nn_test <- predict(nn_t, newdata=test, type="class")
```

nnet ではクラス分類をするには因子に変更

predict で予測, type でclass を指定

```
> ggplot(test,aes(x=x,y=y)) +geom_point(aes(color=nn_test), size=0.5 )
```

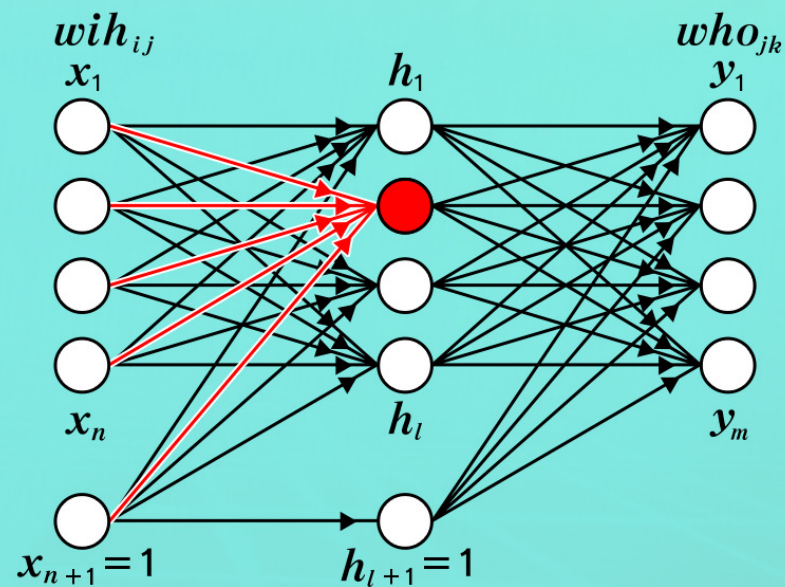
階層型ネットワーク



$$u = \sum_{i=1}^{n+1} w_i x_i$$

$$y = f(u)$$

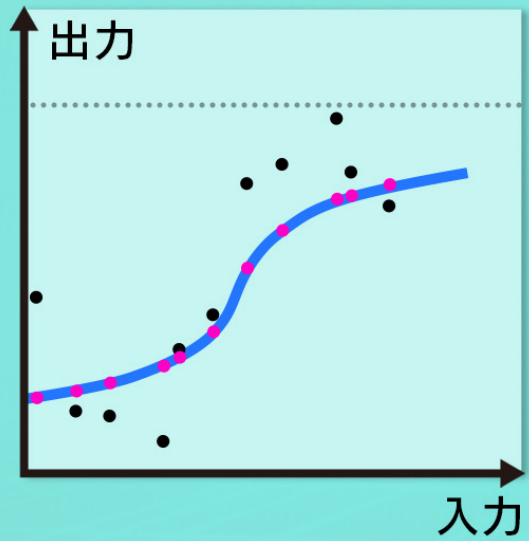
入力から出力へ方向のみ



多入力多出力の関数

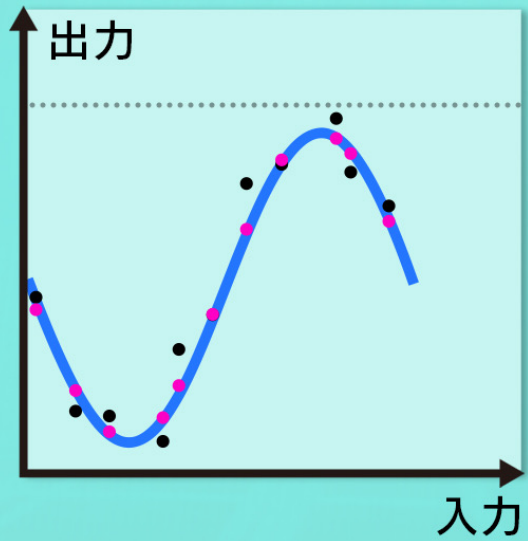
汎化と過学習 (3)

a

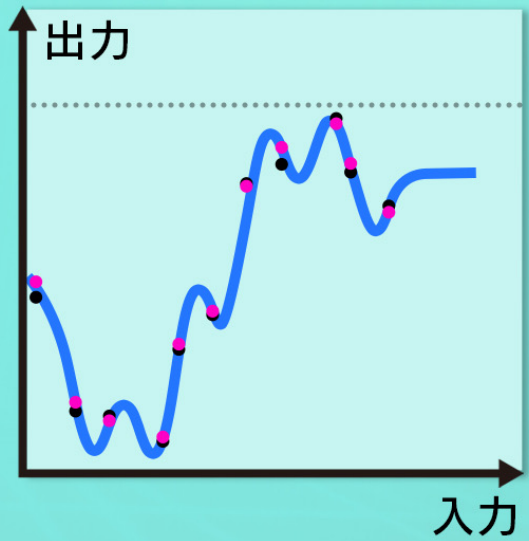


学習が不十分

b



c



過学習

Rによる演習(3)

```
> library(rpart)
> rp_t <- rpart(data=t0, cl_t ~ x + y, method="class")
> library(rpart.plot)
> rpart.plot(rp_t)
> rp_test <- predict(rp_t, newdata=test, type="class")
```

分類木

predict で予測, type でclass を指定

```
> ggplot(test,aes(x=x,y=y)) + geom_point(aes(color=rp_test), size=0.5 )
```

```
> library(nnet)
> t1 <- t0 %>% mutate(cl_t = as.factor(cl_t) )
> nn_t <- nnet(data=t1, cl_t ~ x+y, size=5)
> nn_test <- predict(nn_t, newdata=test, type="class")
```

nnet ではクラス分類をするには因子に変更

predict で予測, type でclass を指定

```
> ggplot(test,aes(x=x,y=y)) +geom_point(aes(color=nn_test), size=0.5 )
```


判定を図る4つの指標

		判定結果	
		陽性	陰性
真の値	陽性	真陽性 (True Positive) TP	偽陰性 (False Negative) FN
	陰性	偽陽性 (False Positive) FP	真陰性 (True Negative) TN

$$\text{正解率 (Accuracy)} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$\text{適合率 (Precision)} = \frac{TP}{TP + FP}$$

$$\text{再現率 (Recall)} = \frac{TP}{TP + FN}$$

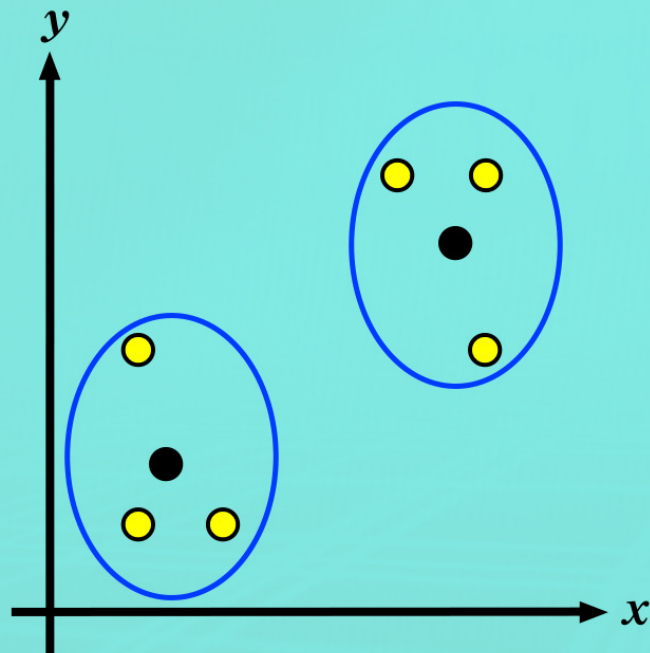
$$\begin{aligned} \text{F値 (F measure)} &= \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \\ &= \frac{2 (\text{Precision} \cdot \text{Recall})}{\text{Precision} + \text{Recall}} \end{aligned}$$

データの分析と

知識発見

k 平均法

6個の点を2つのグループに分ける



◆ 代表となる点を選び直す

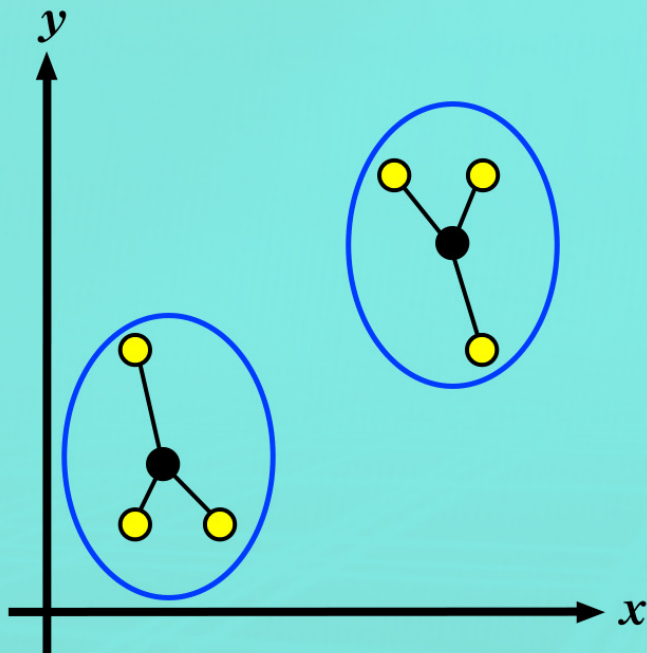
◆ 距離を計算する

◆ 近い方のグループの点とする

という作業を繰り返す

k 平均法

6個の点を2つのグループに分ける



◆ 代表となる点を選び直す

◆ 距離を計算する

◆ 近い方のグループの点とする

という作業を繰り返す