

知的情報処理 3. 最近傍法

櫻井彰人
慶應義塾大学理工学部

本日の目的

- 「k-最近傍法」or「最近傍法」とは
 - 機械学習の最も基本的なアイデア
 - 事例をたくさん記憶する
 - 必要になったら、事例そのものか、近い事例を思い出す。
- もう少し具体的には
- 事例(入力 x_i と出力 y_i の対)を記憶する
 - 問合せ(入力 x に対する出力 y は何?)に対し、
 - 答える(x に近い事例を参照して、最適な y を作る)
- そして
- その基本と問題(実に一般的な問題です)

Rの基礎勉強もします

整理しよう

件名	件名	件名	件名
1. 最近傍法	2. k-最近傍法	3. 線形回帰	4. 線形判別
5. 線形回帰	6. 線形判別	7. 線形回帰	8. 線形判別
9. 線形回帰	10. 線形判別	11. 線形回帰	12. 線形判別
13. 線形回帰	14. 線形判別	15. 線形回帰	16. 線形判別
17. 線形回帰	18. 線形判別	19. 線形回帰	20. 線形判別
21. 線形回帰	22. 線形判別	23. 線形回帰	24. 線形判別
25. 線形回帰	26. 線形判別	27. 線形回帰	28. 線形判別
29. 線形回帰	30. 線形判別	31. 線形回帰	32. 線形判別
33. 線形回帰	34. 線形判別	35. 線形回帰	36. 線形判別
37. 線形回帰	38. 線形判別	39. 線形回帰	40. 線形判別
41. 線形回帰	42. 線形判別	43. 線形回帰	44. 線形判別
45. 線形回帰	46. 線形判別	47. 線形回帰	48. 線形判別
49. 線形回帰	50. 線形判別	51. 線形回帰	52. 線形判別
53. 線形回帰	54. 線形判別	55. 線形回帰	56. 線形判別
57. 線形回帰	58. 線形判別	59. 線形回帰	60. 線形判別
61. 線形回帰	62. 線形判別	63. 線形回帰	64. 線形判別
65. 線形回帰	66. 線形判別	67. 線形回帰	68. 線形判別
69. 線形回帰	70. 線形判別	71. 線形回帰	72. 線形判別
73. 線形回帰	74. 線形判別	75. 線形回帰	76. 線形判別
77. 線形回帰	78. 線形判別	79. 線形回帰	80. 線形判別
81. 線形回帰	82. 線形判別	83. 線形回帰	84. 線形判別
85. 線形回帰	86. 線形判別	87. 線形回帰	88. 線形判別
89. 線形回帰	90. 線形判別	91. 線形回帰	92. 線形判別
93. 線形回帰	94. 線形判別	95. 線形回帰	96. 線形判別
97. 線形回帰	98. 線形判別	99. 線形回帰	100. 線形判別

目次

- 丸暗記(?)
- 作ってみよう(ぴったり一致するもの)
 - 簡単なプログラム例
- 少し発展させよう(ぴったりでなくとも)
 - 近さとは
 - 次元と単位と距離
 - どの事例を参考に
 - 複数個と多数決
 - プログラム例

丸暗記



丸暗記 --- オウムのように

- 考えてみよう。
- 最近傍法に必要なものは?
- 答え:
- 教えられたことを全部、間違いなく覚えること、すなわち、丸暗記
- 問合せがあったら、記憶の中から、最も似たものを探し出す。

最近傍法のポイント

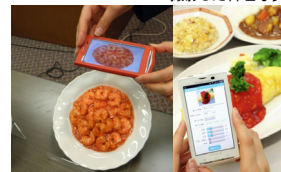
http://teflworldwiki.com/index.php/Rote_Learning
<http://sozaishu.seesaa.net/article/115935835.html>

丸暗記学習するプログラム

- 丸暗記すること自体は、コンピュータにとって難しいことではない(そのようなプログラムを書くことは、難しいことではない)
- しかし、
 - 思い出すときに、速くできるようにする、
 - データが大量でも対応できるようにする
 のは、容易ではない。
- 今回は、これらの問題は考えない(つまり、基本技術だけ考える)

こんなアプリ

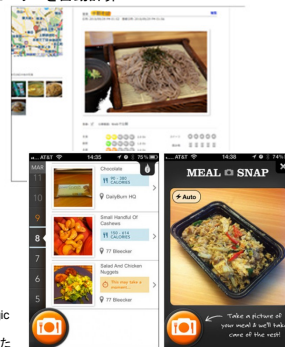
撮影した料理写真でカロリーを自動計算



NTTコミュニケーションズ, NTTレゾナント, foo.log

Meal Snap - Calorie Counting Magic

Amazon の Mechanical Turk を使っているという噂はあった



目次

- 丸暗記(?)
- 作ってみよう(ぴったりに一致するもの)
 - 簡単なプログラム例
- 少し発展させよう(ぴったりでなくとも)
 - 近さとは
 - 次元と単位と距離
 - どの事例を参考に
 - 複数個と多数決
 - プログラム例

作ってみよう

- 課題
- 四字熟語のリストがある。
- 1文字欠けたある四字熟語が与えられたとき、そこに一文字入れて、正しい(つまり与リストにある)四字熟語にしろ。

整理しよう

四字熟語	四字熟語	四字熟語	四字熟語
一期一会	一五一十	一言一句	一汁一菜
一日一善			

概略仕様

- プログラムはRで記述する
- 四字熟語のリストは、sjisコード(Rで読める)で作成してある
- 質問の「1文字欠けたある四字熟語」は、例えば、「縦横無〇」のように、欠け字は「〇」で表現してあるものとする。
- 解は複数ある可能性を排除しない。解の四字熟語を配列(リスト)で返す。
- 簡単のため、
 - プログラム内で四字熟語のリストを読み、
 - 質問の「1文字欠けたある四字熟語」は定数で変数に代入するとする。

何を行えばよいか

for-ループです

```

"D:/R/Sample" を working directory にする
x <- read.table("cwlist.txt", as.is=1)
y <- "縦横無〇"

for ( i in 1:length(x[,1]) ) {
  score <- 0
  for ( j in 1:4 ) {
    if ( x[i,j] == y[j] ) score <- score + 1
  }
  if ( score >= 3 ) { print( i ); print( x[i,] ) }
}
    
```

出力

```

[1] 1219
[1] "縦横無尽"
    
```

同じ文字の個数が3以上なら、印字する

プログラム完成

for-ループです

列を指定する方法です

文字列として読む

```

setwd("D:/R/Sample")
x <- read.table("cwlist.txt", as.is=1)
y <- "縦横無〇"

for ( i in 1:length(x[,1]) ) {
  score <- 0
  for ( j in 1:4 ) {
    if ( substr(x[i,j],j,j) == substr(y,j,j) ) score <- score + 1
  }
  if ( score >= 3 ) { print( i ); print( x[i,] ) }
}
    
```

出力

```

[1] 1219
[1] "縦横無尽"
    
```

同じ位置に同じ文字がある個数を求める

その個数が3以上なら、印字する

例1: 結果を記憶する

```

setwd("D:/R/Sample")
x <- read.table("cwlist.txt", as.is=1)
y <- "縦横無〇"
    
```

```

s <- c()
for ( i in 1:length(x[,1]) ) {
  score <- 0
  for ( j in 1:4 ) {
    if ( substr(x[i,j],j,j) == substr(y,j,j) ) score <- score + 1
  }
  if ( score >= 3 ) s <- c(s,i)
}
print( s ); print( x[s,] )
    
```

出力

```

[1] 1219
[1] "縦横無尽"
    
```

ベクトルの最後尾にiを追加する

これはskipしましょう

例2: Rの機能活用

```
setwd("D:/R/Sample")
x <- read.table("cwl1st.txt", as.is=1)
y <- "縦横無尽"

s <- c( )
for ( i in 1:length(x[,1]) ) {
  if ( length(agrep( y, x[i,] )) != 0 ) s <- c(s,i)
}
print( s ); print( x[s,] )
```

第一引数の文字列にほぼ等しい(defaultでは一文字間違いまで)文字列が引数2にあるとき、第二引数での位置を返す。ないとき、長さ0のベクトルを返す。

出力

```
[1] 1219
[1] "縦横無尽"
```

これは学習か？

- 単なるプログラム！
 - 否定するわけではないが
- 学習の基本であることは確かである。
人間でも、
 - 文字を覚える
 - 数を覚える
 - 九九を覚える、
- このプログラムでは、学習らしくない。
- しかし、後で使いやすいように、データの記憶方法を工夫をすると、学習らしくなる。
 - つまり、記憶するにあたって、データの性質を(データの内的構造を)利用する。これこそ、学習である。

目次

- 丸暗記(?)
- 作ってみよう(ぴったり一致するもの)
 - 簡単なプログラム例
- 少し発展させよう(ぴったりでなくとも)
 - 近さとは
 - 次元と単位と距離
 - どの事例を参考に
 - 複数個と多数決
 - プログラム例

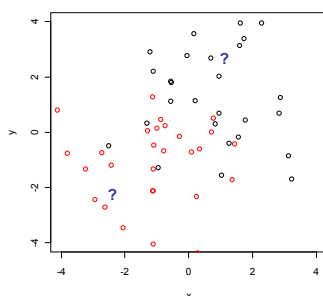
少し発展させよう

- 右の問題で、穴埋めする時
- 説明変数値が
 - 知っているものと同じか
 - かなり近いときに、被説明変数の値を答えにするという方法を考えてみよう

整理しよう

	穴埋め1	穴埋め2	穴埋め3	穴埋め4
説明変数	1	2	3	4
被説明変数	5	6	7	8
説明変数	9	10	11	12
被説明変数	13	14	15	16
説明変数	17	18	19	20
被説明変数	21	22	23	24
説明変数	25	26	27	28
被説明変数	29	30	31	32
説明変数	33	34	35	36
被説明変数	37	38	39	40
説明変数	41	42	43	44
被説明変数	45	46	47	48
説明変数	49	50	51	52
被説明変数	53	54	55	56
説明変数	57	58	59	60
被説明変数	61	62	63	64
説明変数	65	66	67	68
被説明変数	69	70	71	72
説明変数	73	74	75	76
被説明変数	77	78	79	80
説明変数	81	82	83	84
被説明変数	85	86	87	88
説明変数	89	90	91	92
被説明変数	93	94	95	96
説明変数	97	98	99	100
被説明変数	101	102	103	104
説明変数	105	106	107	108
被説明変数	109	110	111	112
説明変数	113	114	115	116
被説明変数	117	118	119	120
説明変数	121	122	123	124
被説明変数	125	126	127	128
説明変数	129	130	131	132
被説明変数	133	134	135	136
説明変数	137	138	139	140
被説明変数	141	142	143	144
説明変数	145	146	147	148
被説明変数	149	150	151	152
説明変数	153	154	155	156
被説明変数	157	158	159	160
説明変数	161	162	163	164
被説明変数	165	166	167	168
説明変数	169	170	171	172
被説明変数	173	174	175	176
説明変数	177	178	179	180
被説明変数	181	182	183	184
説明変数	185	186	187	188
被説明変数	189	190	191	192
説明変数	193	194	195	196
被説明変数	197	198	199	200

図で見ると分りやすい



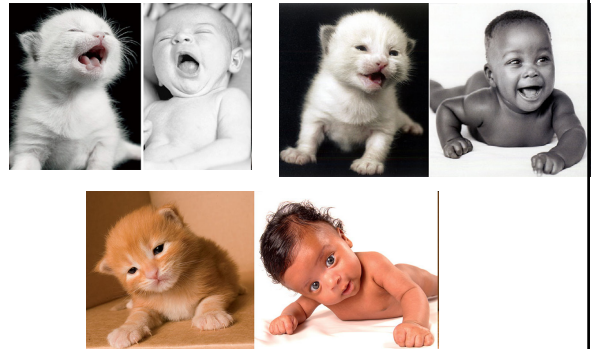
目次

- 丸暗記(?)
- 作ってみよう(ぴったり一致するもの)
 - 簡単なプログラム例
- 少し発展させよう(ぴったりでなくとも)
 - 近さとは
 - 次元と単位と距離
 - どの事例を参考に
 - 複数個と多数決
 - プログラム例

小目次

- 近さとは
 - 定義することの困難さ
 - 人間は、似ているものの発見が得意
 - しかし、言葉にできない
 - 「距離」でいいのか？
 - 単位の問題
 - 次元の問題
 - 正規化で解決

「似ている」の困難さ



http://gigazine.net/news/20071006_cat_similar_baby/

なかなか難しい



脱線: 逆手にとって: CAPTCHA

Completely Automated Public Turing Test To Tell Computers and Humans Apart



この画像の中に表示されている文字を入力してください

画像: D5B0J82U
*文字の入力:

<http://www.captcha.net/>

似ているとは近いこと(?)

- 似ているとは近いことと、考えよう
- 近いとは「距離」が小さいこと。つまり、「近さ」とは「距離」である。
 - 同語反復。でもない。
 - 「距離」というと、数学では、距離の公理を満たす概念。
 - 距離の公理とは？
 - 非負性: $d(x, y) \geq 0$
 - 同一性: $x = y \Leftrightarrow d(x, y) = 0$
 - 対称性: $d(x, y) = d(y, x)$
 - 三角不等式: $d(x, y) + d(y, z) \geq d(x, z)$
 - 代表的なものは:
 - ユークリッド距離: $\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2}$
 - マンハッタン距離: $|x_1 - y_1| + |x_2 - y_2| + |x_3 - y_3|$

小目次

- 近さとは
 - 定義することの困難さ
 - 人間は、似ているものの発見が得意
 - しかし、言葉にできない
 - 「距離」でいいのか？
 - 単位の問題
 - 次元の問題
 - 正規化で解決

「距離」でいいのか？

- 例として、身長と体重で、AさんとBさんの近さをみることにしよう
 - Aさんは (h_A, w_A) , Bさんは (h_B, w_B) としよう
- 距離は？
 - $\sqrt{(h_A - h_B)^2 + (w_A - w_B)^2}$? or $|h_A - h_B| + |w_A - w_B|$? or ...
- そもそも単位は？
 - 身長は、cm? mm? μm ? m? km?
 - 体重は、kg? g? mg? pg? t?

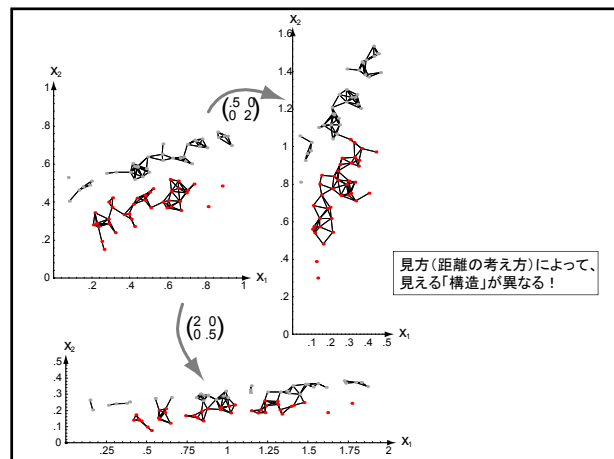
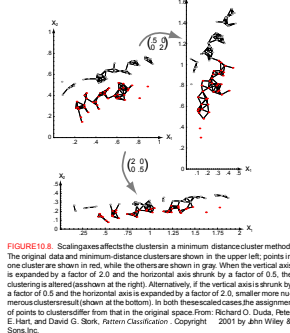
閑話休題。
√の横棒がないけどおかしくはない。そもそもその横棒って何だ？

単位の問題とは？

- 身長と体重を2個の属性とし、次の三人の距離をみよう
 - A: (170cm, 65kg)
 - B: (180cm, 60kg)
 - C: (175cm, 70kg)
- マンハッタン距離を考えよう
 - $d(A, B) = 10 + 5 = 15$
 - $d(B, C) = 5 + 10 = 15$
 - $d(C, A) = 5 + 5 = 10$
- cm と kg を用いるのは不公平。単位の名称からいって、m と g だろう。となると、
 - $d(A, B) = 0.01 + 5000 = 5000.01$
 - $d(B, C) = 0.005 + 10000 = 10000.005$
 - $d(C, A) = 0.005 + 5000 = 5000.005$
- 他にも、いくらでも、考えられる。さて、どうしたものか。

え？ そう？

単位の問題だけではない



つまり: 次元と単位

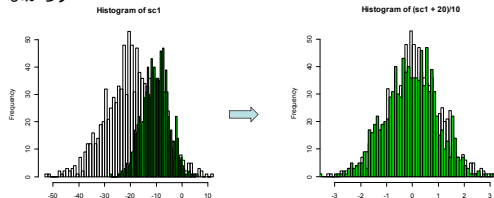
- 問題は2つある: 次元と単位
- 次元: 異なる物理量を区別する概念。次元が異なる物理量は比較できない。
 - 「1kgと1mとは、どちらがより偉いか？」は有意味か
 - なお、この例に表れる次元は3個。重さ、長さ、偉さ。
 - 「偉さ」が物理量かという、まあ、そうではないが。
- 単位: 物理量を数値で表すとき、数値の1に対応する物理量
 - 長さの単位の例: m, mile, 尺, Å
 - 勿論、1cm を単位にしても、3.14km を単位にしてもよい

次元と単位と距離

- 次元解析: 物理量を、質量、長さ、時間、電荷、温度などの次元の有理数幂で表現し、方程式や仮説の妥当性を調べる
- 目下の話題は、「次元」
 - 異なる次元の量は、比較・加算ができない。
 - 当然、(170cm, 65kg) と (180cm, 60kg) の距離をユークリッド距離的に、また、マンハッタン距離的には定義できない。
 - では、どうする？
- もう一つの話題は、「単位」
 - 異なる単位の量は、比較・加算ができない。
 - 身長と指の長さが属性のとき、(1.7m, 10cm) と (1.8m, 12cm) の距離をユークリッド距離的に、また、マンハッタン距離的には定義できない。
 - では、どうする？

正規化

- もしデータが大量にあれば、そして、各属性値が独立に正規分布していれば(ある確率分布に従っていれば)、「値/標準偏差」という量を使えばよい
 - この量は、無次元量。なぜなら、「値」もその「標準偏差」も同じ次元。
 - 「無次元量だから比較・加算できる」という訳ではないが、
 - 無次元の上に、各属性値の標準偏差が1に規格化されたのだから、よからう

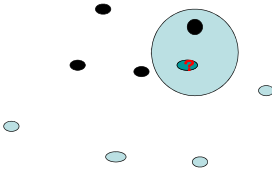


目次

- 丸暗記(?)
- 作ってみよう(ぴったり一致するもの)
 - 簡単なプログラム例
- 少し発展させよう(ぴったりでなくとも)
 - 近さとは
 - 次元と単位と距離
 - どの事例を参考に
 - 複数個と多数決
 - プログラム例

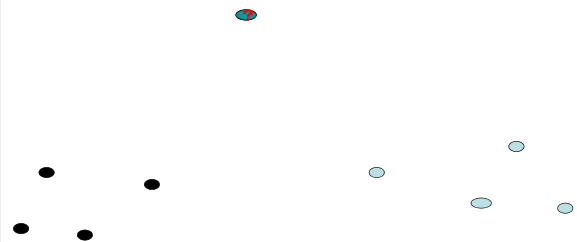
どの事例を参考に？

- 「最も近い事例」にしよう。そうすれば、
- 「どのくらい近ければよい」という基準を考えなくてすむ。
 - つまり、(適当に決めた基準で)見つからなかった場合にも対処できる



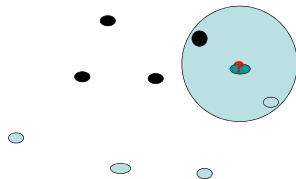
最も近いものが遠いと？

- 「どのくらい近ければよい」という基準を考えなくてすむ。」とはいえ



近い事例が複数あると？

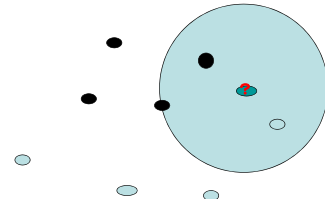
- 「一つ」に決めてしまったが、ちょっと不安。
- 例えば、



– どっち？

多数決はどうか？

- 近い方から3人(とか5人とか7人とか...)の意見を聞くことにしたらどうだろうか。
- そして多数決をとればよいだろう
- 一般には、近くの k 人から聴くことにする、この方法を k -nearest neighbor 法という



目次

- 丸暗記(?)
- 作ってみよう(ぴったり一致するもの)
 - 簡単なプログラム例
- 少し発展させよう(ぴったりでなくとも)
 - 近さとは
 - 次元と単位と距離
 - どの事例を参考に
 - 複数個と多数決
 - プログラム例

深みに: 多数決のプログラムは?

- まず、「多数」を選ばないといけない。
 - プログラムはちょっと複雑
- 例えば、集まった100人のうち、身長が最も低い5人を選ぶ方法を考えることになる。

どうすればいい?

- 一つの方法は、昔とった杵柄(のはず)のバブルソートです。
- 100人をバブルソートする途中で、つまり、上位5人が決まった時点で、バブルソートを中止すればよい。

深みに: 多数決のプログラムは? (続)

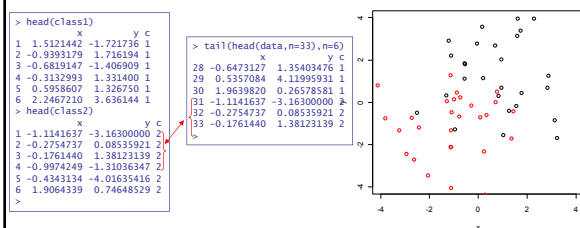
- しかし、長いプログラムを書くのはサバることにしよう。
 - その代わりプログラム実行時間はかかるようになるのだが、まあ、我慢。
- テストする x との距離を、すべての学習データに対して計算し、それを、**全部(ここが無駄)ソートしてしま**い(R内の関数を使えばプログラムは書かなくてすむので)、それから先頭の必要個数(今回は k 個)を使えばよい。
- 朗報! R で sort の help を見てください。部分的にソートしてくれそうです。partial=5 とすると、最小値5個が得られたところで、ソートは終了するようです。
- しかし、index (何番目の要素かを示す値)は返してくれない。返すように指定する(index.return=TRUE)のは完全ソートのときのみ可能
- partial は諦めて、index.return=TRUE としよう。
- なお、rank という関数もある。

実プログラム例: 2次元の場合

- 2次元かつ 0 or 1 に分類する場合の k -NN プログラムを作ってみよう。
- 入力の近さ判定にはユークリッド距離、出力の近似には、多数決を使う
- 実験用のデータとして、平均が (1,1) と (-1,-1)、分散共分散行列が $1.5^2 I$ (I は単位行列)である2個の正規分布を考える。二つのクラスの事例は、それぞれの正規分布に従い生成されると考える。
- 何個かデータを生成する

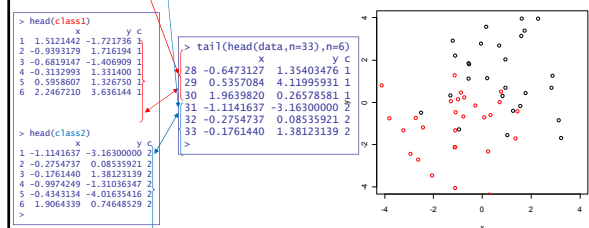
実験用データの作成

```
class1<-data.frame(x=rnorm(30,mean=1,sd=1.5),y=rnorm(30,mean=1,sd=1.5),c=rep(1,30))
class2<-data.frame(x=rnorm(30,mean=-1,sd=1.5),y=rnorm(30,mean=-1,sd=1.5),c=rep(2,30))
data<-rbind(class1,class2)
plot(subset(data,data$c==1)[,1:2],xlim=c(-4,4),ylim=c(-4,4))
points(subset(data,data$c==2)[,1:2],col="red",xlim=c(-4,4),ylim=c(-4,4))
```



何をやっているか

`class1<-data.frame(x=rnorm(30,mean=1,sd=1.5),y=rnorm(30,mean=1,sd=1.5),c=rep(1,30))`
 データフレームにするにあたり x という名の列 y という名の列に平均1、標準偏差1.5の正規乱数30個
`data<-rbind(class1,class2)` 行方向につなげる。
`plot(subset(data,data$c==1)[,1:2],xlim=c(-4,4),ylim=c(-4,4))`
`points(subset(data,data$c==2)[,1:2],col="red",xlim=c(-4,4),ylim=c(-4,4))`



多数決はどうしよう？

- table() 関数を用いて、数えることにする。
 - table() は分割表 (contingency table, cross tabulation, or cross tab) を作る関数

```
> x <- c(0,1,0,1,1,0,0,0)
> (tbl<-table(x))
 0 1 
5 3 
> as.data.frame(tbl)
  x Freq
1 0     5
2 1     3
> (count1<-as.data.frame(tbl)$Freq[as.data.frame(tbl)$x==1])
[1] 5
> (count2<-as.data.frame(tbl)$Freq[as.data.frame(tbl)$x==2])
[1] 3
>
```

多数決

- 頻度を求め、その名前 (factor の level) からクラスを表す数値を取り出したい。

```
# majority(c(0,1,2,2,0,1,0,2,1,2,0,1,1)) = 1 となるような関数
majority <- function( x ) {
  return( as.numeric(names(which.max(table(x)))) )
}
```

まとめれば k-NN

多数決関数の定義

```
# majority(c(0,1,2,2,0,1,0,2,1,2,0,1,1)) = 1 となるような関数
majority <- function( x ) {
  return( as.numeric(names(which.max(table(x)))) )
}
```

実験用データの作成

```
class1<-data.frame(x1=rnorm(30,mean=1,sd=1.5),
                   x2=rnorm(30,mean=1,sd=1.5),c=rep(1,30))
class2<-data.frame(x1=rnorm(30,mean=-1,sd=1.5),
                   x2=rnorm(30,mean=-1,sd=1.5),c=rep(2,30))
data<-rbind(class1,class2)
```

Main program

```
# k-NN majority
k <- 3
x1<- 0 # テスト点の x1 座標
x2<- 0 # テスト点の x2 座標
sorted <- sort( ( x1 - data$x1 )^2 + ( x2 - data$x2 )^2, index.return=TRUE)
majority( data$sorted$ix[ 1:k ] )
```

境界線が表示できます

```
# k of k-NN
k <- 3

# test data which, in this case, are used for the plot
xtest1 <- seq(min(data$x1),max(data$x1),length=30)
xtest2 <- seq(min(data$x2),max(data$x2),length=30)
xtest <- expand.grid(xtest1,xtest2)

# use k-NN to get predictions on the grid
cPredicted <- rep(0,30*30)
for ( i in 1:length(xtest[,1]) ) {
  x1 <- xtest$x1[i]
  x2 <- xtest$x2[i]
  sorted <- sort( ( x1 - data$x1 )^2 + ( x2 - data$x2 )^2, index.return=TRUE)
  cPredicted[i] <- majority( data$sorted$ix[ 1:k ] )
}

# cPredicted in a matrix form to be used for "plot"
cPredictedMatrix <- matrix( cPredicted, 30, 30 )
```

何を行おうとしているか？へのヒント

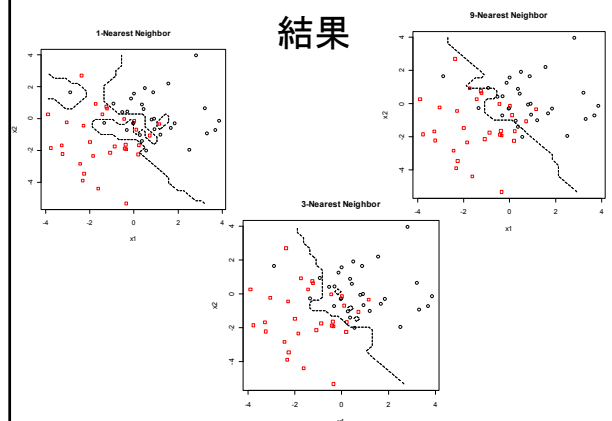
```
plot(subset(data,data$c==1)[,1:2],xlim=c(-4,4),ylim=c(-4,4))
points(subset(data,data$c==2)[,1:2],col="red",xlim=c(-4,4),ylim=c(-4,4))
points(xtest,pch=".",col="blue")
```

プログラム(続)

```
# plot "data"
plot(data[1:2],
      xlab="x1", ylab="x2", # labels of axes
      main=paste(k,"-Nearest Neighbor",sep=""), # title of the graph
      col=data$c, pch=as.real(data$c)+20 ) # color and characters to be used

# draw boundaries
for(i in 1:max(cPredicted))
  contour(
    xtest1, xtest2, cPredictedMatrix==i, # x, y, and z where z is 0 or 1
    levels=0.5, # levels=0.5 i.e. the middle of 0 and 1
    drawlabels=FALSE, add=TRUE,
    col=1, lty=3) # color is black and line type is broken
```

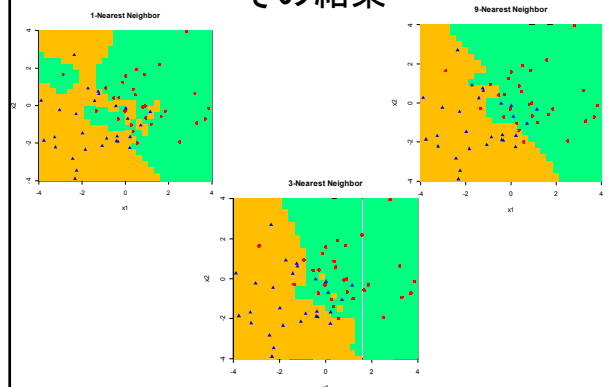
結果



別のプロット

```
# another way to plot
# the following is for just two classes
image( xtest1, xtest2, cPredictedMatrix==2,
       col=c("#00FF80", "#FFB000"), axes=TRUE,
       xlim=c(-4,4), ylim=c(-4,4),
       xlab="x1", ylab="x2", main=paste(k, "-Nearest Neighbor", sep="") )
points( data[1:2], col=ifelse(data$c==1,"red","blue"),
        pch=ifelse(data$c==1,16,17) )
```

その結果



まとめにかえて(後半)

目次

- どの「穴埋め」か？
- k-NN (k nearest neighbor, k 最近傍法)とは
- なぜ「境界」がでこぼこか？
- k はどう決めるか

解いた問題

- 穴埋め問題1
- 穴埋め問題2-3
- 説明変数値が
 - ・ 知っているものと同じか
 - ・ かなり近いときに、
被説明変数の値を答えにする
という方法

[illegible]

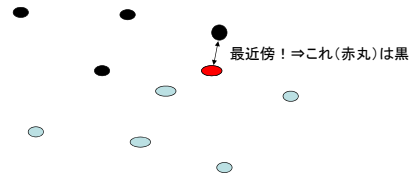
目次

- どの「穴埋め」か？
- k-NN (k nearest neighbor, k 最近傍法)とは
- なぜ「境界」がでこぼこか？
- k はどう決めるか

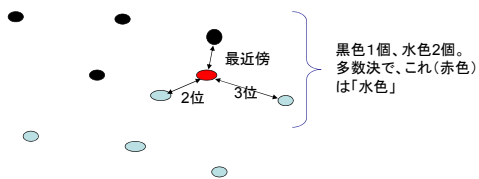
最近傍法

- 最近傍法 (Nearest neighbor)
 - 問合せ x_q に対し、最近接の x_n を見つけ、
 $f(x_q) \leftarrow f(x_n)$ と (f を近似) する
- k -Nearest neighbor
 - k 個の最近接データの間で、多数決、または
 - k 個の最近接データの間で、平均値、または、、、

1-Nearest Neighbor



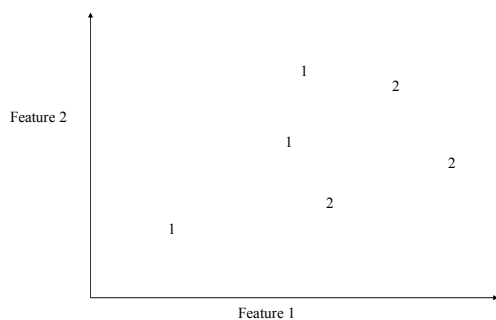
3-Nearest Neighbor



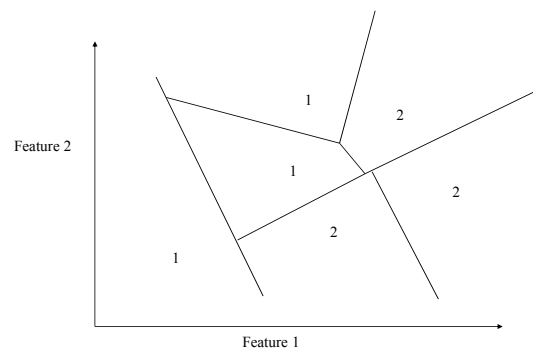
目次

- どの「穴埋め」か?
- k -NN (k nearest neighbor, k 最近傍法) とは
- なぜ「境界」がでこぼこか?
- k はどう決めるか

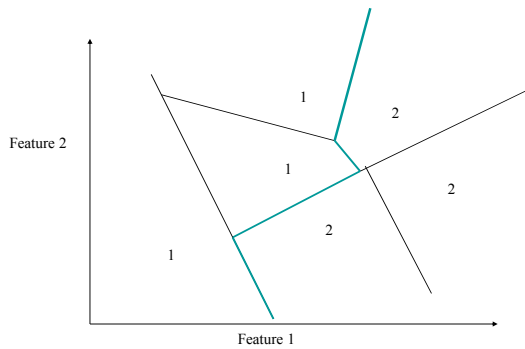
1-NNの境界線の幾何的解釈



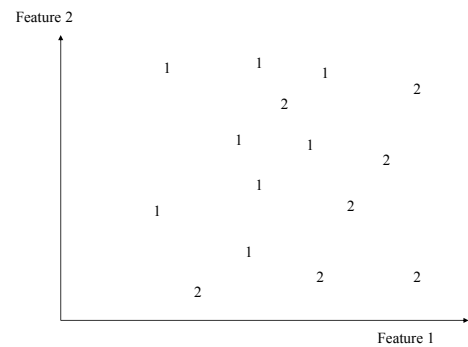
境界



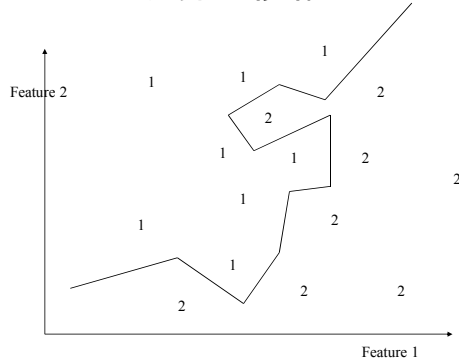
境界



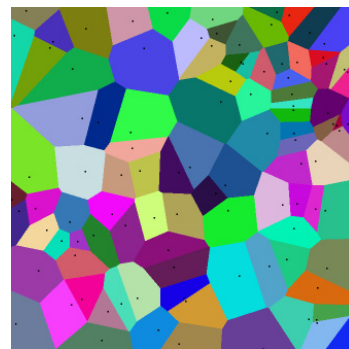
データ点が多いとき



境界は複雑となる



Voronoi図



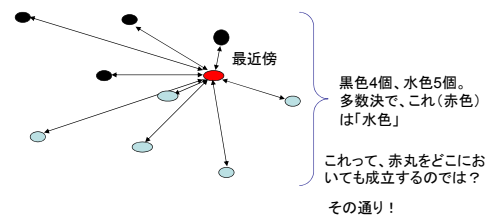
目次

- どの「穴埋め」か？
- k-NN (k nearest neighbor, k 最近傍法) とは
- なぜ「境界」がでこぼこか？
- k はどう決めるか

k が大きすぎると

9-Nearest Neighbor

三年前、いい点に気が付いた人がいました。



つまり、全員に意見を聞くと、各員の答えはいつも同じだから、答えはいつも(どこで聞いても)同じ

k の選択

- こうしたパラメータの設定(この講義でもこれから頻繁に出てくる)は、一般に、難しい
- 理論的根拠をもって、または、データ(validation data set)を用いて決めることになる
- k-NNについては、bootstrapによる方法が提案されている

Peter Hall, Byeong U. Park, and Richard J. Samworth.
Choice of neighbor order in nearest-neighbor classification.
The Annals of Statistics, 2008, Vol. 36, No. 5, 2135–2152.

即レポ2

練習問題

- この講義で作成した k-NN を用いて、データもこの講義で作成したものを用いて、k を大きくしたとき、境界は、 $(1,1)$ と $(-1,-1)$ の垂直二等分線に近くなるか？
 - k を変えてみるだけでなく、データ数も変えてみよ。
 - なぜ「この特別な2点の垂直二等分線」なのだろうか、その理由を考えてみてください。

一昨年、全然垂直二等分線らしくないのに、「垂直二等分線になりました」とレポートに書いた人がいました。これは「白を黒という」に等しい。