

主成分分析大全

 生信技能树

关注

2 2018.12.21 11:58:51 字数 3,657 阅读 9,444

1 背景

主成分分析法是数据挖掘中常用的一种降维算法,是Pearson在1901年提出的,再后来由hotelling在1933年加以发展提出的一种多变量的统计方法,其最主要的用途在于“降维”,通过析取主成分显出的最大的个别差异,也可以用来削减回归分析和聚类分析中变量的数目,与因子分析类似。

所谓降维,就是把具有相关性的变量数目减少,用较少的变量来取代原先变量。如果原始变量互相正交,即没有相关性,则主成分分析没有效果。

在生物信息学的实际应用情况中,通常是得到了成百上千个基因的信息,这些基因相互之间会有影响,通过主成分分析后,得到有限的几个主成分就可以代表它们的基因了。也就是所谓的降维。

R语言有非常多的途径做主成分分析,比如自带的princomp()和psych包的principal()函数,还有gmodels包的fast.prcomp函数。

2 拆解主成分分析步骤

实际应用时我们通常会选择主成分分析函数,直接把input数据一步分析到位,只需要看懂输出结果即可。但是为了加深理解,这里一步步拆解主成分分析步骤,讲解原理。

2.1 测试数据

数据集USJudgeRatings包含了律师对美国高等法院法官的评分。数据框包含43个样本,12个变量!

下面简单看一看这12个变量是什么,以及它们的相关性。

```
library(knitr)
kable(head(USJudgeRatings))
```

	CO NT	INT G	DM NR	DIL G	CF MG	DE CI	PR EP	FA MI	OR AL	WR IT	PH YS	RTE N
AARONSON, L.H.	5.7	7.9	7.7	7.3	7.1	7.4	7.1	7.1	7.1	7.0	8.3	7.8
ALEXANDER,J .M.	6.8	8.9	8.8	8.5	7.8	8.1	8.0	8.0	7.8	7.9	8.5	8.7

推荐阅读

- 群落多样性之Beta多样性（三）
阅读 314
- R 稳健PCA分析
阅读 165
- R数据分析：主成分分析及可视化
阅读 717
- 基于R的线性混合效应模型分析
阅读 518
- 无监督学习-线性方法|深度学习（李宏毅）（十七）
阅读 448



BERDON,R.I.	6.8	8.8	8.5	8.8	8.3	8.5	8.7	8.7	8.4	8.5	8.8	8.7
BRACKEN,J.J.	7.3	6.4	4.3	6.5	6.0	6.2	5.7	5.7	5.1	5.3	5.5	4.8
BURNS,E.B.	6.2	8.8	8.7	8.5	7.9	8.0	8.1	8.0	8.0	8.0	8.6	8.6

这12个变量的介绍如下：

[,1]	CONT	Number of contacts of lawyer with judge.
[,2]	INTG	Judicial integrity.司法诚实性
[,3]	DMNR	Demeanor.风度；举止；行为
[,4]	DILG	Diligence.勤奋，勤勉；注意的程度
[,5]	CFMG	Case flow managing.
[,6]	DECI	Prompt decisions.
[,7]	PREP	Preparation for trial.
[,8]	FAMI	Familiarity with law.
[,9]	ORAL	Sound oral rulings.
[,10]	WRIT	Sound written rulings.
[,11]	PHYS	Physical ability.
[,12]	RTEN	Worthy of retention.

这些是专业领域的用词，大家可以先不用纠结具体细节。

2.2 为什么要做主成分分析

不管三七二十一就直接套用统计方法都是耍流氓，做主成分分析并不是拍脑袋决定的。在这个例子里面，我们拿到了这43个法官的12个信息，就可以通过这12个指标来对法官进行分类，但也许实际情况下收集其他法官的12个信息比较麻烦，所以我们希望只收集三五个信息即可，然后也想达到比较好的分类效果。或者至少也得剔除几个指标吧，这个时候主成分分析就能派上用场啦。这12个变量能得到12个主成分，如果前两个主成分可以揭示85%以上的变异度，也就是说我们可以用两个主成分来代替那12个指标。

在生物信息学领域，比如我们测了1000个病人的2万个基因的表达矩阵，同时也有他们的健康状态信息。那么我们想仔细研究这些数据，想得到基因表达与健康状态的某种关系。这样我就可以对其余几十亿的人检测基因表达来预测其健康状态。如果我们进行了主成分分析，就可以选择解释度比较高的主成分对应的基因，可能就几十上百个而已，大幅度的降低广泛的基因检测成本。

2.3 step1:数据标准化(中心化)

```
dat_scale=scale(USJudgeRatings,scale=F)
options(digits=4, scipen=4)
kable(head(dat_scale))
```

	CO NT	INT G	DM NR	DIL G	CF MG	DEC I	PRE P	FA MI	OR AL	WRI T	PHY S	RTE N
AARONSON,L.H.	-1.7 372	-0.1 209	0.18 37	-0.39 3	-0.3 791	-0.1 651	-0.3 674	-0.3 884	-0. 193	-0.3 837	0.36 51	0.19 77
ALEXANDE	-0.6	0.87	1.28	0.8	0.32	0.53	0.53	0.51	0.5	0.51	0.56	1.09



写下你的评论...

评论0

赞35



推荐阅读

群落多样性之Beta多样性（三）
阅读 314

R 稳健PCA分析
阅读 165

R数据分析：主成分分析及可视化
阅读 717

基于R的线性混合效应模型分析
阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）
阅读 448

主成分分析大全

ARMENTA NO,A.J.	-0.2 372	0.07 91	0.28 37	0.1 07	0.02 09	0.03 49	0.03 26	0.01 16	0.0 07	0.01 63	-0.0 349	0.19 77
BERDON,R. I.	-0.6 372	0.77 91	0.98 37	1.1 07	0.82 09	0.93 49	1.23 26	1.21 16	1.1 07	1.11 63	0.86 51	1.09 77
BRACKEN,J .J.	-0.1 372	-1.6 209	-3.2 163	-1. 19 3	-1.4 791	-1.3 651	-1.7 674	-1.7 884	-2. 193	-2.0 837	-2.4 349	-2.8 023
BURNS,E.B.	-1.2 372	0.77 91	1.18 37	0.8 07	0.42 09	0.43 49	0.63 26	0.51 16	0.7 07	0.61 63	0.66 51	0.99 77

scale()是对数据中心化的函数，当参数scale=F时，表示将数据按列减去平均值，scale=T表示按列进行标准化，公式为 $(x - \text{mean}(x)) / \text{sd}(x)$ ，本例采用中心化。

2.4 step2:求相关系数矩阵

```
dat_cor=cor(dat_scale)
options(digits=4, scipen=4)
kable(dat_cor)
```

	CON T	INT G	DM NR	DIL G	CF MG	DEC I	PRE P	FAM I	ORA L	WRI T	PHY S	RTE N
CO NT	1.00 00	-0.13 32	-0.15 37	0.01 24	0.13 69	0.08 65	0.01 15	-0.02 56	-0.01 20	-0.04 38	0.05 42	-0.03 36
INT G	-0.13 32	1.00 00	0.96 46	0.87 15	0.81 41	0.80 28	0.87 78	0.86 89	0.91 14	0.90 88	0.74 19	0.937 3
DM NR	-0.15 37	0.96 46	1.00 00	0.83 69	0.81 34	0.80 41	0.85 58	0.84 12	0.90 68	0.89 31	0.78 87	0.943 7
DIL G	0.01 24	0.87 15	0.83 69	1.00 00	0.95 88	0.95 62	0.97 86	0.95 74	0.95 45	0.95 93	0.81 29	0.930 0
CF MG	0.13 69	0.81 41	0.81 34	0.95 88	1.00 00	0.98 11	0.95 79	0.93 55	0.95 06	0.94 22	0.87 95	0.927 1
DEC I	0.08 65	0.80 28	0.80 41	0.95 62	0.98 11	1.00 00	0.95 71	0.94 28	0.94 83	0.94 61	0.87 18	0.925 0
PRE P	0.01 15	0.87 78	0.85 58	0.97 86	0.95 79	0.95 71	1.00 00	0.98 99	0.98 31	0.98 68	0.84 87	0.950 3
FA MI	-0.02 56	0.86 89	0.84 12	0.95 74	0.93 55	0.94 28	0.98 99	1.00 00	0.98 13	0.99 07	0.84 37	0.941 6
ORA L	-0.01 20	0.91 14	0.90 68	0.95 45	0.95 06	0.94 83	0.98 31	0.98 13	1.00 00	0.99 34	0.89 12	0.982 1
WRI T	-0.04 38	0.90 88	0.89 31	0.95 93	0.94 22	0.94 61	0.98 68	0.99 07	0.99 34	1.00 00	0.85 59	0.967 6
PHY S	0.05 42	0.74 19	0.78 87	0.81 29	0.87 95	0.87 18	0.84 87	0.84 37	0.89 12	0.85 59	1.00 00	0.906 5
RTE N	-0.03 36	0.93 73	0.94 37	0.93 00	0.92 71	0.92 50	0.95 03	0.94 16	0.98 21	0.96 76	0.90 65	1.000 0



主成分分析大全



生信技能树

关注

赞赏支持

利用eigen函数计算相关系数矩阵的特征值和特征向量。

```
dat_eigen=eigen(dat_cor)
dat_var=dat_eigen$values ## 相关系数矩阵的特征值
options(digits=4, scipen=4)
dat_var
```

```
## [1] 10.133504 1.104147 0.332902 0.253847 0.084453 0.037286 0.019683
## [8] 0.015415 0.007833 0.005612 0.003258 0.002060
```

```
pca_var=dat_var/sum(dat_var)
pca_var
```

```
## [1] 0.8444586 0.0920122 0.0277418 0.0211539 0.0070377 0.0031072 0.0016402
## [8] 0.0012846 0.0006528 0.0004676 0.0002715 0.0001717
```

```
pca_cvar=cumsum(dat_var)/sum(dat_var)
pca_cvar
```

```
## [1] 0.8445 0.9365 0.9642 0.9854 0.9924 0.9955 0.9972 0.9984 0.9991 0.9996
## [11] 0.9998 1.0000
```

可以看出，PC1(84.4%)和PC2(9.2%)共可以解释这12个变量的93.6的程度，除了CONT外的其他的11个变量与PC1都有较好的相关性，所以PC1与这11个变量基本斜交，而CONT不能被PC1表示，所以基本与PC1正交垂直，而PC2与CONT基本平行，表示其基本可以表示CONT。

2.6 step4:崖低碎石图和累积贡献图

```
library(ggplot2)
p=ggplot(aes(x=1:12,y=pca_var))
p1=ggplot(aes(x=1:12,y=pca_cvar))
p+geom_point(pch=2,lwd=3,col=2)+geom_line(col=2,lwd=1.2)
```

推荐阅读

群落多样性之Beta多样性（三）

阅读 314

R 稳健PCA分析

阅读 165

R数据分析：主成分分析及可视化

阅读 717

基于R的线性混合效应模型分析

阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）

阅读 448

写下你的评论...

评论0

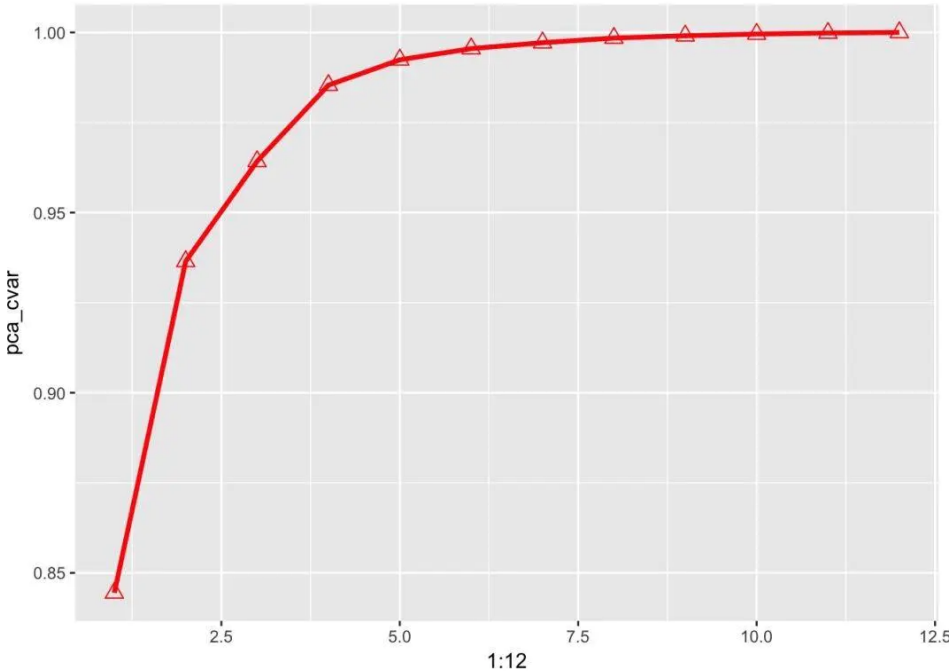
赞35

...



image

```
p1+geom_point(pch=2,lwd=3,col=2)+geom_line(col=2,lwd=1.2)
```



image

推荐阅读

群落多样性之Beta多样性（三）

阅读 314

R 稳健PCA分析

阅读 165

R数据分析：主成分分析及可视化

阅读 717

基于R的线性混合效应模型分析

阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）

阅读 448

崖低碎石图（scree plot）即贡献率图，是希望图形一开始很陡峭，如悬崖一般，而剩下的数值都很小，如崖底的碎石一样。

崖低碎石图和累积贡献率图是对主成分贡献率和累积贡献率的一种直观表示，用以作为选择主成分个数的参考。本例中第一个主成解释总变异的84.4%，可以只选择第一个主成分，但第二主成分也不小，因此选择前两个主成分。

主成分的个数选择没有一定之规，需按实际情况具体分析，一般要求累积贡献率大于85%或特征值大于1。

但是在实际的生物信息学应用中，通常达不到这个要求。

2.7 step5:主成分载荷

主成分载荷表示各个主成分与原始变量的相关系数。

```
pca_vect= dat_eigen$vector ## 相关系数矩阵的特征向量
loadings=sweep(pca_vect,2,sqrt(pca_var),"*")
```

写下你的评论...

评论0

赞35

...



CONT	0.0028	0.2830
INTG	-0.2652	-0.0552
DMNR	-0.2636	-0.0599
DILG	-0.2797	0.0110
CFMG	-0.2780	0.0511
DECI	-0.2774	0.0388
PREP	-0.2843	0.0098
FAMI	-0.2818	-0.0004
ORAL	-0.2874	-0.0011
WRIT	-0.2858	-0.0095
PHYS	-0.2580	0.0270
RTEN	-0.2847	-0.0119

推荐阅读

群落多样性之Beta多样性（三）
阅读 314

R 稳健PCA分析
阅读 165

R数据分析：主成分分析及可视化
阅读 717

基于R的线性混合效应模型分析
阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）
阅读 448

结果表明，CONT指标跟其它指标表现完全不一样，第一个主成分很明显跟除了CONT之外的所有其它指标负相关，而第二个主成分则主要取决于CONT指标。

2.8 step6:主成分得分计算和图示

将中心化的变量dat_scale乘以特征向量矩阵即得到每个观测值的得分。

```
pca_score=as.matrix(dat_scale)%*%pca_vect
head(pca_score[,1:2])
```

```
##           [,1]      [,2]
## AARONSON,L.H.  0.5098 -1.7080
## ALEXANDER,J.M. -2.2676 -0.8508
## ARMENTANO,A.J. -0.2267 -0.2903
## BERDON,R.I.   -3.4058 -0.5997
## BRACKEN,J.J.   6.5937  0.2478
## BURNS,E.B.    -2.3336 -1.3563
```

将两个主成分以散点图形式展示,看看这些样本被这两个主成分是如何分开的

```
p=ggplot(aes(x=pca_score[,1],y=pca_score[,2]))+geom_point(color=USJudgeRatings[,1],pch=USJudgeRati
p
```

https://mmbiz.qpic.cn/mmbiz_png/cZNhZQ6j4wytdsnNenoKyYFGJJUTcgm3Eicuqs0GA0dyHv5mXZ0yibbrEfrMxjT9NLC9LrRsX8WQcNXeBYaJicDgg/640?wx_fmt=png&wxfrom=5&wx_lazy=1&wx_co=1

对于主成分分析，不同数据会有不同的分析方法，应具体情况具体分析。

总结一下PCA的算法步骤：



写下你的评论...

评论0

赞35

...

主成分分析大全



生信技能树

关注

赞赏支持

- 3) 求出协方差矩阵
- 4) 求出协方差矩阵的特征值及对应的特征向量
- 5) 将特征向量按对应特征值大小从上到下按行排列成矩阵，取前k行组成矩阵P
- 6) $Y=PX$ 即为降维到k维后的数据

PCA本质上是将方差最大的方向作为主要特征，并且在各个正交方向上将数据“离相关”，也就是让它们在不同正交方向上没有相关性。

PCA也存在一些限制，例如它可以很好的解除线性相关，但是对于高阶相关性就没有办法了，对于存在高阶相关性的数据，可以考虑Kernel PCA，通过Kernel函数将非线性相关转为线性相关，关于这点就不展开讨论了。另外，PCA假设数据各主特征是分布在正交方向上，如果在非正交方向上存在几个方差较大的方向，PCA的效果就大打折扣了。

最后需要说明的是，PCA是一种无参数技术，也就是说面对同样的数据，如果不考虑清洗，谁来做结果都一样，没有主观参数的介入，所以PCA便于通用实现，但是本身无法个性化的优化。

3 实战一

比如你要做一项分析人的糖尿病的因素有哪些，这时你设计了10个你觉得都很重要的指标，然而这10个指标对于你的分析确实太过繁杂，这时你就可以采用主成分分析的方法进行降维。10个指标之间会有这样那样的联系，相互之间会有影响，通过主成分分析后，得到三五个主成分指标。此时这几个主成分指标既涵盖了你10个指标中的绝大部分信息，这让你的分析得到了简化（从10维降到3、5维）。

下面是442个糖尿病人相关的数据集，具体如下：

- x a matrix with 10 columns (自变量)
- y a numeric vector (因变量)

```
library(lars)
library(glmnet)
data(diabetes)
attach(diabetes)
summary(x)
```

```
##      age          sex          bmi
##  Min.   :-0.10723   Min.   :-0.0446   Min.   :-0.09028
##  1st Qu.: -0.03730   1st Qu.: -0.0446   1st Qu.: -0.03423
##  Median : 0.00538   Median : -0.0446   Median : -0.00728
##  Mean   : 0.00000   Mean    : 0.0000   Mean    : 0.00000
##  3rd Qu.: 0.03808   3rd Qu.: 0.0507   3rd Qu.: 0.03125
##  Max.   : 0.11073   Max.    : 0.0507   Max.    : 0.17056
##      map          tc          ldl
##  Min.   :-0.11240   Min.   :-0.12678   Min.   :-0.11561
##  1st Qu.: -0.03666   1st Qu.: -0.03425   1st Qu.: -0.03036
##  Median : -0.00567   Median : -0.00432   Median : -0.00382
##  Mean   : 0.00000   Mean    : 0.00000   Mean    : 0.00000
##  3rd Qu.: 0.03564   3rd Qu.: 0.02836   3rd Qu.: 0.02984
##  Max.   : 0.13204   Max.    : 0.15201   Max.    : 0.10770
```

写下你的评论...

评论0

赞35

...

推荐阅读

群落多样性之Beta多样性（三）
阅读 314

R 稳健PCA分析
阅读 165

R数据分析：主成分分析及可视化
阅读 717

基于R的线性混合效应模型分析
阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）
阅读 448

主成分分析大全



生信技能树

关注

赞赏支持

```
## 3rd Qu.: 0.02931 3rd Qu.: 0.03431 3rd Qu.: 0.03243
## Max. : 0.18118 Max. : 0.18523 Max. : 0.13360
## glu
## Min. : -0.13777
## 1st Qu.: -0.03318
## Median : -0.00108
## Mean : 0.00000
## 3rd Qu.: 0.02792
## Max. : 0.13561
```

dim(x)

```
## [1] 442 10
```

length(y)

```
## [1] 442
```

summary(y)

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 25 87 140 152 212 346
```

boxplot(y)

推荐阅读

群落多样性之Beta多样性（三）

阅读 314

R 稳健PCA分析

阅读 165

R数据分析：主成分分析及可视化

阅读 717

基于R的线性混合效应模型分析

阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）

阅读 448



写下你的评论...



评论0



赞35



主成分分析大全

其中x是包含10个变量，分别为：age sex bmi map tc tgl hdl ldl tg glu 它们都在一定程度上或多或少的会影响个体糖尿病状态。

数据的详细介绍见 [Efron, Hastie, Johnstone and Tibshirani \(2003\) "Least Angle Regression" \(with discussion\) Annals of Statistics;](#)

一步法主成分分析

上面我们把整个主成分分析步骤拆解开来讲解具体原理，但是实际数据处理过程中我们通常是用现成的函数一步法完成主成分分析，而且这个是非常高频的分析，所以R里面自带了一个函数 `princomp()` 来完成主成分分析，如下：

```
data=x ## 这里的x是上面的 diabetes 数据集里面的 442个病人的10个生理指标
pca<-princomp(data, cor=FALSE)
```

cor是逻辑变量,当cor=TRUE表示用样本的相关矩阵R做主成分分析,当cor=FALSE表示用样本的协方差阵S做主成分分析。

```
summary(pca)

## Importance of components:
##              Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6
## Standard deviation  0.09542 0.05811 0.05223 0.04649 0.03871 0.03693
## Proportion of Variance 0.40242 0.14923 0.12060 0.09555 0.06622 0.06027
## Cumulative Proportion 0.40242 0.55165 0.67225 0.76780 0.83402 0.89429
##              Comp.7  Comp.8  Comp.9  Comp.10
## Standard deviation  0.03484 0.03132 0.01331 0.00440
## Proportion of Variance 0.05366 0.04337 0.00783 0.00085
## Cumulative Proportion 0.94794 0.99131 0.99914 1.00000
```

可以看到前三个主成份的信息量也只有67.2%，达不到我们前面说到85%，所以很难说可以用这三个主成分去代替这10个生理指标来量化病人的状态。

值得一提的是，如果你看懂了前面的主成分分析的拆解步骤，就应该明白有多少个变量就有多少个主成分，只是并不是所有的主成分都有意义，理想状态下我们希望有限的几个主成分就可以代替数量居多的变量，尤其是生物信息学里面的基因表达矩阵，两三万个基因的表达情况我们希望几十个基因就可以替代它们，因为那些基因之间是相互关联的。

碎石图

也可以画出主成分的碎石图，来决定选几个主成分。

```
screeplot(pca, type='lines')
```



生信技能树

关注

赞赏支持

推荐阅读

- 群落多样性之Beta多样性（三）
阅读 314
- R 稳健PCA分析
阅读 165
- R数据分析：主成分分析及可视化
阅读 717
- 基于R的线性混合效应模型分析
阅读 518
- 无监督学习-线性方法|深度学习（李宏毅）（十七）
阅读 448

写下你的评论...

评论0 赞35 ...



推荐阅读

群落多样性之Beta多样性（三）

阅读 314

R 稳健PCA分析

阅读 165

R数据分析：主成分分析及可视化

阅读 717

基于R的线性混合效应模型分析

阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）

阅读 448

image

由碎石图可以看出，第5个主成分之后，图线变化趋于平稳，因此可以选择前5个主成分做分析。

样本分布的散点图

根据前两个主成分画出样本分布的散点图。

```
| biplot(pca)
```

image

上面这个图似乎意义不大，因为大部分情况下都是需要结合样本的分组信息来看看这些主成分是否可以把样本比较不错的分开。



写下你的评论...

评论0

赞35



Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6	Comp.7	Comp.8	Comp.9	Comp.10
-0.0279	-0.0926	0.0280	0.0039	0.0122	0.0481	-0.0086	-0.0360	-0.0086	-0.0023
0.1347	0.0653	0.0013	0.0224	0.0068	0.0482	0.0107	0.0090	0.0240	0.0021
-0.0129	-0.0778	0.0352	0.0376	0.0554	0.0529	-0.0220	-0.0401	-0.0012	-0.0026
-0.0023	0.0182	-0.0958	-0.0653	-0.0122	-0.0212	0.0229	0.0175	-0.0066	-0.0035

kable(data[1:4,])

age	sex	bmi	map	tc	ldl	hdl	tch	ltg	glu
0.038076	0.050680	0.061696	0.021872	-0.044223	-0.034821	-0.043401	-0.002592	0.019908	-0.017646
-0.001882	-0.044642	-0.051474	-0.026328	-0.008449	-0.019163	0.074412	-0.039493	-0.068330	-0.092204
0.085299	0.050680	0.044451	-0.005671	-0.045599	-0.034194	-0.032356	-0.002592	0.002864	-0.025930
-0.089063	-0.044642	-0.011595	-0.036656	0.012191	0.024991	-0.036038	0.034309	0.022692	-0.009362

预测主成份的值，这里用的就是训练数据，所以得出训练数据主成分的值。

4 实战二

R中自带数据集 `data(Harman23.cor)` 数据集中包含305名受试者的8个身体测量指标

```
data(Harman23.cor)
kable(Harman23.cor[1:5])

## Warning in kable_markdown(x = structure(c("0", "0", "0", "0", "0", "0", :
## The table should have a header (column names)

## Warning in kable_markdown(x = structure("305", .Dim = c(1L, 1L), .Dimnames
## = list(: The table should have a header (column names)
```

5 进阶的主成分分析-psych包

正文中的`princomp()`函数为基础包中进行主成分分析的函数。另外，R中`psych`包中提供了一些更加丰富有用的函数，这里列出几个相关度较高的函数，以供读者了解。

主成分分析大全

principal()	含多种可选的力左旋转力法的主成分力力
fa()	可用主轴、最小残差、加权最小平方或最大似然法估计的因子分析
fa.parallel()	含平行分析的碎石图
factor.plot()	绘制因子分析或主成分分析的结果
fa.diagram()	绘制因子分析或主成分的载荷矩阵
scree()	因子分析和主成分分析的碎石图

还有很多主成分分析结果可视化包，在直播我的基因组里面都提到过。

6 再推荐一个R包factoextra

factoextra是一个R包，易于提取和可视化探索性多变量数据分析的输出，包括：

- 主成分分析（PCA），用于通过在不丢失重要信息的情况下降低数据的维度来总结连续（即定量）多变量数据中包含的信息。
- 对应分析（CA）是适用于分析由两个定性变量（或分类数据）形成的大型应变表的主成分分析的扩展。
- 多重对应分析（MCA），它是CA对包含两个以上分类变量的数据表的适应。
- 多因素分析（MFA）专门用于数据集，其中变量被组织成组（定性和/或定量变量）。
- 分层多因素分析（HMFA）：在将数据组织成层次结构的情况下，MFA的扩展。
- 混合数据因子分析（FAMD）是MFA的一个特例，专门用于分析包含定量和定性变量的数据集。

有许多R包实现主要组件方法。这些包包括：FactoMineR，ade4，stats，ca，MASS和ExPosition。然而，根据使用的包，结果呈现不同。为了帮助解释和多变量分析的可视化（如聚类分析和维数降低分析），所以作者开发了一个名为factoextra的易于使用的R包。

7.主成分分析的生物信息学应用

比如对千人基因组计划的对VCF突变数据进行主成分分析来区分人种：

<https://www.biostars.org/p/185869/>

8. 主成分分析的其它可视化方法

看到一个包 `ropls` 可视化做的不错，本来以为 `ropls` 肯定是一个正常的r包，应该是在cran里面，结果

```
> install.packages('ropls')
Warning in install.packages :
  package 'ropls' is not available (for R version 3.4.3)
Warning in install.packages :
  cannot open URL 'https://cran.rstudio.com/bin/macosx/el-capitan/contrib/3.4/PACKAGES.rds': HTTP s
> BiocInstaller::biocLite('ropls')
BioC_mirror: https://bioconductor.org
Using Bioconductor 3.6 (BiocInstaller 1.28.0), R 3.4.3 (2017-11-30).
Installing package(s) 'ropls'
trying URL 'https://bioconductor.org/packages/3.6/bioc/bin/macosx/el-capitan/contrib/3.4/ropls_1.10
Content type 'application/x-gzip' length 122650 bytes (120 KB)
```



推荐阅读

群落多样性之Beta多样性（三）
阅读 314

R 稳健PCA分析
阅读 165

R数据分析：主成分分析及可视化
阅读 717

基于R的线性混合效应模型分析
阅读 518

无监督学习-线性方法|深度学习（李宏毅）（十七）
阅读 448



后来仔细看标题，终于明白了。

ropis: PCA, PLS(-DA) and OPLS(-DA) for multivariate analysis and feature selection of omics data

构建就是组学数据

9.参考资料：

- http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf
- https://www.cs.princeton.edu/picasso/mats/PCA-Tutorial-Intuition_jp.pdf
- <http://www.cs.umd.edu/~samir/498/PCA.pdf>
- http://www.yale.edu/ceo/Documentation/PCA_Outline.pdf
- <http://people.tamu.edu/~alawing/materials/ESSM689/pca.pdf>（R相关）
- <http://www2.dc.ufscar.br/~cesar.souza/publications/pca-tutorial.pdf>（2012）



35人点赞 >



孟孟



更多精彩内容，就在简书APP



"小礼物走一走，来简书关注我"

赞赏支持

还没有人赞赏，支持一下



生信技能树 生信技能树，生信菜鸟团，jimmy
总资产153 (约10.91元) 共写了59.5W字 获得3,429个赞 共6,878个粉丝

关注

张家界之三日游，张家界风景你看了吗



写下你的评论...



评论0



赞35





收入即学习



R语言



Web生物数据分析



生物信息学



生物信息学



数据-R语言...



生物信息学与算法

展开更多

推荐阅读

R语言主成分和因子分析篇

转载自 R语言主成分和因子分析篇另可参考 R in action读书笔记 (19) 第十四章 主成分和因子分析 主成分...



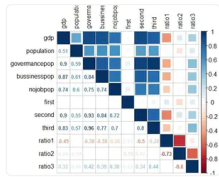
风雨如晦 阅读 2,221 评论 0 赞 6

应用统计学与R语言实现学习笔记 (十二) ——主成分分析

Chapter 12 Principle Component Analysis 本篇是第十二章,内容是主成分分析。这...



G小调的Qing歌 阅读 4,667 评论 0 赞 14

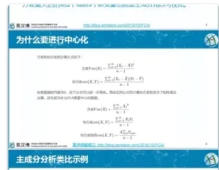


PCA主成分分析实战和可视化 | 附R代码和测试数据

一文看懂PCA主成分分析中介绍了PCA分析的原理和分析的意义(基本简介如下,更多见博客),今天就用数据来实际操作一...



生信宝典 阅读 17,740 评论 11 赞 61



【PCA-1】主成分分析

欢迎关注公众号: oddxix 主成分分析简介 主成分分析 (PCA, principal component an...



oddxix 阅读 2,950 评论 1 赞 25

	T16_131100	T16_131101	T16_131102	T16_131103	T16_131104
1	T16_131100	T16_131101	T16_131102	T16_131103	T16_131104
2	T16_131101	T16_131102	T16_131103	T16_131104	T16_131105
3	T16_131102	T16_131103	T16_131104	T16_131105	T16_131106
4	T16_131103	T16_131104	T16_131105	T16_131106	T16_131107
5	T16_131104	T16_131105	T16_131106	T16_131107	T16_131108
6	T16_131105	T16_131106	T16_131107	T16_131108	T16_131109
7	T16_131106	T16_131107	T16_131108	T16_131109	T16_131110
8	T16_131107	T16_131108	T16_131109	T16_131110	T16_131111
9	T16_131108	T16_131109	T16_131110	T16_131111	T16_131112
10	T16_131109	T16_131110	T16_131111	T16_131112	T16_131113
11	T16_131110	T16_131111	T16_131112	T16_131113	T16_131114
12	T16_131111	T16_131112	T16_131113	T16_131114	T16_131115
13	T16_131112	T16_131113	T16_131114	T16_131115	T16_131116
14	T16_131113	T16_131114	T16_131115	T16_131116	T16_131117
15	T16_131114	T16_131115	T16_131116	T16_131117	T16_131118
16	T16_131115	T16_131116	T16_131117	T16_131118	T16_131119
17	T16_131116	T16_131117	T16_131118	T16_131119	T16_131120
18	T16_131117	T16_131118	T16_131119	T16_131120	T16_131121
19	T16_131118	T16_131119	T16_131120	T16_131121	T16_131122
20	T16_131119	T16_131120	T16_131121	T16_131122	T16_131123
21	T16_131120	T16_131121	T16_131122	T16_131123	T16_131124
22	T16_131121	T16_131122	T16_131123	T16_131124	T16_131125
23	T16_131122	T16_131123	T16_131124	T16_131125	T16_131126
24	T16_131123	T16_131124	T16_131125	T16_131126	T16_131127
25	T16_131124	T16_131125	T16_131126	T16_131127	T16_131128
26	T16_131125	T16_131126	T16_131127	T16_131128	T16_131129
27	T16_131126	T16_131127	T16_131128	T16_131129	T16_131130
28	T16_131127	T16_131128	T16_131129	T16_131130	T16_131131
29	T16_131128	T16_131129	T16_131130	T16_131131	T16_131132
30	T16_131129	T16_131130	T16_131131	T16_131132	T16_131133
31	T16_131130	T16_131131	T16_131132	T16_131133	T16_131134
32	T16_131131	T16_131132	T16_131133	T16_131134	T16_131135
33	T16_131132	T16_131133	T16_131134	T16_131135	T16_131136
34	T16_131133	T16_131134	T16_131135	T16_131136	T16_131137
35	T16_131134	T16_131135	T16_131136	T16_131137	T16_131138
36	T16_131135	T16_131136	T16_131137	T16_131138	T16_131139
37	T16_131136	T16_131137	T16_131138	T16_131139	T16_131140
38	T16_131137	T16_131138	T16_131139	T16_131140	T16_131141
39	T16_131138	T16_131139	T16_131140	T16_131141	T16_131142
40	T16_131139	T16_131140	T16_131141	T16_131142	T16_131143
41	T16_131140	T16_131141	T16_131142	T16_131143	T16_131144
42	T16_131141	T16_131142	T16_131143	T16_131144	T16_131145
43	T16_131142	T16_131143	T16_131144	T16_131145	T16_131146
44	T16_131143	T16_131144	T16_131145	T16_131146	T16_131147
45	T16_131144	T16_131145	T16_131146	T16_131147	T16_131148
46	T16_131145	T16_131146	T16_131147	T16_131148	T16_131149
47	T16_131146	T16_131147	T16_131148	T16_131149	T16_131150
48	T16_131147	T16_131148	T16_131149	T16_131150	T16_131151
49	T16_131148	T16_131149	T16_131150	T16_131151	T16_131152
50	T16_131149	T16_131150	T16_131151	T16_131152	T16_131153
51	T16_131150	T16_131151	T16_131152	T16_131153	T16_131154
52	T16_131151	T16_131152	T16_131153	T16_131154	T16_131155
53	T16_131152	T16_131153	T16_131154	T16_131155	T16_131156
54	T16_131153	T16_131154	T16_131155	T16_131156	T16_131157
55	T16_131154	T16_131155	T16_131156	T16_131157	T16_131158
56	T16_131155	T16_131156	T16_131157	T16_131158	T16_131159
57	T16_131156	T16_131157	T16_131158	T16_131159	T16_131160
58	T16_131157	T16_131158	T16_131159	T16_131160	T16_131161
59	T16_131158	T16_131159	T16_131160	T16_131161	T16_131162
60	T16_131159	T16_131160	T16_131161	T16_131162	T16_131163
61	T16_131160	T16_131161	T16_131162	T16_131163	T16_131164
62	T16_131161	T16_131162	T16_131163	T16_131164	T16_131165
63	T16_131162	T16_131163	T16_131164	T16_131165	T16_131166
64	T16_131163	T16_131164	T16_131165	T16_131166	T16_131167
65	T16_131164	T16_131165	T16_131166	T16_131167	T16_131168
66	T16_131165	T16_131166	T16_131167	T16_131168	T16_131169
67	T16_131166	T16_131167	T16_131168	T16_131169	T16_131170
68	T16_131167	T16_131168	T16_131169	T16_131170	T16_131171
69	T16_131168	T16_131169	T16_131170	T16_131171	T16_131172
70	T16_131169	T16_131170	T16_131171	T16_131172	T16_131173
71	T16_131170	T16_131171	T16_131172	T16_131173	T16_131174
72	T16_131171	T16_131172	T16_131173	T16_131174	T16_131175
73	T16_131172	T16_131173	T16_131174	T16_131175	T16_131176
74	T16_131173	T16_131174	T16_131175	T16_131176	T16_131177
75	T16_131174	T16_131175	T16_131176	T16_131177	T16_131178
76	T16_131175	T16_131176	T16_131177	T16_131178	T16_131179
77	T16_131176	T16_131177	T16_131178	T16_131179	T16_131180
78	T16_131177	T16_131178	T16_131179	T16_131180	T16_131181
79	T16_131178	T16_131179	T16_131180	T16_131181	T16_131182
80	T16_131179	T16_131180	T16_131181	T16_131182	T16_131183
81	T16_131180	T16_131181	T16_131182	T16_131183	T16_131184
82	T16_131181	T16_131182	T16_131183	T16_131184	T16_131185
83	T16_131182	T16_131183	T16_131184	T16_131185	T16_131186
84	T16_131183	T16_131184	T16_131185	T16_131186	T16_131187
85	T16_131184	T16_131185	T16_131186	T16_131187	T16_131188
86	T16_131185	T16_131186	T16_131187	T16_131188	T16_131189
87	T16_131186	T16_131187	T16_131188	T16_131189	T16_131190
88	T16_131187	T16_131188	T16_131189	T16_131190	T16_131191
89	T16_131188	T16_131189	T16_131190	T16_131191	T16_131192
90	T16_131189	T16_131190	T16_131191	T16_131192	T16_131193
91	T16_131190	T16_131191	T16_131192	T16_131193	T16_131194
92	T16_131191	T16_131192	T16_131193	T16_131194	T16_131195
93	T16_131192	T16_131193	T16_131194	T16_131195	T16_131196
94	T16_131193	T16_131194	T16_131195	T16_131196	T16_131197
95	T16_131194	T16_131195	T16_131196	T16_131197	T16_131198
96	T16_131195	T16_131196	T16_131197	T16_131198	T16_131199
97	T16_131196	T16_131197	T16_131198	T16_131199	T16_131200
98	T16_131197	T16_131198	T16_131199	T16_131200	T16_131201
99	T16_131198	T16_131199	T16_131200	T16_131201	T16_131202
100	T16_131199	T16_131200	T16_131201	T16_131202	T16_131203

写给女儿的一封信

宝贝你好,你即将要过四岁生日了,我知道你期待这一天很久了,你经常在家里问我和妈妈,你什么时候过生日,可以邀请你的同...



自律就是自由 阅读 163 评论 0 赞 1

