

IR実務者のためのR入門



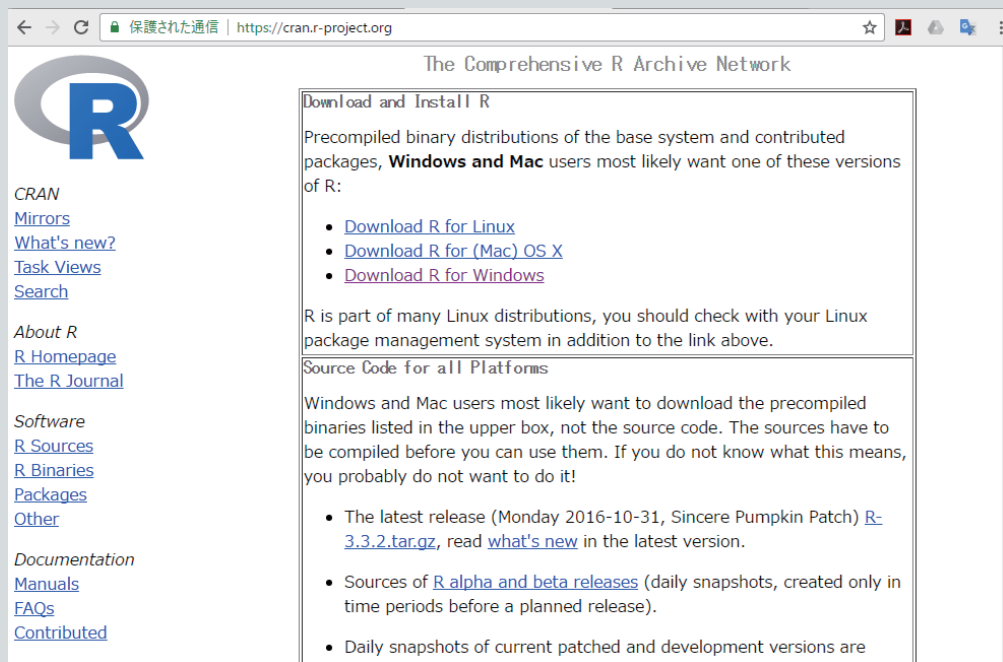
山本義郎（東海大学）

日本計算機統計学会

本日の内容

- RとRstudioのインストール
- Rを使ったら何ができるか
- Rの基本的な操作
- Rを使う環境
 - Rコマンダー（コマンドを知らなくてもOK）
 - Rstudio（Rを使いこなすなら是非）
- 実データの解析例

Rのインストール



CRAN
Rstudio を検索

IR実務者向けのR入門

2017/2/10

3

Rを使ったら何ができるか(1)

• データ解析



Excelから



Rへ

Excelが得意なこと

- ・集計(ピボットテーブル)
- ・グラフ(棒、折れ線、円グラフ)
- ・関数を使ったデータの変換
- ・フィルター
- ・データベース関数
- ...

Rが得意なこと

- ・大量データの処理
- ・データ解析(統計解析)
- ・グラフ・統計グラフ
- ・GIS(地理情報処理)
- ・プログラミング
- ・シミュレーション
- ・ルーチンワーク

IR実務者向けのR入門

2017/2/10

4

Rを使ったら何ができるか(2)

- レポーティング



Officeから



Rへ

- WebページやPDFファイルに

- インタラクティブな機能も

EXCELよりRを使うべき理由

1. 大量データの処理
2. データ解析(統計解析)
3. グラフ・統計グラフ
4. GIS(地理情報処理)
5. プログラミング
6. シミュレーション
7. ルーチンワーク

Rを使うべき理由(1)

1. 大量データの処理

- Excelのデータも読み込める
- データベースからのデータのインポート
- データの抽出(フィルタ)などは縦横無尽

Rを使うべき理由(2)

2. データ解析(統計解析)

- 検定、分散分析、多重比較
- 統計データの視覚化(様々な統計グラフ)
- 多変量データ解析
 - 回帰分析、クラスター分析、主成分分析、、
- データマイニング
 - アソシエーションルール分析、決定木分析、、

Rを使うべき理由(3)

3. グラフ・統計グラフ

- 統計グラフ
 - 箱ひげ図、散布図行列、、、
- 統計処理のグラフ
 - 回帰直線、デンドログラム、対応分析、、、
- 大量のグラフ作成もプログラミングで容易に
- インタラクティブなグラフ作成も可能

Rを使うべき理由(4)

4. GIS(地理情報処理)

- 地理情報システム並の処理ができる
 - 塗り分け地図(コロプレスマップ)
- GoogleMapやGoogleEarthの利用も
- QGISなどのGISソフトウェアとの連携も

Rを使うべき理由(5)

5. プログラミング

- CやJavaよりも気軽にプログラミングができる
- シミュレーション
- ルーチンワークにマクロとして使える

Rを使うべき理由(6)

6. シミュレーション

- プログラミング言語でのルーチンワーク

Rを使うべき理由(7)

7. ルーチンワーク

- 定形処理が得意

Rの起動と終了



スタートメニュー（デスクトップ）から、「R x64 3.3.2」を起動

```
RGui (64-bit)
ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ

R Console

R version 3.2.2 (2015-08-14) -- "Fire Safety"
Copyright (C) 2015 The R Foundation for Statistical Computing
Platform: x86_64-mingw32/x64 (64-bit)

R は、自由なソフトウェアであり「完全に無保証」です。
一定の条件に従えば、自由にこれを再配布することができます。
配布条件の詳細に関しては、'license()' あるいは 'licence()' と入力してください。

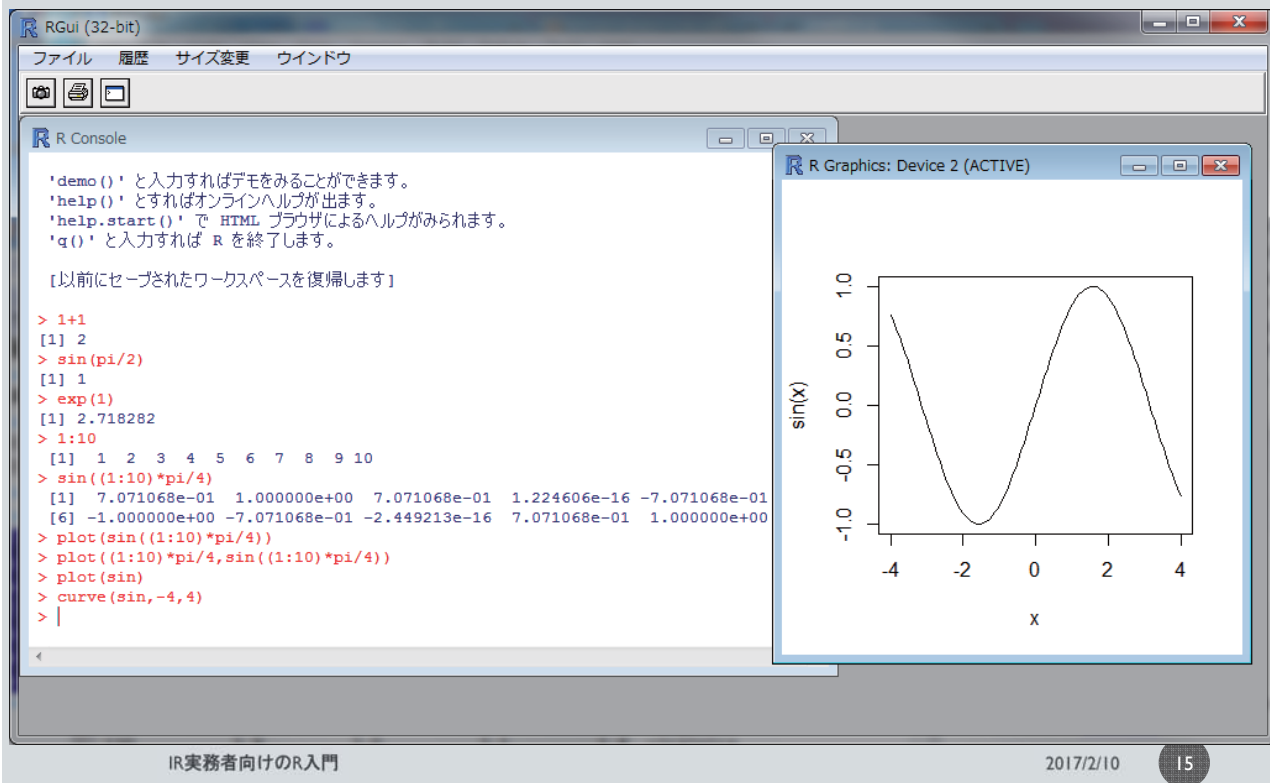
R は多くの貢献者による共同プロジェクトです。
詳しくは 'contributors()' と入力してください。
また、R や R のパッケージを出版物で引用する際の形式については
'citation()' と入力してください。

'demo()' と入力すればデモをみることができます。
'help()' とすればオンラインヘルプが出ます。
'help.start()' で HTML ブラウザによるヘルプがみられます。
'q()' と入力すれば R を終了します。

> 1+2
[1] 3
> q()
```

終了命令

Rの基本(1) 関数電卓・グラフ

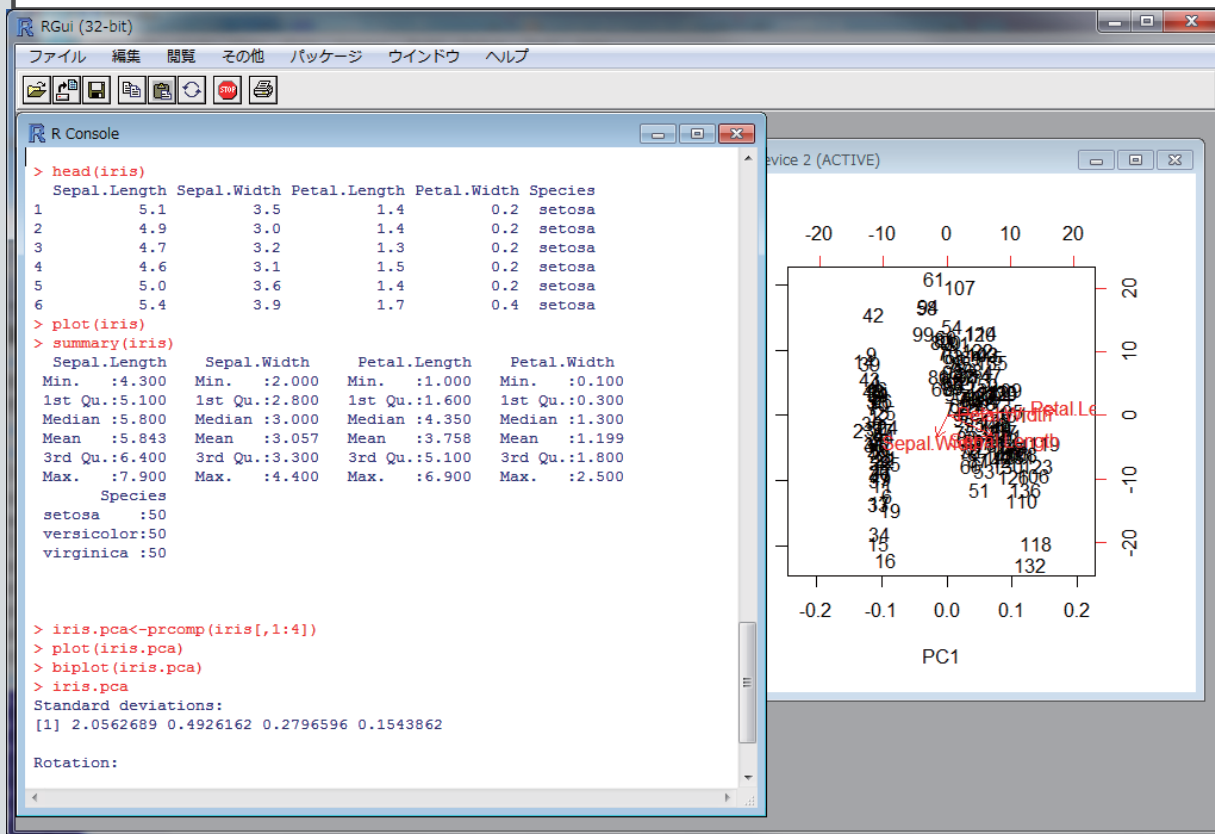


Rの基本(2)統計処理・データ解析

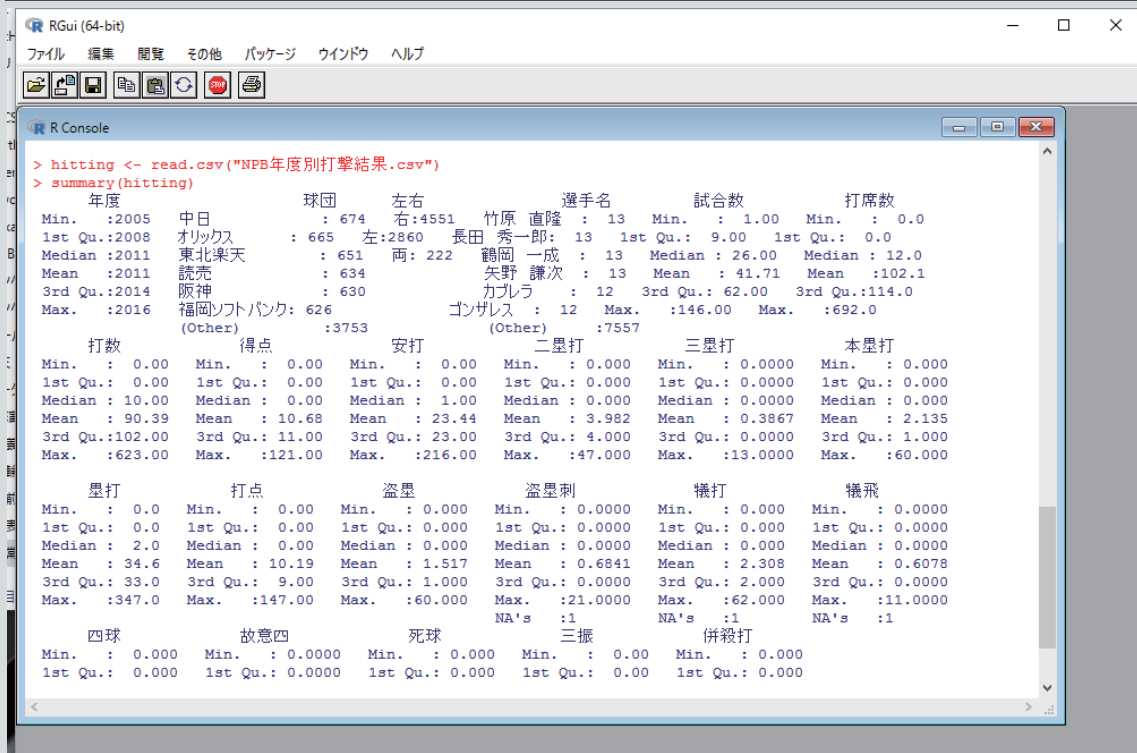
```
R Console

> 1+2
[1] 3
> mean(c(148, 160, 159, 153, 151, 140, 158, 137, 149, 160))
[1] 151.5
> height <- c(148, 160, 159, 153, 151, 140, 158, 137, 149, 160)
> height
[1] 148 160 159 153 151 140 158 137 149 160
> mean(height)
[1] 151.5
> sex <- c("F", "M", "M", "F", "F", "F", "M", "F", "M", "F")
> tapply(height, sex, mean)
      F      M
148.1667 156.5000
> mydata <- data.frame(height, sex)
> mydata
  height sex
1    148  F
2    160  M
3    159  M
4    153  F
5    151  F
6    140  F
7    158  M
8    137  F
9    149  M
10   160  F
> mydata$height
[1] 148 160 159 153 151 140 158 137 149 160
> mean(mydata$height)
[1] 151.5
> |
```


Rの基本(2) 統計処理・データ解析

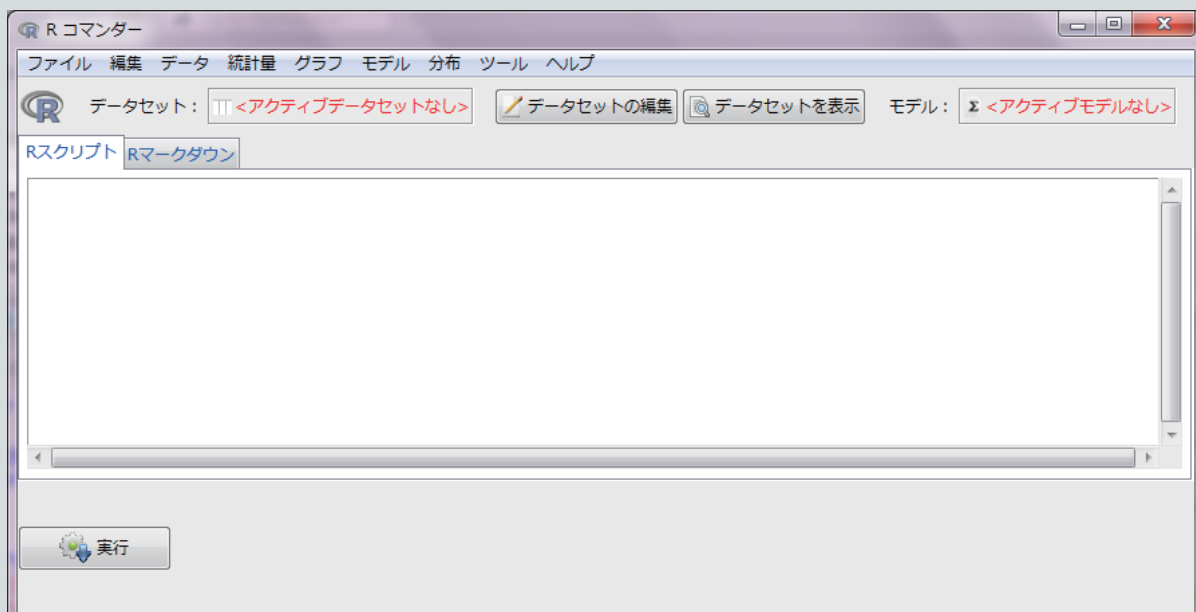


Rの基本(3) データの読み込み



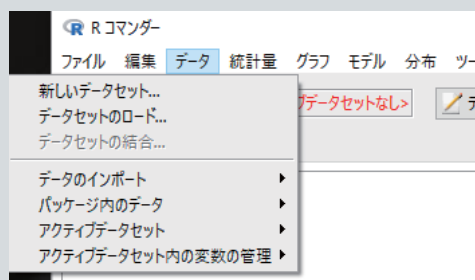
Rコマンダー

- Rcmrパッケージ
- コマンドを覚えなくてもメニューで使える

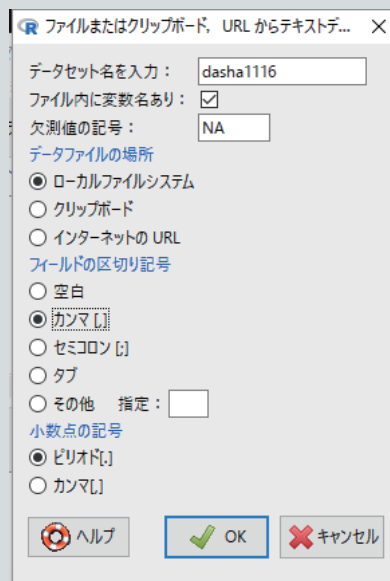


Rコマンダーでできること(1)

- データの作成

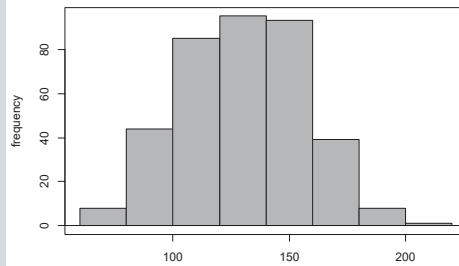
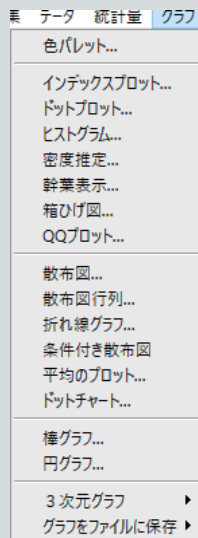


インポート

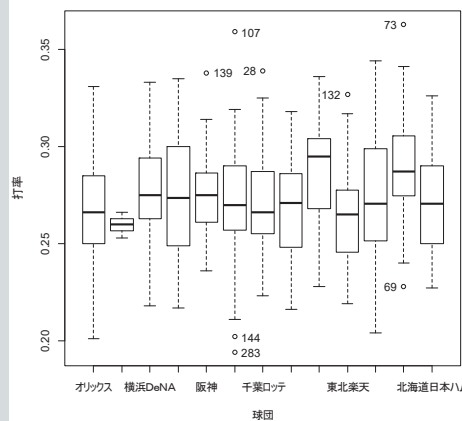


Rコマンダーでできること(2)

• グラフ



ヒストグラム

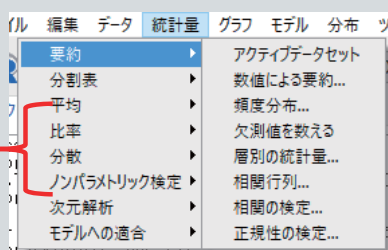


箱ひげ図

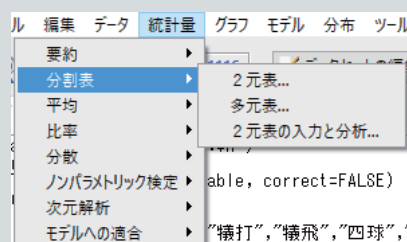
Rコマンダーでできること(3)

• 記述統計・検定

検定



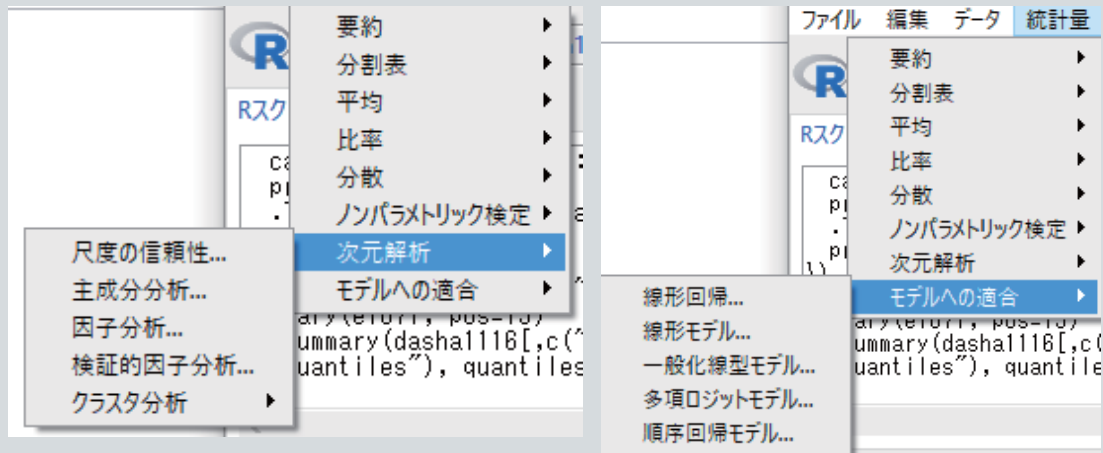
量的変数の要約



質的変数の要約

Rコマンダーでできること(4)

• 多変量データ解析

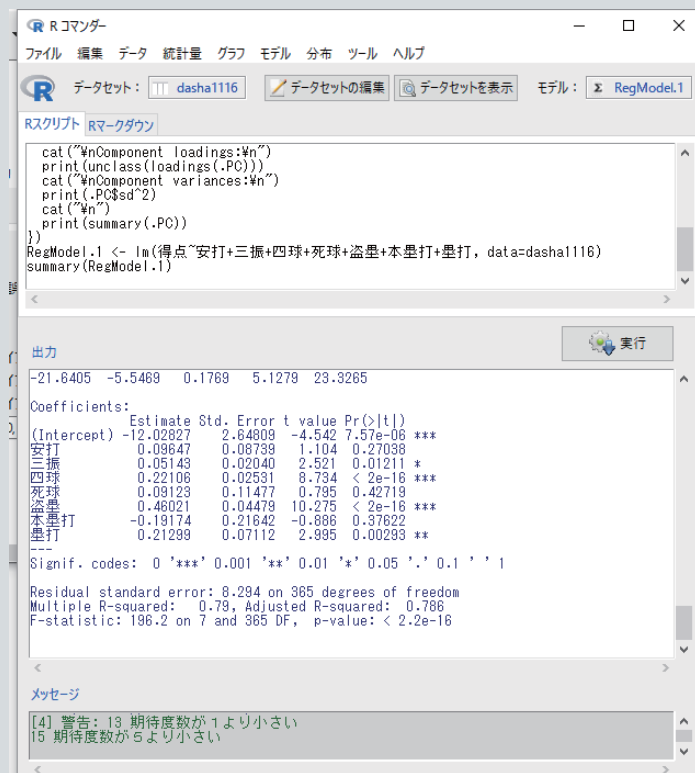


分類（クラスタリング）
次元縮約

予測

Rコマンダーでできること(4)

• 線形重回帰



Rコマンダーでできること(4)

線形重回帰モデル

The screenshot shows the R Commander interface. The 'Model' menu is open, displaying options for selecting a model, comparing coefficients, and saving results. The 'Model' field shows 'RegModel.1'. The output window displays the results of a linear regression model, including the formula, coefficients, and a table of statistics.

Model: Σ RegModel.1

Call:
lm(formula = 得点 ~ 安打 + 三振 + 四球 + 盗塁 + 塁打, data = dashall116)

Coefficients:
(Intercept) 安打 三振 四球 盗塁 塁打
-11.88111 0.16795 0.05315 0.22149 0.46463 0.15332

	Df	Sum of Sq	RSS	AIC
<none>			25200	1807.0
- 三振	1	470.2	25670	1808.0
+ 本塁打	1	45.9	25154	1812.2
+ 死球	1	35.3	25184	1812.4
- 安打	1	1801.3	27001	1826.8
- 塁打	1	4374.8	29574	1860.8
- 四球	1	5272.0	30472	1871.9
- 盗塁	1	7464.5	32664	1897.8

Rコマンダーでできること(5)

レポートの作成(Rマークダウン)

The screenshot shows the R Commander interface. The 'R Markdown' menu is open, displaying options for creating a new R Markdown file, opening an existing one, and saving. The output window displays the results of a linear regression model, including the formula, coefficients, and a table of statistics.

Model: Σ RegModel.1

Call:
lm(formula = 得点 ~ 安打 + 三振 + 四球 + 盗塁 + 塁打, data = dashall116)

Coefficients:
(Intercept) 安打 三振 四球 盗塁 塁打
-11.88111 0.16795 0.05315 0.22149 0.46463 0.15332

	Df	Sum of Sq	RSS	AIC
<none>			25200	1807.0
- 三振	1	470.2	25670	1808.0
+ 本塁打	1	45.9	25154	1812.2
+ 死球	1	35.3	25184	1812.4
- 安打	1	1801.3	27001	1826.8
- 塁打	1	4374.8	29574	1860.8
- 四球	1	5272.0	30472	1871.9
- 盗塁	1	7464.5	32664	1897.8

データマイニング（目的と手法）

目的	手法	適用例(マーケティング)
予測 (判別予測 と数値予測)	<ul style="list-style-type: none"> ・決定木 ・ニューラルネット ・ロジスティック回帰 ・重回帰分析(線形・非線形) ・判別分析 	<ul style="list-style-type: none"> ・スコアリング ・チャーン ・不正発見 ・リスク管理
分類 セグメンテーション	<ul style="list-style-type: none"> ・自己組織化マップ ・クラスター分析 ・ニューラルネット ・主成分分析 ・コレスポンデンス分析 ・決定木 	<ul style="list-style-type: none"> ・優良顧客の属性 ・市場セグメンテーション
関連性の発見 (リンク分析)	<ul style="list-style-type: none"> ・連関規則 ・時系列パターン分析 ・主成分分析 ・コレスポンデンス分析 	<ul style="list-style-type: none"> ・マーケットバスケット分析

27

データマイニング・統計手法・OLAP

	データマイニング	統計解析	OLAP
分析目的	探索的解析 目的志向解析	目的志向解析 (探索的解析)	履歴データ
分析プロセス	パターンの発見と 検証	仮説検証	主観に基づく集計
難易度	簡易に実行可能 だが、各技術の基 本的知識が必要。 (半自動)	各手法の基本的 知識および統計に 関する知識が必要。 (対話的) モデルが扱い易い	容易(対話的)
手法	機械学習(NN) リンク分析 グラフ化	多変量解析 グラフ化	集計 抽出

データマイニング・統計手法・OLAP

データマイニング手法 (機械学習)

遺伝アルゴリズム
最適化

ニューラルネット

クラスター分析

決定木

関連規則

パターン分析

グラフ化

集計

統計解析

回帰分析

時系列分析

対応分析

主成分分析

数量化法

決定木

因子分析

正準相関分析

判別分析

OLAP

抽出

Rでのデータマイニング

Ohmsha
Publisher of Science and Engineering Books

Home

特約書店一覧 | ダウンロード | 購入案内 | カートを見る

Home > 理工学専門書 > 理学 > 数学 > Rによるデータマイニング入門

Rによるデータマイニング入門

著者：山本 義郎 藤野 友和 久保田貴文 共著

定価：3,132 円(本体2,900 円+税)

A5 244頁

ISBN 978-4-274-21817-0

発売日 2015/11

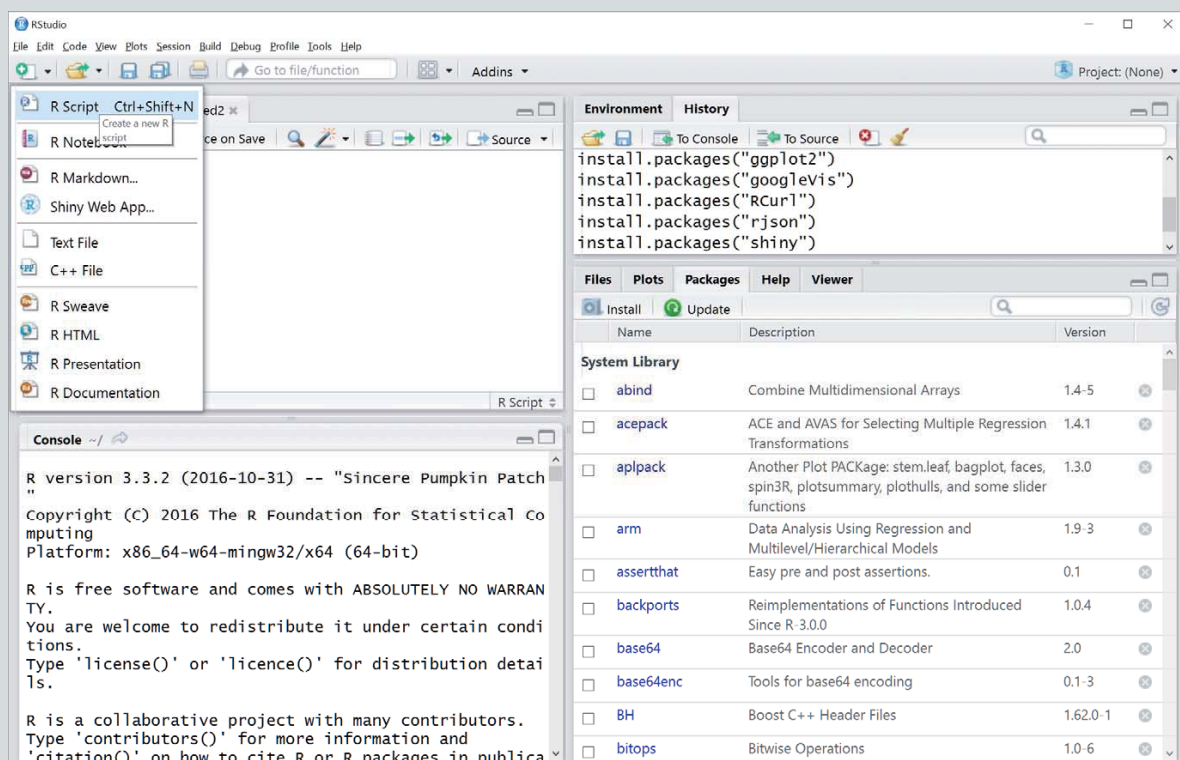
カートに入れる

※本体価格は変更される場合があります。
※通常2〜3日以内に発送いたします。
合計5000円（税別）以上のご注文の場合、配送料は無料となります。
本書で取り上げたデータとRの例題ファイル
◆PDF版はこちら

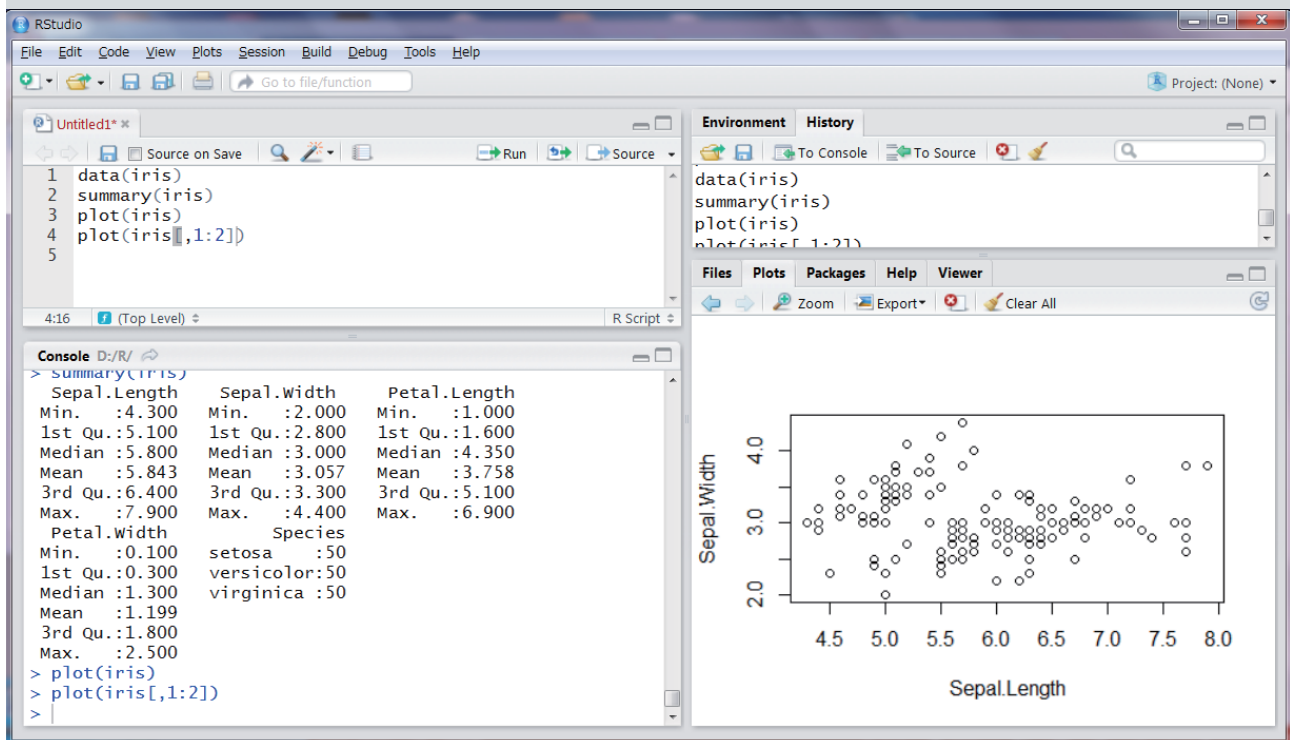
RSTUDIOを使う

- Rstudioの長所
 - 作業の記録(プログラムの編集)
 - グラフ
 - 複数参照できる
 - インタラクティブ
 - ドキュメンテーション機能が充実

RSTUDIOの起動



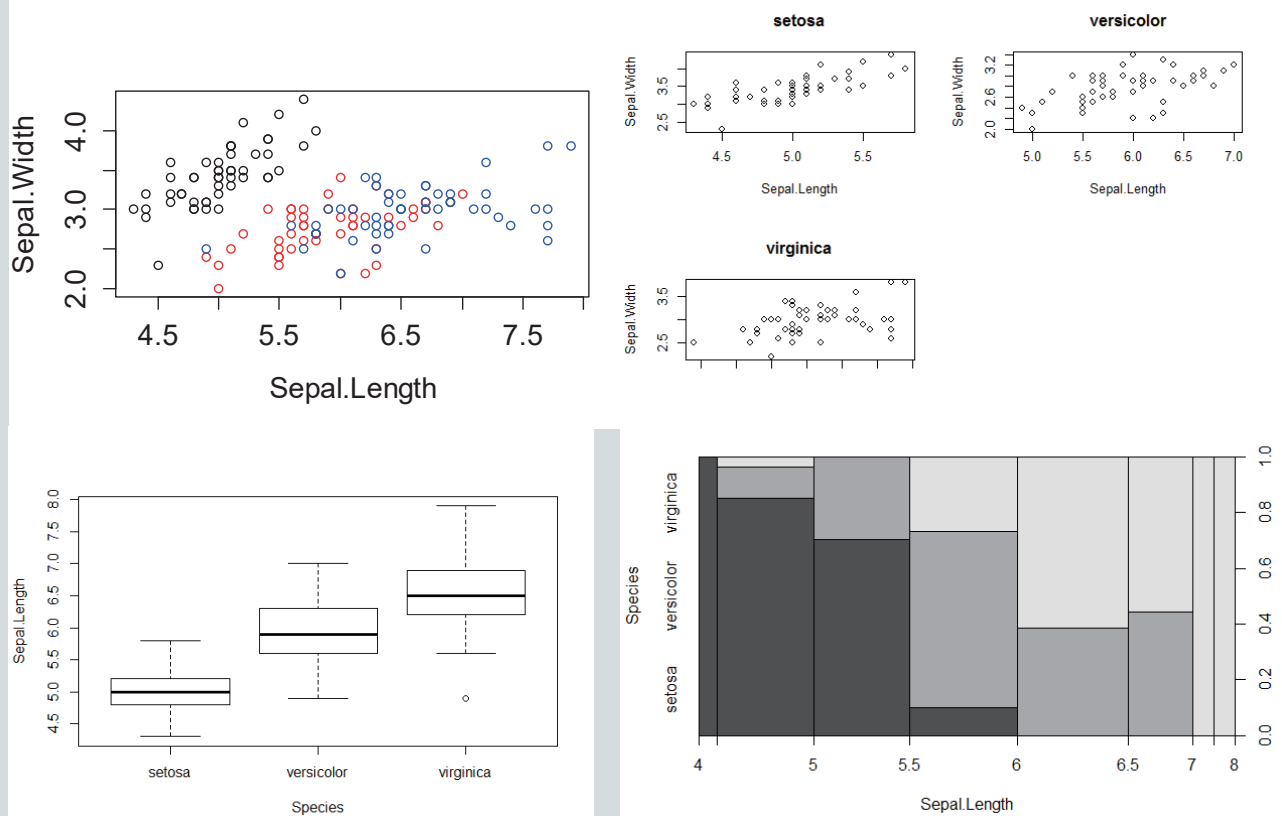
RSTUDIOでの作業



Rでのグラフ描画

```
plot(iris[,1:2],type="n")
points(iris[iris$Species=="setosa",1:2])
points(iris[iris$Species=="versicolor",1:2],col="red")
points(iris[iris$Species=="virginica",1:2],col="blue")
plot(Sepal.Width~Sepal.Length,data=iris)
plot(Sepal.Width~Sepal.Length,data=iris,col=Species)
legend("topright",legend=levels(iris$Species),pch=1,col=1:3)
plot(Sepal.Width~Sepal.Length,data=iris,col=Species)
legend(x=6.8,y=4.5,legend=levels(iris$Species),pch=1,col=1:3)
par(mfrow=c(2,2))
plot(iris[iris$Species=="setosa",1:2],main="setosa")
plot(iris[iris$Species=="versicolor",1:2],main="versicolor")
plot(iris[iris$Species=="virginica",1:2],main="virginica")
par(mfrow=c(1,1))
plot(Sepal.Length~Species,data=iris)
```

Rでのグラフ描画



LATTICEグラフ

```
xyplot(Sepal.Width~Sepal.Length,data=iris)
xyplot(Sepal.Width~Sepal.Length|Species,data=iris)
xyplot(Sepal.Width~Sepal.Length,group=Species,data=iris)
xyplot(Sepal.Width~Sepal.Length,group=Species,data=iris,auto.key=list(title="Species",column=3))

xyp.iris2 <- xyplot(Sepal.Width~Sepal.Length,group=Species,data=iris)
update(xyp.iris2,auto.key=list(title="Species",column=3))

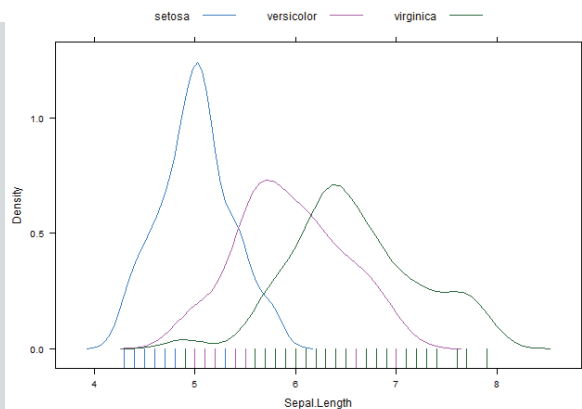
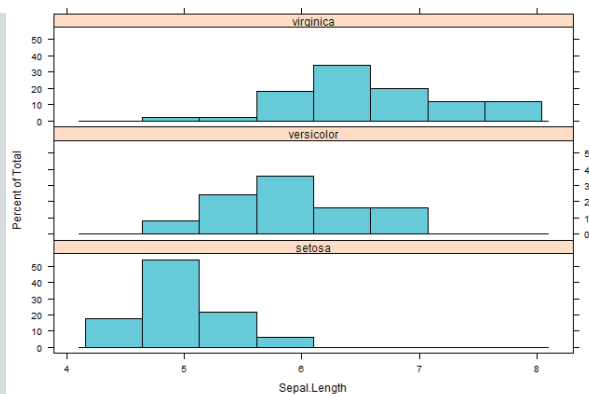
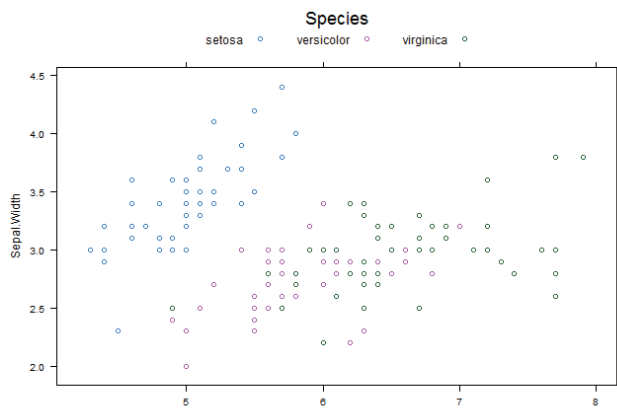
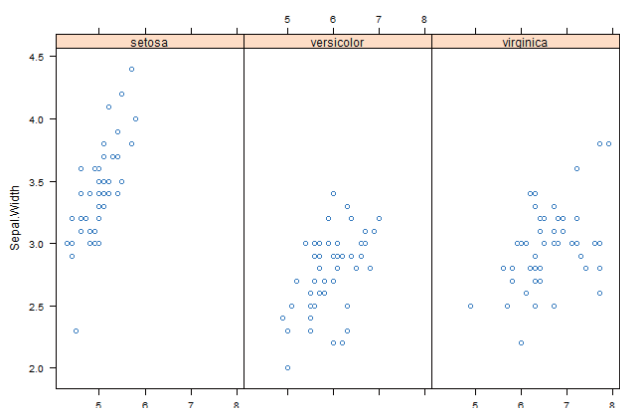
histogram(~Sepal.Length,data=iris)
histogram(~Sepal.Length|Species,data=iris)
histogram(~Sepal.Length|Species,data=iris,layout=c(1,3))

densityplot(~Sepal.Length,data=iris)
densityplot(~Sepal.Length|Species,data=iris,plot.points="rug",layout=c(1,3))
densityplot(~Sepal.Length,group=Species,data=iris,plot.points="rug",auto.key=list(column=3))

bwplot(Species~Sepal.Length,data=iris)
bwplot(Sepal.Length~Species,data=iris)

stripplot(Sepal.Length~Species,data=iris)
stripplot(Sepal.Length~Species,data=iris,jitter.data=TRUE)
```

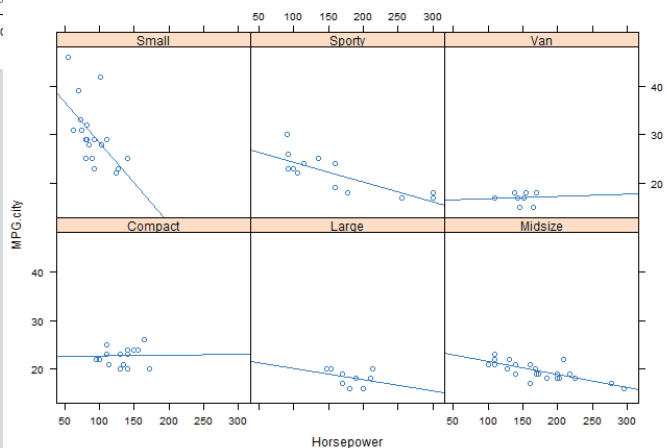
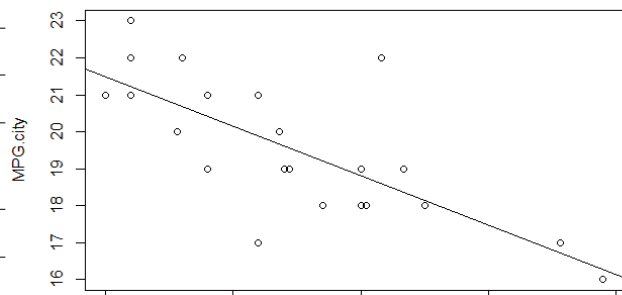
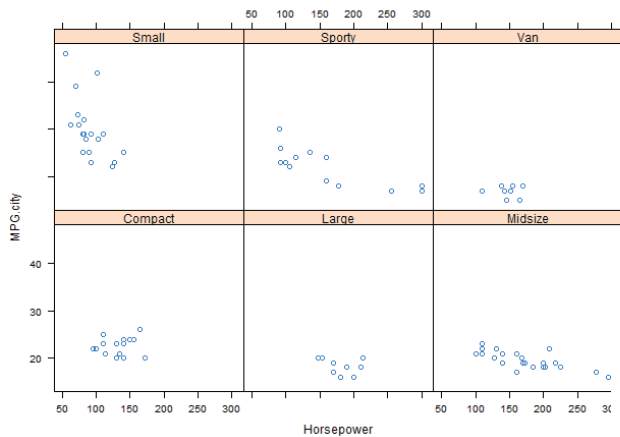
LATTICE グラフ



データの抽出利用

```
data(Cars93, package="MASS")
names(Cars93)
xyplot(MPG.city~Horsepower,data=Cars93)
xyplot(MPG.city~Horsepower|Type,data=Cars93)
Cars93[Cars93$Type=="Midsize",c(7,13)]
plot(Cars93[Cars93$Type=="Midsize",c(7,13)])
subset(Cars93,Type=="Midsize",c(Horsepower,MPG.city))
with(subset(Cars93,Type=="Midsize",c(Horsepower,MPG.city)),
      plot(MPG.city~Horsepower))
with(subset(Cars93,Type=="Midsize",c(Horsepower,MPG.city)),
      {plot(MPG.city~Horsepower)
       lm.mid<-lm(MPG.city~Horsepower)
       summary(lm.mid)
      })
Cars93mid<-subset(Cars93,Type=="Midsize",c(Horsepower,MPG.city))
plot(MPG.city~Horsepower,data=Cars93mid)
lm.mid<-lm(MPG.city~Horsepower,data=Cars93mid)
summary(lm.mid)
```

データの抽出利用



IR実務者向けのR入門

GGPLOT2グラフ

```
library(ggplot2)

ggplot(Cars93, aes(x=Horsepower,y=MPG.city))+geom_point()

gp<-ggplot(Cars93, aes(x=Horsepower,y=MPG.city))

gp+geom_point()

gp+geom_point()+stat_smooth(method=lm)

gp+geom_point()+stat_smooth(method=lm,se=FALSE)

ggplot(Cars93, aes(x=Horsepower,y=MPG.city,colour=Type))+geom_point()

gpType<-ggplot(Cars93, aes(x=Horsepower,y=MPG.city,colour=Type))

gpType+geom_point()+stat_smooth(method=lm,se=FALSE)

gpType2<-ggplot(Cars93, aes(x=Horsepower,y=MPG.city))

gpType2+geom_point()+stat_smooth(method=lm,se=FALSE)+facet_grid(Type~.)

gpType2+geom_point()+stat_smooth(method=lm,se=FALSE)+facet_wrap(~Type)

gpType2+geom_point()+stat_smooth(method=lm,se=FALSE)+facet_wrap(~Type,ncol=2)

ggplot(Cars93,aes(x=Horsepower,colour=Type,fill=Type))+geom_histogram()

ggplot(Cars93,aes(x=Horsepower,colour=Type,fill=Type))+geom_histogram(position="identity",alpha=0.4)

ggplot(Cars93,aes(x=Horsepower,colour=Type,fill=Type))+geom_density(alpha=0.4)

ggplot(Cars93,aes(x=Horsepower))+geom_histogram()+facet_grid(Type~.)

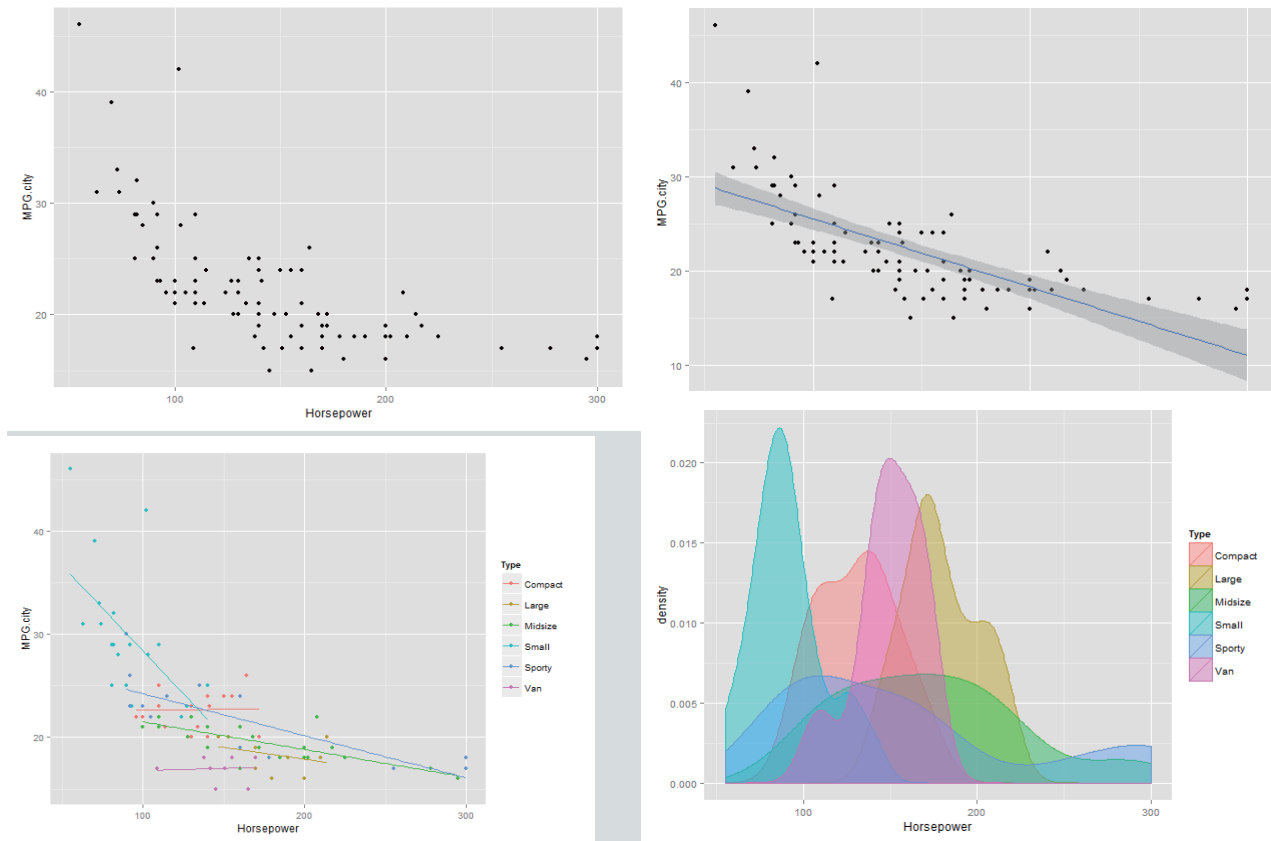
ggplot(Cars93,aes(x=Horsepower))+geom_density()+facet_grid(Type~.)
```

IR実務者向けのR入門

2017/2/10

40

GGPLOT2グラフ



DPLYRによるデータの抽出

```
head(select(Cars93,Horsepower,MPG.city,Type))
head(filter(Cars93,Type=="Midsize"))

Cars93 %>%
  select(Horsepower,MPG.city,Type) %>%
  filter(Type=="Midsize") %>%
  ggplot(aes(x=Horsepower,y=MPG.city))+geom_point()+stat_smooth(method=lm)

Cars93mid <- Cars93 %>%
  select(Horsepower,MPG.city,Type) %>%
  filter(Type=="Midsize")

ggplot(Cars93mid,aes(x=Horsepower,y=MPG.city))+geom_point()+stat_smooth(method=lm)
lm.Mid<-lm(MPG.city~Horsepower,data=Cars93mid)
summary(lm.Mid)

head(mutate(Cars93mid,var1=MPG.city/Horsepower))
head(arrange(Cars93mid,Horsepower)) #desc(Horsepower)

summarize(Cars93mid,n=n(),m.Hp=mean(Horsepower),s.Hp=sd(Horsepower))

grouped.Cars93 <- group_by(Cars93,Type)
summarize(grouped.Cars93,n=n(),m.Hp=mean(Horsepower),s.Hp=sd(Horsepower))
```

実データの解析例

- プロ野球の打者の成績
 - Excelデータ
 - フィルタ、新たな変数の作成
 - Rでのコマンドによる作業
 - 毎年のデータの読み込みと結合
 - Rstudioの機能
 - manipulateパッケージ
 - R ドキュメンテーション
 - Shiny