# 資料科學與機器學習

國立臺灣大學共同教育中心

蔡芸琤

Statistics

Pattern Recognition

Neurocomputing

Machine Learning

AI

Data Mining

Databases

KDD

Knowledge Discovery and Data Mining

# 機器學習是甚麼？

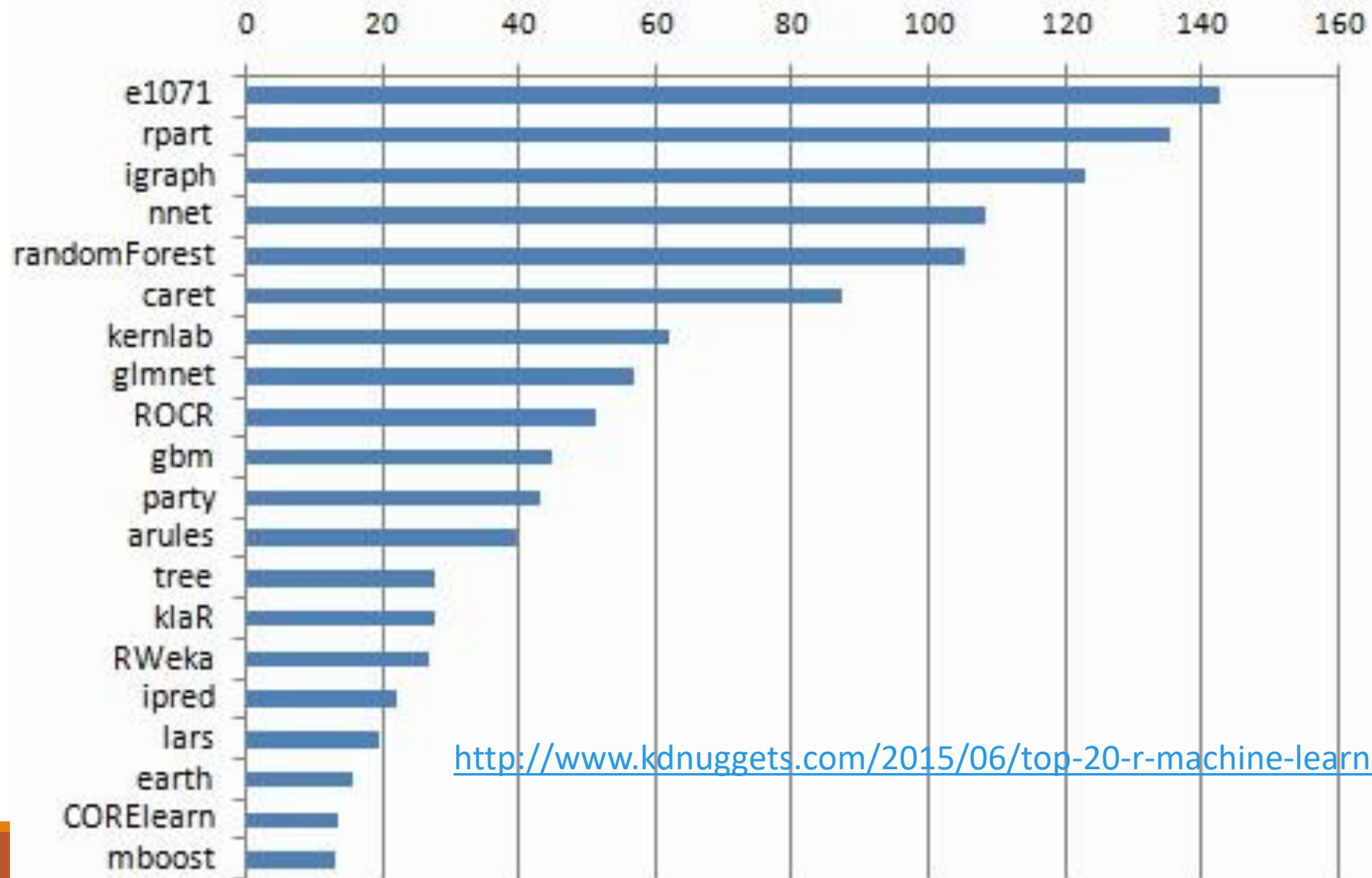# Machine Learning Algorithms *(sample)*

## Unsupervised

**Continuous**

- Clustering & Dimensionality Reduction
  - SVD
  - PCA
  - K-means

**Categorical**

- Association Analysis
  - Apriori
  - FP-Growth
- Hidden Markov Model

## Supervised

- Regression
  - Linear
  - Polynomial
- Decision Trees
- Random Forests

- Classification
  - KNN
  - Trees
  - Logistic Regression
  - Naive-Bayes
  - SVM

http://www.cc.ntu.edu.tw/chinese/epaper/0031/20141220_3105.html

# Top 20 R Machine Learning packages, by Downloads (000) from CRAN

# Principal Components Analysis

# 套件安裝

https://github.com/vqv/ggbiplot

library(devtools)

install_github("ggbiplot", "vqv")

library(scales)

library(grid)

library(ggbiplot)

# Step 2.1

| i | $X_1$ | $X_2$ |
|---|---|---|
| A | 1 | 1 |
| B | 1 | 0 |
| C | 0 | 2 |
| D | 2 | 4 |
| E | 3 | 5 |

$\overline{X}_1^1 = (1, 0.5)$

$\overline{X}_2^1 = (1.7, 3.7)$

| i | 1 | 2 |
|---|---|---|
| A | 0.5 | 2.7 |
| B | 0.5 | 3.7 |
| C | 1.8 | 2.4 |
| D | 3.6 | 0.5 |
| E | 4.9 | 1.9 |

Compute distances between each of the cluster means and all other points.

https://youtu.be/mtkWR8sx0NA

# 以 TF-IDF 進行 K-Means 分群

Go to file/function          Addins ▾                                              Project: (None) ▾

DLTM.R ✕ | MLDM.R ✕ | as.data.frame(kmeansOut$cluster) ✕ | kmeansOut$cluster ✕ | log.ir ✕ | PCA.R ✕ | doc.tfidf ✕ | inspect(tdm[1:9, 1:10]) ✕ | s.tdm ✕ | austen_bigrams ✕

Source on Save                                                                Run    Source ▾

```r
1  source('MLDM.R')
2
3  testTfidf = doc.tfidf
4  tfidf.pca <- prcomp(testTfidf)
5  tfidf.kmeans <- as.factor(kmeansOut$cluster)
6
7  g <- ggbiplot(tfidf.pca, obs.scale = 1, var.scale = 1,
8                groups = tfidf.kmeans, ellipse = TRUE,
9                circle = TRUE, labels = rownames(testTfidf))
10 g <- g + scale_color_discrete(name = '')
11 g <- g + theme(legend.direction = 'horizontal',
12                legend.position = 'top')
13 print(g)
```

17:1    (Top Level)                                                              R Script ≑

Console  C:/Users/pecu6/Desktop/Word2Vec/

```r
> g <- g + scale_color_discrete(name = '')
> g <- g + theme(legend.direction = 'horizontal',
+                legend.position = 'top')
> print(g)
>
```

Environment  History

Import Dataset ▾                                    List ▾

Global Environment ▾

| g | Large gg (… |
| idf | Named num … |
| kme… | List of 9 |
| loc | int [1:581… |
| mix… | Environment |
| N | 10L |
| non… | Named int … |
| ord… | Named num … |
| q | chr [1:581… |
| res… | chr [1:15]… |

Files  Plots  Packages  Help  Viewer

Zoom    Export ▾

DLTM.R ×   MLDM.R ×   as.data.frame(kmeansOut$cluster) ×   kmeansOut$cluster ×   log.ir ×   PCA.R ×   doc.tfidf ×   inspect(tdm[1:9, 1:10]) ×   s.tdm ×   austen_bigrams ×

Source on Save   Run   Source
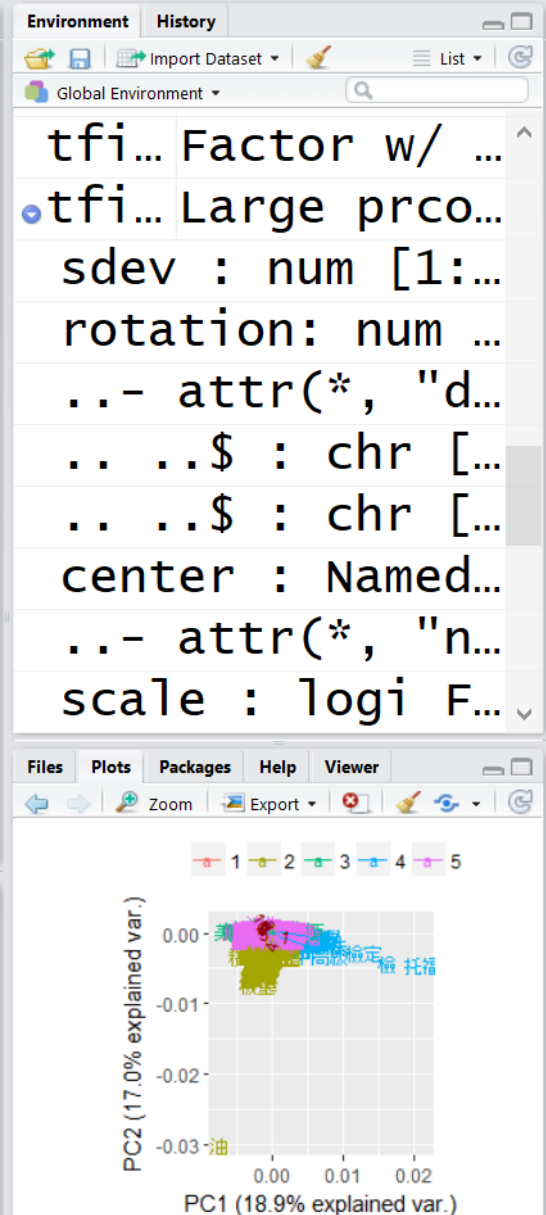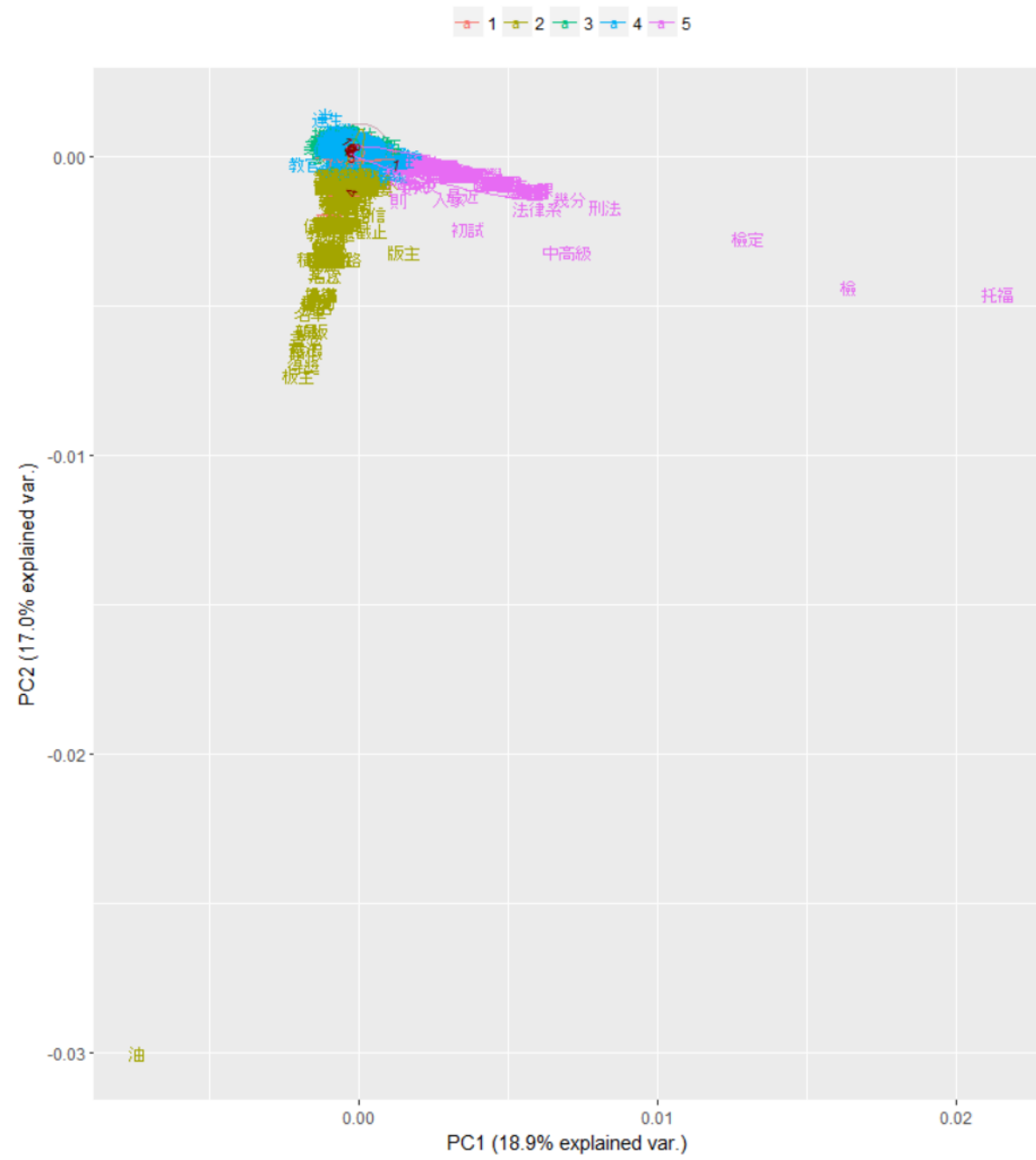
```r
1  source('MLDM.R')
2
3  tfidf.pca <- prcomp(doc.tfidf)
4  tfidf.kmeans <- as.factor(kmeansOut$cluster)
5
6  g <- ggbiplot(tfidf.pca, obs.scale = 1, var.scale = 1,
7                groups = tfidf.kmeans, ellipse = TRUE,
8                circle = TRUE, labels = rownames(doc.tfidf))
9  g <- g + scale_color_discrete(name = '')
10 g <- g + theme(legend.direction = 'horizontal',
11                legend.position = 'top')
12 print(g)
13
```

3:1   (Top Level)   R Script

**Environment**  History

Global Environment

tfi… Factor w/ …
tfi… Large prco…
  sdev : num [1:…
  rotation: num …
  ..- attr(*, "d…
  .. ..$ : chr […
  .. ..$ : chr […
  center : Named…
  ..- attr(*, "n…
  scale : logi F…

**Files**  **Plots**  Packages  Help  Viewer

Zoom   Export

Console  C:/Users/pecu6/Desktop/Word2Vec/

>

Plot Zoom

1  2  3  4  5

PC2 (17.0% explained var.)

PC1 (18.9% explained var.)
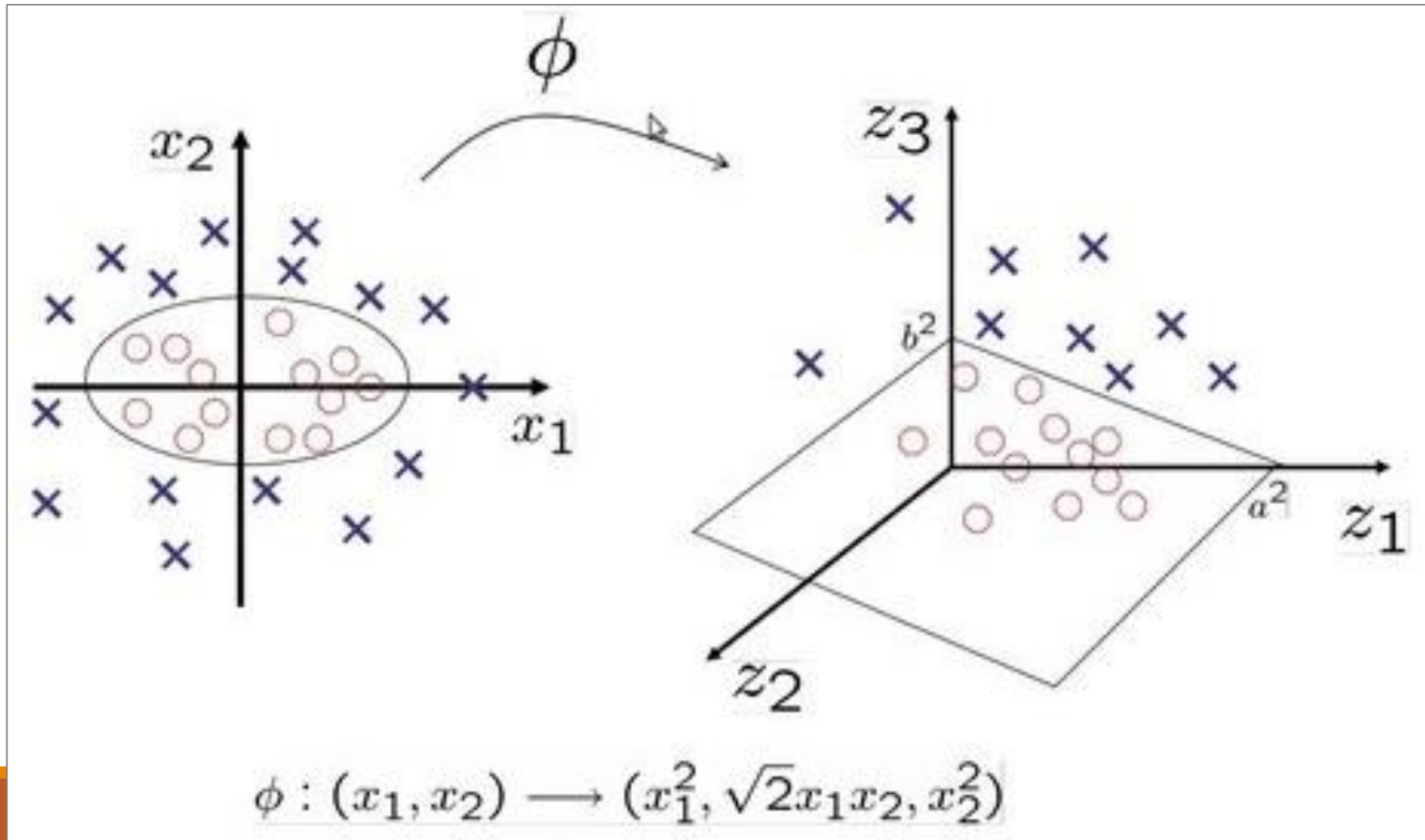
0.00

-0.01

-0.02

-0.03

0.00   0.01   0.02

教育

進步

則 入 幕近 法律系 分 刑法
初試
中高級
檢定
檢
托福
信 止
版主
路
名單
板主

油

# Support Vector Machine



$$\phi : (x_1, x_2) \longrightarrow (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

# Support Vector Regression with R

Linear 資料準備

# 使用套件：library(e1071)

RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Go to file/function          Addins ▾                          Project: (None) ▾

DLTM.R ×   SVM.R* ×   dat ×   as.data.frame(kmeansOut$cluster) ×   kmeansOut$cluster ×   log.ir ×   PCA.R ×   doc.tfidf ×   inspect(tdm[1:9, 1:10]) ×   s.tdm ×   austen_bigrams ×   longley ×
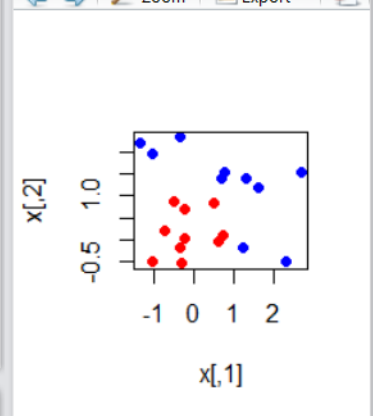
Environment   History

```r
1  set.seed(10111)
2  x = matrix(rnorm(40), 20, 2)
3  y = rep(c(-1, 1), c(10, 10))
4  x[y == 1, ] = x[y == 1, ] + 1
5  plot(x, col = y + 3, pch = 19)
6
7  library(e1071)
8  dat = data.frame(x, y = as.factor(y))
9  svmfit = svm(y ~ ., data = dat,
10             kernel = "linear",
11             cost = 10, scale = FALSE)
12 print(svmfit)
13 plot(svmfit, dat)
14
15 make.grid = function(x, n = 75) {
16   grange = apply(x, 2, range)
17   x1 = seq(from = grange[1, 1], to = grange[2, 1], length = n)
18   x2 = seq(from = grange[1, 2], to = grange[2, 2], length = n)
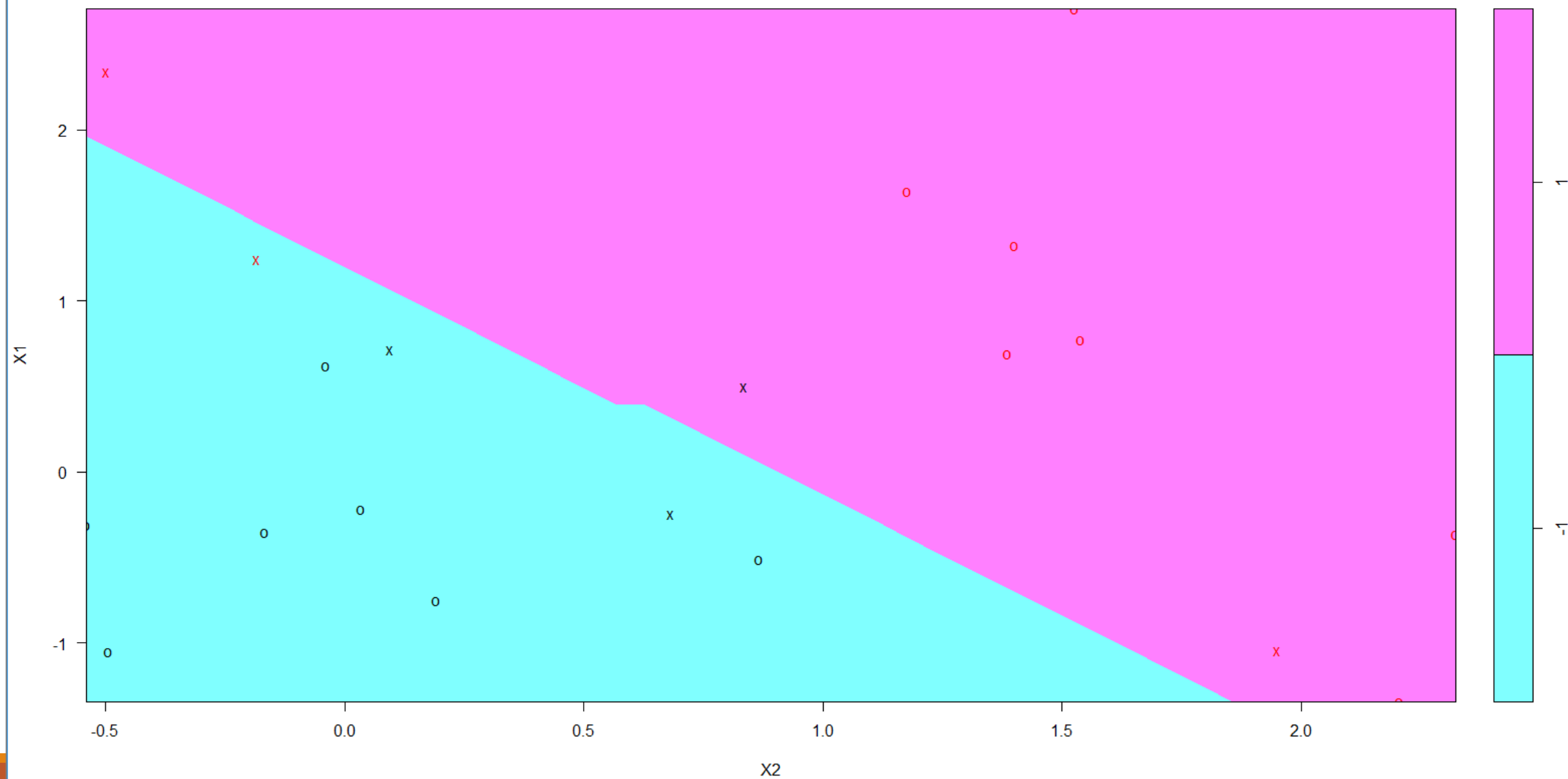```

Global Environment ▾

g...  num  [...
i...  150 o...
n...  Large...
s...  Large...
t...  Large...
x    num  [...
x... 5625 ...

Values
a...  chr  [...
a...  List  ...

Files  Plots  Packages  Help  Viewe
Zoom  Export  

14:1   (Top Level) ÷                                    R Script ÷

Console

SVM classification plot

Non-linear Data

https://gate.ac.uk/sale/tao/