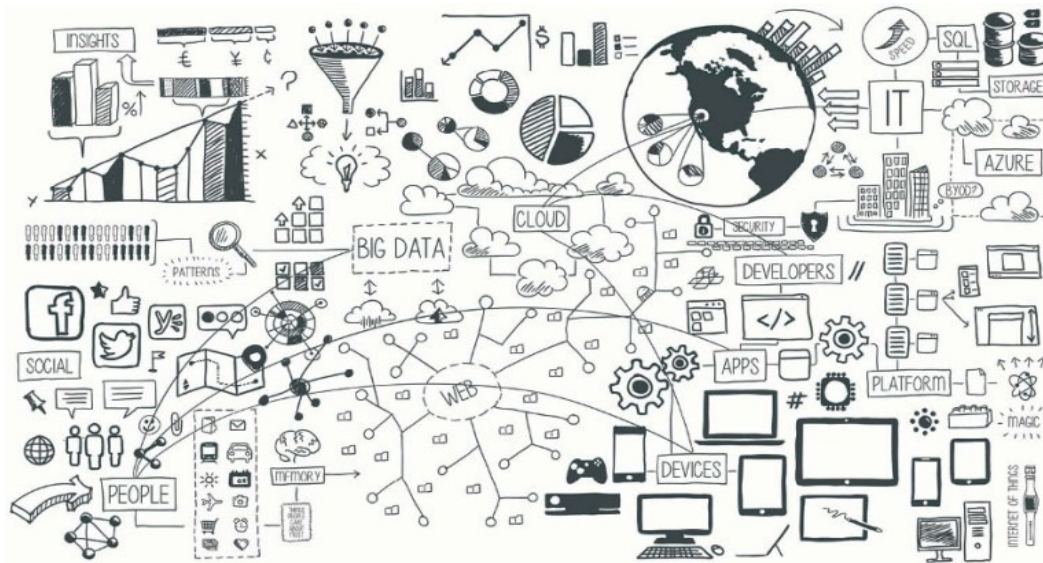


DATASETS & DATA CHALLENGES



Maialen Iturbide

José Manuel Gutiérrez

Grupo de Meteorología

Univ. de Cantabria – CSIC
MACC / IFCA

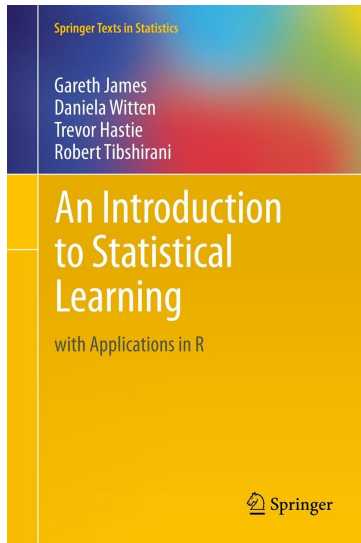


DATASETS & DATA CHALLENGES

NOTA: Las líneas de código de R en esta presentación se muestran sobre un fondo gris.

Oct	30	Aplazada (sesión de refuerzo)
Nov	6	Presentación, introducción y perspectiva histórica
	8	Paradigmas, problemas canónicos y data challenges
	13	Reglas de asociación
	15	Practica: Reglas de asociación
	20	Evaluación, sobreajuste y crossvalidacion
	22	Practica: Crossvalidacion
	27	Arboles de clasificacion y decision
	29	Practica: Arboles de clasificación
		T01. Datos discretos
Dic	4	Técnicas de vecinos cercano (k-NN)
	11	Práctica: Vecinos cercanos
	13	Reducción de dimensión lineal
	18	Practica: LDA y PCA
	20	Reducción no lineal
		T02. Clasificación
Ene	8	Arboles de clasificación y regresion (CART)
	10	Practica: CART
	15	Ensembles: Bagging and Boosting
	17	Practica Random Forests
		T03. Prediccion
	22	Practica Gradient boosting
	24a	Técnicas de agrupamiento
	24b	Practica: Técnicas de agrupamiento
	29a	Practica: El paquete CARET
	29b	Examen

1



An Introduction to Statistical Learning: With Applications in R

James, G., Witten, D., Hastie, T., Tibshirani, R.

Springer (2013)

<http://www-bcf.usc.edu/~gareth/ISL>

```
install.packages("ISLR")  
library("ISLR")  
library(help = "ISLR")
```

2

```
library(help = "datasets")
```

3

kaggle

<https://www.kaggle.com/datasets>

4



<https://archive.ics.uci.edu/ml/datasets.html>

New York City Taxi Trip Duration

Share code and data to improve ride time predictions

\$30,000

Prize Money



Kaggle · 1,257 teams · 3 months ago

[Overview](#)

[Data](#)

[Kernels](#)

[Discussion](#)

[Leaderboard](#)

[Rules](#)

[Late Submission](#)

Overview

Description

Evaluation

Prizes

Timeline

In this competition, Kaggle is challenging you to build a model that predicts the total ride duration of taxi trips in New York City. Your primary dataset is one released by the NYC Taxi and Limousine Commission, which includes pickup time, geo-coordinates, number of passengers, and several other variables.

Longtime Kagglers will recognize that this competition objective is similar to the [ECML/PKDD trip time challenge](#) we hosted in 2015. But, this challenge comes with a twist. Instead of awarding prizes to the top finishers on the leaderboard, this playground competition was created to reward collaboration and collective learning.



<https://www.kaggle.com/headsortails/nyc-taxi-eda-update-the-fast-the-curious/notebook>

Listado de datasets utilizados en el curso

EN FUNCIÓN DE LA NATURALEZA DE LOS DATOS PODEMOS CLASIFICARLAS COMO

SÓLO CATEGÓRICAS (FACTORES)

- **Groceries.** Disponible en kaggle y en el paquete {arulesViz} de R.
- **Mushroom.** Disponible en kaggle y UCI.

MIXTOS (CONTINUOS Y FACTORES)

- **Iris.** Disponible en kaggle, UCI y el paquete {datasets} de R.
- **MNIST.** Disponible en <https://pjreddie.com/projects/mnist-in-csv/>
- **Gene expression (Golub et al).** Disponible en kaggle.
- **Meteo** (Santander Meteorology Group).
- **The fruits dataset by Dr. Iain Murray.** Disponible en <https://towardsdatascience.com/solving-a-simple-classification-problem-with-python-fruits-lovers-edition-d20ab6b071d2>



Todos estos datasets están disponibles en **gitHub**:

<https://github.com/SantanderMetGroup/Master-Data-Science>

Instacart Market Basket Analysis

Which products will an Instacart consumer purchase again?

\$25,000

Prize Money



Instacart · 2,623 teams · 4 months ago

Overview Data **Kernels** Discussion Leaderboard Rules

New Kernel

Public

Your Work

Favorites

Sort by Hotness

<https://www.kaggle.com/philippsp/exploratory-analysis-instacart>

Search kernels



588



Exploratory Analysis - Instacart

5mo ago

intermediate, eda, data visualization



Rmd

119



Instacart XGBoost Starter - LB 0.3791

En el curso utilizaremos un dataset más pequeño, “Groceries”, disponible en el paquete de R **arulesViz**.

Attribute characteristics	Categorical
Number of instances	9835
Number of attributes	169

```
install.packages("arulesViz")
data("Groceries")
```

Mushroom Classification

Safe to eat or deadly poison?

<https://www.kaggle.com/uciml/mushroom-classification/data>



UCI Machine Learning • last updated a year ago

Overview

Data

Kernels

Discussion

Activity

Download (30 KB)

New Kernel

<http://archive.ics.uci.edu/ml/datasets/Mushroom>

Data Set Characteristics:	Multivariate	Number of Instances:	8124	Area:	Life
Attribute Characteristics:	Categorical	Number of Attributes:	22	Date Donated	1987-04-27
Associated Tasks:	Classification	Missing Values?	Yes	Number of Web Hits:	298439

Attribute Information: (classes: edible=e, poisonous=p)

cap-shape: bell=b,conical=c,convex=x,flat=f, knobbed=k,sunken=s

cap-surface: fibrous=f,grooves=g,scaly=y,smooth=s

cap-color: brown=n,buff=b,cinnamon=c,gray=g,green=r,pink=p,purple=u,...

bruises: bruises=t,no=f

odor: almond=a,anise=l,creosote=c,fishy=y,foul=f,musty=m,none=n,...

...

```
mush <- read.csv("Data_mining/datasets/mushrooms.csv")
str(mush)
```

```
'data.frame': 8124 obs. of 23 variables:
 $ class          : Factor w/ 2 levels "e","p": 2 1 1 2 1 1 1 1 2 1 ...
 $ cap.shape      : Factor w/ 6 levels "b","c","f","k",...: 6 6 1 6 6 6 1 1 6 1 ...
 $ cap.surface    : Factor w/ 4 levels "f","g","s","y": 3 3 3 4 3 4 3 4 4 3 ...
 $ cap.color      : Factor w/ 10 levels "b","c","e","g",...: 5 10 9 9 4 10 9 9 9 10 ...
 $ bruises        : Factor w/ 2 levels "f","t": 2 2 2 2 1 2 2 2 2 2 ...
 $ odor           : Factor w/ 9 levels "a","c","f","l",...: 7 1 4 7 6 1 1 4 7 1 ...
 ...
```

Featured Dataset

428

Iris Species

Classify iris plants into three species in this classic dataset

UCI Machine Learning • last updated a year ago

[Overview](#)
[Data](#)
[Kernels](#)
[Discussion](#)
[Activity](#)
[Download \(4 KB\)](#)
[New Kernel](#)

Sort by Hotness

<http://archive.ics.uci.edu/ml/datasets/Iris>

Data Set Characteristics:	Multivariate	Number of Instances:	150	Area:	Life
Attribute Characteristics:	Real	Number of Attributes:	4	Date Donated	1988-07-01
Associated Tasks:	Classification	Missing Values?	No	Number of Web Hits:	1549312

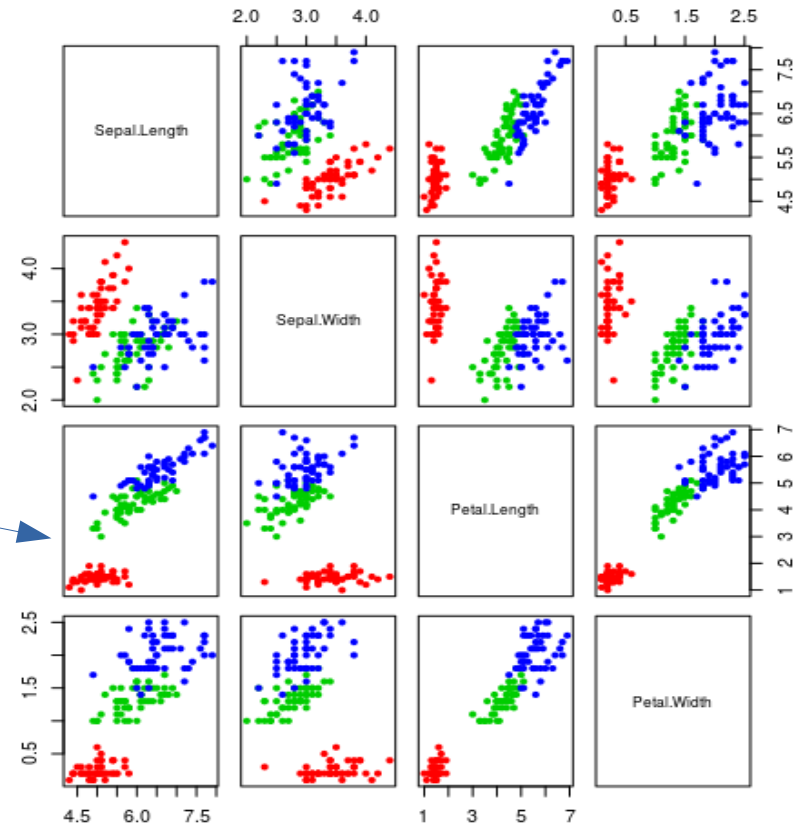


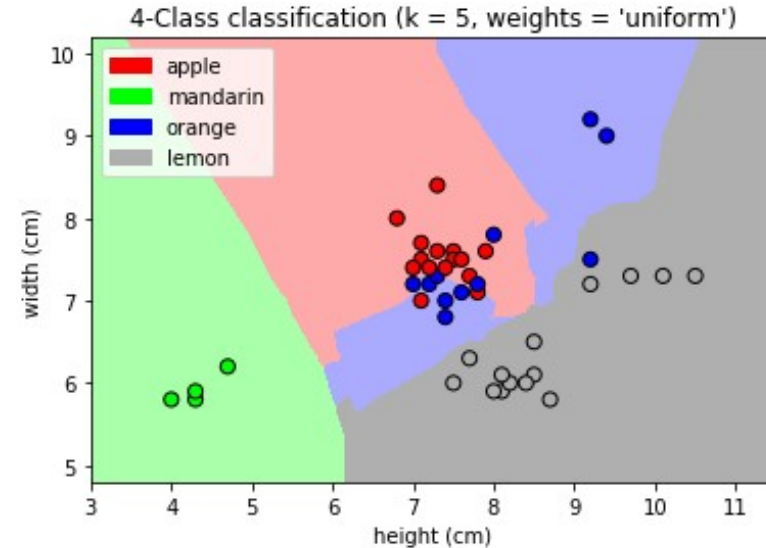
```
data("iris")
str(iris)

'data.frame':  150 obs. of  5 variables:
 $ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
 $ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
 $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
 $ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
 $ Species      : Factor w/ 3 levels "setosa","versicolor",...: 1 1 ...
```

```
library(graphics)
pairs(iris[1:4],
      main = "Anderson's Iris species",
      pch = 20,
      col = c("red", "green3", "blue")[unclass(iris$Species)])
```

Anderson's Iris Data -- 3 species





<https://towardsdatascience.com/solving-a-simple-classification-problem-with-python-fruits-lovers-edition-d20ab6b071d2>

```
fruits <- read.table("Data_mining/datasets/fruits.txt", header = TRUE)
str(fruits)
```

```
'data.frame':  59 obs. of  7 variables:
 $ fruit_label  : int  1 1 1 2 2 2 2 2 1 1 ...
 $ fruit_name   : Factor w/ 4 levels "apple","lemon",...: 1 1 1 3 3 3 3 3 1 1 ...
 $ fruit_subtype: Factor w/ 10 levels "braeburn","cripps_pink",...: 4 4 4 5 5 5 5 5 1 1 ...
 $ mass        : int  192 180 176 86 84 80 80 76 178 172 ...
 $ width       : num  8.4 8 7.4 6.2 6 5.8 5.9 5.8 7.1 7.4 ...
 $ height      : num  7.3 6.8 7.2 4.7 4.6 4.3 4.3 4 7.8 7 ...
 $ color_score  : num  0.55 0.59 0.6 0.8 0.79 0.77 0.81 0.81 0.92 0.89 ...
```

Gene expression dataset (Golub et al.)

Molecular Classification of Cancer by Gene Expression Monitoring



Chris Crawford • last updated 4 months ago

Overview **Data** Kernels Discussion Activity

Download (1 MB)

New Kernel

Optimization Based Tumor Classification from Microarray Gene Expression Data

Onur Dagliyan¹, Fadime Uney-Yuksektepe², I. Halil Kavakli¹, Metin Turkey^{3*}

An important use of data obtained from microarray measurements is the classification of tumor types with respect to genes that are either up or down regulated in specific cancer types.

Table 1. Cancer data sets used in this study.

Data set	Samples	Genes	Classes	Reference
Leukemia	72	7129	2	Golub et al. (1999)
Prostate cancer	102	12600	2	Singh et al. (2002)
Prostate outcome	21	12600	2	Singh et al. (2002)
DLBCL	77	7129	2	Shipp et al. (2002)

Gene expression dataset (Golub et al.)

Molecular Classification of Cancer by Gene Expression Monitoring



Chris Crawford • last updated 4 months ago

Overview Data Kernels Discussion Activity

Download (1 MB)

New Kernel

Optimization Based Tumor Classification from Microarray Gene Expression Data

Onur Dagliyan¹, Fadime Uney-Yuksektepe², I. Halil Kavakli¹, Metin Turkey^{3*}

```
gene <- read.csv("Data_mining/datasets/gene_trainDF.csv")
str(gene)
```

```
'data.frame': 38 obs. of 7130 variables:
 $ X1 : num 1.1314 1.3258 -2.0812 0.8449 -0.0963 ...
 $ X2 : num 0.459 0.48 -0.332 1.156 0.844 ...
 ...
 $ X7129: num -0.16 0.412 -0.26 -1.504 0.139 ...
 $ label: Factor w/ 2 levels "ALL","AML": 1 1 1 1 1 1 1 1 1 1 ...
```


The highest accuracy is obtained with the optimal gene set consisting of 4 genes:

- Myeloperoxidase (M19507-at),
- adipsin (M84526-at),
- CD33 antigen and
- TCF3 transcription factor 3.

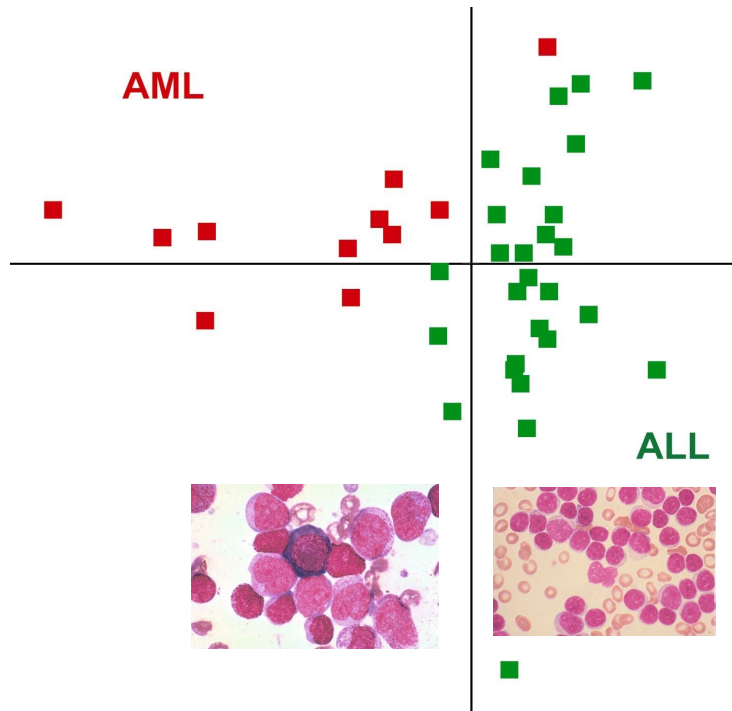
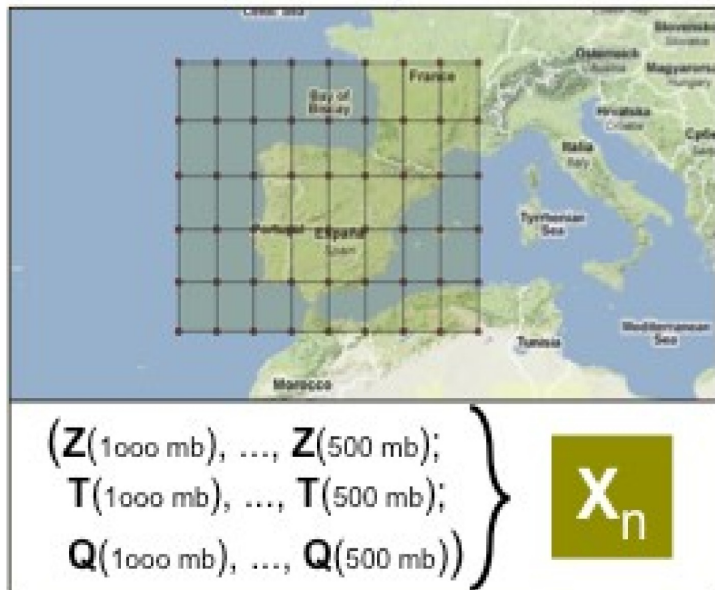


Table 2. Classification results of leukemia (AML-ALL) data set

Classifier	Test Set	10-CV	LOOCV
HBE	100	97.146 0.903	98.61
BayesNet	94.12	95.71	95.83
LibSVM	58.82	86.576 10.44	91.67
SMD	97.06	93.146 0.571	94.44
Logistic Regression	91.18	96.866 1.67	98.61
FBF Network	97.06	97.43 ± 1.07	97.22
IEk	97.06	96.006 1.40	95.83
J48	94.12	89.146 1.94	90.28
Random Forest	94.12	93.146 1.07	90.2

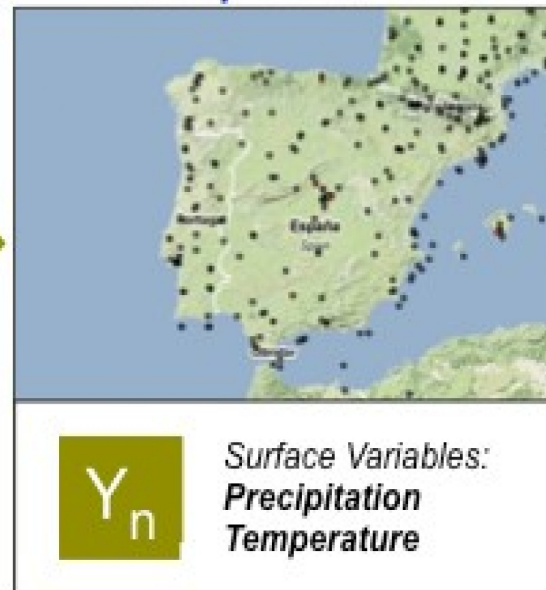
Large scale predictors

Downscaling
Model

$$Y_n = f(X_n)$$

Statistical methods
based on historical
data to link large
scale circulation to
local climates.

Local predictands



Predictors: Z500,T850,T700,T500,2T,Q850,Q500,SLP
lonLim = (-10,4),
latLim = (36,44),
years = 1979:2008

Predictand: precipitation in Lisboa.
LonLim = -9.15
LatLim = 38.7
years = 1979:2008

```
meteo <- read.csv("Data_mining/datasets/meteo.csv")
str(meteo)
```

```
'data.frame': 10958 obs. of 321 variables:
 $ y : num 10.9 0.6 13 0 0 1.2 1.1 0 0 0.7 ...
 $ X1 : num 57043 56963 56523 54628 53584 ...
 $ X2 : num 56535 56493 55971 53980 53391 ...
 $ X3 : num 55884 55931 55304 53494 53310 ...
 $ X4 : num 55176 55340 54498 53073 53293 ...
```

9 6 6 5 4 0 7 4 0 1
3 1 3 4 7 2 7 1 2 1
1 7 4 2 3 5 1 2 4 4

Digit Recognizer

Learn computer vision fundamentals with the famous MNIST data

1,996 teams · 2 years to go

MIXTO

[Overview](#) [Data](#) [Kernels](#) [Discussion](#) [Leaderboard](#) [Rules](#)

<https://www.kaggle.com/c/digit-recognizer#tutorial>

```
mydatadir <- paste0(getwd(), "/MNIST_train.csv")  
train <- read.csv(mydatadir)  
str(train)
```

```
'data.frame':  42000 obs. of  785 variables:  
 $label :int 1 0 1 4 0 0 7 3 5 3 ...  
 $pixel0 :int 0 0 0 0 0 0 0 0 0 0 ...  
 $pixel1 :int 0 0 0 0 0 0 0 0 0 0 ...  
 $pixel2 :int 0 0 0 0 0 0 0 0 0 0 ...  
 ...
```

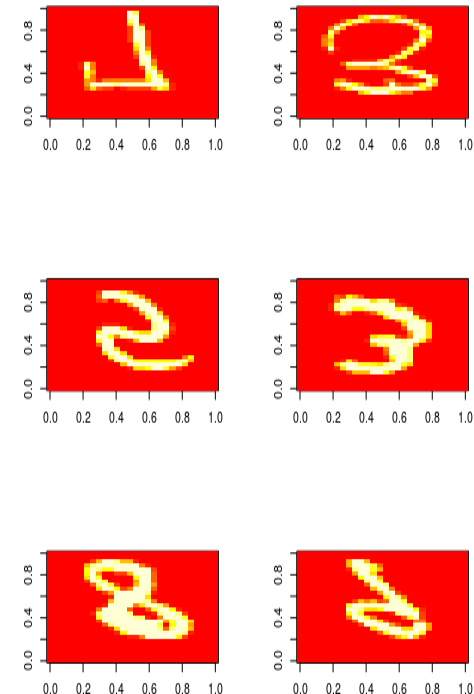
```
# split data into response variable (y) and predictors (x)  
y <- train[,1]; x <- train[,-1]  
dim(x)
```

```
[1] 42000  784
```

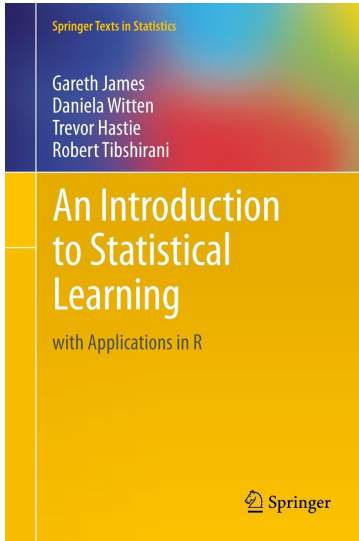
```
par(mfrow = c(3,2))  
image(matrix(as.matrix(x[7,]), nrow = sqrt(784), ncol = sqrt(784)))  
for (i in 8:12) {  
  image(matrix(as.matrix(x[i,]), nrow = sqrt(784), ncol = sqrt(784)))  
}
```

```
Y[7:12]
```

```
[1] 7 3 5 3 8 9
```



1 30-60mins



Echa un vistazo a los datasets que hay en el paquete ISLR.

```
install.packages("ISLR")  
library("ISLR")  
library(help = "ISLR")
```

Analiza la estructura de los datasets: ¿de qué tipo son? ¿para qué tipo de problemas serían adecuados? e.g.

```
data("Hitters")  
str(Hitters)
```

2 60-90mins

Lee con calma el siguiente notebook de kaggle sobre las duraciones de los trayectos de taxi en Nueva York:

<https://www.kaggle.com/headsortails/nyc-taxi-eda-update-the-fast-the-curious/notebook>