

Predictive Analytics

L. Torgo

ltorgo@dal.ca

Faculty of Computer Science / Institute for Big Data Analytics
Dalhousie University

Mar, 2019



Introduction

What is Prediction?

Definition

- Prediction (forecasting) is the ability to anticipate the future.
- Prediction is possible if we assume that there is some regularity in what we observe, i.e. if the observed events are not random.

Example

Medical Diagnosis: given an historical record containing the symptoms observed in several patients and the respective diagnosis, try to forecast the correct diagnosis for a new patient for which we know the symptoms.

Prediction Models

- Are obtained on the basis of the assumption that there is an unknown mechanism that maps the characteristics of the observations into conclusions/diagnoses. The goal of prediction models is to discover this mechanism.
 - Going back to the medical diagnosis what we want is to know how symptoms influence the diagnosis.
- Have access to a data set with “examples” of this mapping, e.g. this patient had symptoms x, y, z and the conclusion was that he had disease p
- Try to obtain, using the available data, an approximation of the unknown function that maps the observation descriptors into the conclusions, i.e. $Prediction = f(Descriptors)$



“Entities” involved in Predictive Modelling

- **Descriptors** of the observation:
 - set of variables that describe the properties (features, attributes) of the cases in the data set
- **Target variable**:
 - what we want to predict/conclude regards the observations
- The goal is to obtain an approximation of the function $Y = f(X_1, X_2, \dots, X_p)$, where Y is the target variable and X_1, X_2, \dots, X_p the variables describing the characteristics of the cases.
- It is assumed that Y is a variable whose values depend on the values of the variables which describe the cases. We just do not know how!
- The goal of the modelling techniques is thus to obtain a good approximation of the unknown function $f()$



How are the Models Used?

Predictive models have two main uses:

1 Prediction

use the obtained models to make predictions regards the target variable of new cases given their descriptors.

2 Comprehensibility

use the models to better understand which are the factors that influence the conclusions.

Types of Prediction Problems

- Depending on the type of the target variable (Y) we may be facing two different types of prediction models:
 - 1 **Classification Problems** - the target variable Y is nominal
e.g. medical diagnosis - given the symptoms of a patient try to predict the diagnosis
 - 2 **Regression Problems** - the target variable Y is numeric
e.g. forecast the market value of a certain asset given its characteristics

Examples of Prediction Problems

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.8	3.0	5.5	2.1	virginica
5.7	2.5	5.0	2.0	virginica
...

Species = $f(\text{Sepal.Length}, \dots)$

Classification Task

Regression Task

season	size	speed	mxPH	mnO2	Cl	Ch1a	...	a1
winter	small	medium	8.00	9.8	60.800	50.0	...	0.0
spring	small	medium	8.35	8.0	57.750	1.3	...	1.4
autumn	small	medium	8.10	11.4	40.020	15.6	...	3.3
spring	small	medium	8.07	4.8	77.364	1.4	...	3.1
autumn	small	medium	8.06	9.0	55.350	10.5	...	9.2
winter	small	high	8.25	13.1	65.750	28.4	...	15.1
...

a1 = $f(\text{season}, \text{size}, \dots)$

TAT
ICA
ICIA

Types of Prediction Models

- There are many techniques that can be used to obtain prediction models based on a data set.
- Independently of the pros and cons of each alternative, all have some key characteristics:
 - 1 They assume a certain **functional form** for the unknown function $f()$
 - 2 Given this assumed form the methods try to obtain the best possible model based on:
 - 1 the given **data set**
 - 2 a certain **preference criterion** that allows comparing the different alternative model variants

Functional Forms of the Models

- There are many variants. Examples include:
 - Mathematical formulae - e.g. linear discriminants
 - Logical approaches - e.g. classification or regression trees, rules
 - Probabilistic approaches - e.g. naive Bayes
 - Other approaches - e.g. neural networks, SVMs, etc.
 - Sets of models (ensembles) - e.g. random forests, adaBoost
- These different approaches entail different compromises in terms of:
 - Assumptions on the unknown form of dependency between the target and the predictors
 - Computational complexity of the search problem
 - Flexibility in terms of being able to approximate different types of functions
 - Interpretability of the resulting model
 - etc.



Which Models or Model Variants to Use?

- This question is often known as the **Model Selection** problem
- The answer depends on the goals of the final user - i.e. the **Preference Biases** of the user
- Establishing which are the preference criteria for a given prediction problem allows to compare and select different models or variants of the same model



Regression Evaluation Metrics

Regression Problems

The setting

- Given data set $\{ \langle \mathbf{x}_i, y_i \rangle \}_{i=1}^N$, where \mathbf{x}_i is a feature vector $\langle x_1, x_2, \dots, x_p \rangle$ and $y_i \in \mathcal{R}$ is the value of the numeric variable Y
- There is an unknown function $Y = f(\mathbf{x})$

The approach

- Assume a functional form $h_\theta(\mathbf{x})$ for the unknown function $f()$, where θ are a set of parameters
- Assume a preference criterion over the space of possible parameterizations of $h()$
- Search for the “optimal” $h()$ according to the criterion and the data set

Preference Criteria

How to evaluate the performance of a regression model?

The Main Criteria

- Prediction error
- Computational complexity
- Model interpretability

Measuring Regression Error

Mean Squared Error

- Given a set of test cases N_{test} we can obtain the predictions for these cases using some regression model.
- The *Mean Squared Error (MSE)* measures the average squared deviation between the predictions and the true values.
- In order to calculate the value of *MSE* we need to have both the predictions and the true values of the N_{test} cases.

Measuring Regression Error

Mean Squared Error (cont.)

- If we have such information the *MSE* can be calculated as follows,

$$MSE = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} (\hat{y}_i - y_i)^2$$

where \hat{y}_i is the prediction of the model under evaluation for the case i and y_i the respective true target variable value.

- Note that the *MSE* is measured in a unit that is squared of the original variable scale. Because of the this is sometimes common to use the *Root Mean Squared Error (RMSE)*, defined as

$$RMSE = \sqrt{MSE}$$



Measuring Regression Error

Mean Absolute Error

- The *Mean Absolute Error (MAE)* measures the average absolute deviation between the predictions and the true values.
- The value of the *MAE* can be calculated as follows,

$$MAE = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} |\hat{y}_i - y_i|$$

where \hat{y}_i is the prediction of the model under evaluation for the case i and y_i the respective true target variable value.

- Note that the *MAE* is measured in the same unit as the original variable scale.



Relative Error Metrics

- Relative error metrics are unit less which means that their scores can be compared across different domains.
- They are calculated by comparing the scores of the model under evaluation against the scores of some baseline model.
- The relative score is expected to be a value between 0 and 1, with values nearer (or even above) 1 representing performances as bad as the baseline model, which is usually chosen as something too naive.



Relative Error Metrics (cont.)

- The most common baseline model is the constant model consisting of predicting for all test cases the average target variable value calculated in the training data.
- The *Normalized Mean Squared Error (NMSE)* is given by,

$$NMSE = \frac{\sum_{i=1}^{N_{test}} (\hat{y}_i - y_i)^2}{\sum_{i=1}^{N_{test}} (\bar{y} - y_i)^2}$$

- The *Normalized Mean Absolute Error (NMAE)* is given by,

$$NMAE = \frac{\sum_{i=1}^{N_{test}} |\hat{y}_i - y_i|}{\sum_{i=1}^{N_{test}} |\bar{y} - y_i|}$$



Relative Error Metrics (cont.)

- The *Mean Average Percentage Error (MAPE)* is given by,

$$MAPE = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} \frac{|\hat{y}_i - y_i|}{y_i}$$

- The *Symmetric Mean Absolute Percentage Error (sMAPE)* is given by,

$$sMAPE = \frac{1}{n} \sum_{i=1}^{N_{test}} \frac{|\hat{y}_i - y_i|}{|\hat{y}_i| + |y_i|}$$



Relative Error Metrics (cont.)

- The *Correlation* between the predictions and the true values ($\rho_{\hat{y},y}$) is given by,

$$\rho_{\hat{y},y} = \frac{\sum_{i=1}^{N_{test}} (\hat{y}_i - \bar{\hat{y}})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{N_{test}} (\hat{y}_i - \bar{\hat{y}})^2 \sum_{i=1}^{N_{test}} (y_i - \bar{y})^2}}$$



An Example in R

```
trueVals <- c(10.2,-3,5.4,3,-43,21,
              32.4,10.4,-65,23)
preds <- c(13.1,-6,0.4,-1.3,-30,1.6,
           3.9,16.2,-6,20.4)
mse <- mean((trueVals-preds)^2)
mse

## [1] 493.991

rmse <- sqrt(mse)
rmse

## [1] 22.22591

mae <- mean(abs(trueVals-preds))
mae

## [1] 14.35

nmse <- sum((trueVals-preds)^2) /
        sum((trueVals-mean(trueVals))^2)
nmse
```

```
nmae <- sum(abs(trueVals-preds)) /
        sum(abs(trueVals-mean(trueVals)))
nmae

## [1] 0.65633

mape <- mean(abs(trueVals-preds)/trueVals)
mape


## [1] 0.290773

smape <- 1/length(preds) * sum(abs(preds - trueVals) /
                              (abs(preds)+abs(trueVals)))
smape

## [1] 0.5250418

corr <- cor(trueVals,preds)
corr

## [1] 0.6745381
```

 DE VALENCIA

Multiple Linear Regression

Multiple Linear Regression

- Multiple linear regression is probably the most used statistical method
- It is one of the many possible approaches to the multiple regression problem where given a training data set $\mathbf{D} = \{\langle \mathbf{x}_i, y_i \rangle\}_{i=1}^n$ we want to obtain an approximation of the unknown regression function $f()$ that maps the predictors values into a target continuous variable value.
- In matrix notation we have $\mathbf{D} = \mathbf{X}|\mathbf{Y}$, where \mathbf{X} is a matrix $n \times p$, and \mathbf{Y} is a matrix $n \times 1$.



Multiple Linear Regression (cont.)

- A regression model $r_D(.)$ can be seen as a function that transforms a vector of values of the predictors, \mathbf{x} , into a real number, Y . This model is an approximation of the unknown $f()$ function.
- Regression models assume the following relationship, $y_i = r(\beta, \mathbf{x}_i) + \epsilon_i$, where $r(\beta, \mathbf{x}_i)$ is a regression model with parameters β and ϵ_i are observation errors.
- The goal of a learning method is to obtain the model parameters β that minimize a certain preference criterion.



Multiple Linear Regression (cont.)

- In the case of multiple linear regression the functional form that is assumed is the following:

$$Y = \beta_0 + \beta_1 \cdot X_1 + \dots + \beta_p \cdot X_p$$

- The goal is to find the vector of parameters β that minimizes the **sum of the squared errors**

$$\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 \cdot X_1 + \dots + \beta_p \cdot X_p))^2$$



Implementation of Multiple Linear Regression

- The minimization of the squared error leads to solving what are usually known as the **normal equations**,
 $(\mathbf{X}^T \cdot \mathbf{X}) \cdot \beta = \mathbf{X}^T \cdot Y$
- This can be solved through matrix inversion
 $\beta = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot Y$
- Matrix inversion may suffer from numerical instabilities when the matrices are singular.
- A better alternative is to use Singular Value Decomposition (SVD), which may be used to solve equations of the form $\mathbf{A} \cdot \mathbf{X} = \mathbf{b}$

Press,W.; Teukolsky,S.; Vetterling,W. and Flannery,B. (1992) : Numerical Recipes in C. Cambridge University Press.



Multiple Linear Regression

Pros and Cons

- Well-known and over-studied topic with many variants of this simple methodology (e.g. Drapper and Smith, 1981)
- Simple and effective approach when the “linearity” assumption is adequate to the data.
- Form of the model is intuitive - a set of additive effects of each variable towards the prediction
- Computationally very efficient
- Too strong assumptions on the shape of the unknown regression function

Drapper and Smith (1981): Applied Regression Analysis, 2nd edition. Wiley Series in Probability and Mathematical Statistics.



Obtaining Multiple Linear Regression Models in R

```
library(DMwR2)
data(algae)
algae <- algae[-c(62,199),] # the 2 incomplete samples
clean.algae <- knnImputation(algae) # lm() does not handle NAs!
la1 <- lm(a1 ~ ., clean.algae[,1:12])
la1

##
## Call:
## lm(formula = a1 ~ ., data = clean.algae[, 1:12])
##
## Coefficients:
## (Intercept)  seasonspring  seasonsummer  seasonwinter  sizemedium
## 42.942055    3.726978      0.747597      3.692955      3.263728
## sizesmall    speedlow    speedmedium    mxPH          mnO2
## 9.682140     3.922084     0.246764     -3.589118     1.052636
## C1           NO3           NH4           oPO4          PO4
## -0.040172   -1.511235     0.001634     -0.005435     -0.052241
## Chla
## -0.088022
```



Obtaining Multiple Linear Regression Models in R (cont.)

```
summary(la1)

##
## Call:
## lm(formula = a1 ~ ., data = clean.algae[, 1:12])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -37.679 -11.893  -2.567   7.410  62.190
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  42.942055  24.010879   1.788  0.07537 .
## seasonspring  3.726978   4.137741   0.901  0.36892
## seasonsummer  0.747597   4.020711   0.186  0.85270
## seasonwinter  3.692955   3.865391   0.955  0.34065
## sizemedium    3.263728   3.802051   0.858  0.39179
## sizesmall     9.682140   4.179971   2.316  0.02166 *
## speedlow      3.922084   4.706315   0.833  0.40573
## speedmedium  0.246764   3.241874   0.076  0.93941
## mxPH         -3.589118   2.703528  -1.328  0.18598
## mnO2          1.052636   0.705018   1.493  0.13715
## Cl            -0.040172   0.033661  -1.193  0.23426
## NO3          -1.511235   0.551339  -2.741  0.00674 **
## NH4           0.001634   0.001003   1.628  0.10516
## oPO4         -0.005435   0.039884  -0.136  0.89177
## PO4          -0.052241   0.030755  -1.699  0.09109 .
## Chl_a        0.000000   0.000000   1.100  0.27265
## Chl_b        0.000000   0.000000   1.100  0.27265
## Chl_c        0.000000   0.000000   1.100  0.27265
## Chl_d        0.000000   0.000000   1.100  0.27265
## Chl_e        0.000000   0.000000   1.100  0.27265
## Chl_f        0.000000   0.000000   1.100  0.27265
## Chl_g        0.000000   0.000000   1.100  0.27265
## Chl_h        0.000000   0.000000   1.100  0.27265
## Chl_i        0.000000   0.000000   1.100  0.27265
## Chl_j        0.000000   0.000000   1.100  0.27265
## Chl_k        0.000000   0.000000   1.100  0.27265
## Chl_l        0.000000   0.000000   1.100  0.27265
## Chl_m        0.000000   0.000000   1.100  0.27265
## Chl_n        0.000000   0.000000   1.100  0.27265
## Chl_o        0.000000   0.000000   1.100  0.27265
## Chl_p        0.000000   0.000000   1.100  0.27265
## Chl_q        0.000000   0.000000   1.100  0.27265
## Chl_r        0.000000   0.000000   1.100  0.27265
## Chl_s        0.000000   0.000000   1.100  0.27265
## Chl_t        0.000000   0.000000   1.100  0.27265
## Chl_u        0.000000   0.000000   1.100  0.27265
## Chl_v        0.000000   0.000000   1.100  0.27265
## Chl_w        0.000000   0.000000   1.100  0.27265
## Chl_x        0.000000   0.000000   1.100  0.27265
## Chl_y        0.000000   0.000000   1.100  0.27265
## Chl_z        0.000000   0.000000   1.100  0.27265
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 17.65 on 182 degrees of freedom
```

AT
CA
IA

The Diagnostic Information of the summary() Call

- Distribution of the residuals (errors) of the model - should have mean zero and should be normally distributed and as small as possible
- Estimate of each coefficient and respective standard error
- Test of the hypothesis that each coefficient is null, i.e. $H_0 : \beta_i = 0$
 - Uses a t-test
 - Calculates a t-value as β_i / s_{β_i}
 - Presents a column (Pr(>t)) with the level at which the hypothesis is rejected. A value of 0.0001 would mean that we are 99.99% confident that the coefficient is not null
 - Coefficients for which we can reject the hypothesis are tagged with a symbol

The Diagnostic Information of the summary() Call (cont.)

- We are also given the R^2 coefficients (multiple and adjusted). These coefficients indicate the degree of fit of the model to the data, i.e. the proportion of variance explained by the model. Values near 1 (100%) are better. The adjusted coefficient is more demanding as it takes into account the size of the model
- Finally, there is also a test of the hypothesis that there is no dependence of the target variable on the predictors, i.e. $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$. The F -statistic is used with this purpose. R provides the confidence level at which we are sure to reject this hypothesis. A p -level of 0.0001 means that we are 99.99% confident that the hypothesis is not true.



Simplifying the Linear Model

```
final.lal <- step(lal)
```

```
summary(final.lal)
```

```
##
## Call:
## lm(formula = a1 ~ size + mxPH + Cl + NO3 + PO4, data = clean.algae[,
##     1:12])
##
## Residuals:
##     Min       1Q   Median       3Q      Max
## -28.874 -12.732  -3.741   8.424  62.926
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  57.28555   20.96132   2.733  0.00687 **
## sizemedium    2.80050    3.40190    0.823  0.41141
## sizesmall   10.40636    3.82243    2.722  0.00708 **
## mxPH         -3.97076    2.48204   -1.600  0.11130
## Cl           -0.05227    0.03165   -1.651  0.10028
## NO3          -0.89529    0.35148   -2.547  0.01165 *
## PO4          -0.05911    0.01117   -5.291 3.32e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 17.5 on 191 degrees of freedom
## Multiple R-squared:  0.3527, Adjusted R-squared:  0.3324
## F-statistic: 17.35 on 6 and 191 DF,  p-value: 5.554e-16
```

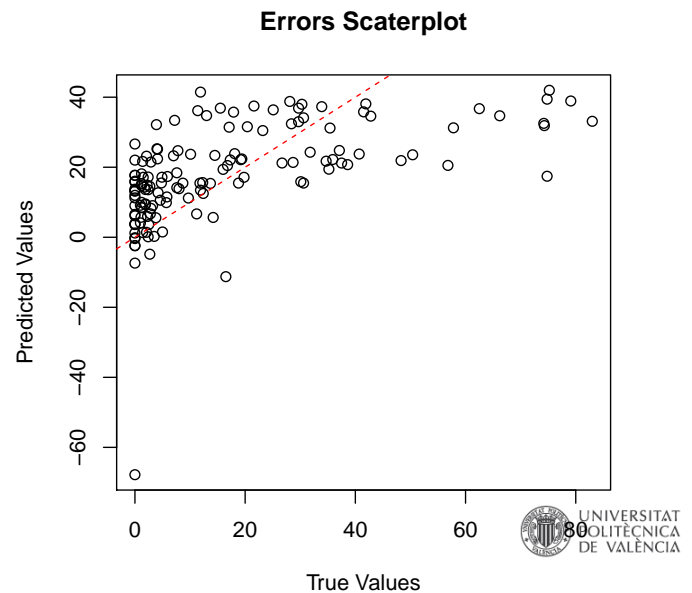

Using the Models for Prediction

```
clean.test.algae <- knnImputation(test.algae,
  k = 10, distData = clean.algae[, 1:11])
preds <- predict(final.lal, clean.test.algae)
mean((preds-algae.sols$al)^2)
```

```
## [1] 296.0934
```

```
plot(algae.sols$al, preds, main='Errors Scaterplot',
  ylab='Predicted Values', xlab='True Values')
abline(0, 1, col='red', lty=2)
```

But there are no negative algae frequencies!...



Hands on Linear Regression - the Boston data set

The data set `Boston` is available in package **MASS**. Load it and explore its help page to grab a minimal understanding of the data and then answer the following questions:

- 1 Obtain a random split of the data into two sub-sets using the proportion 70%-30%.
- 2 Obtain a multiple linear regression model using the larger set.
- 3 Check the diagnostic information provided for the model.
- 4 Obtain the predictions of the obtained model on the smaller set.
- 5 Obtain the mean squared error of these predictions and also an error scatter plot.

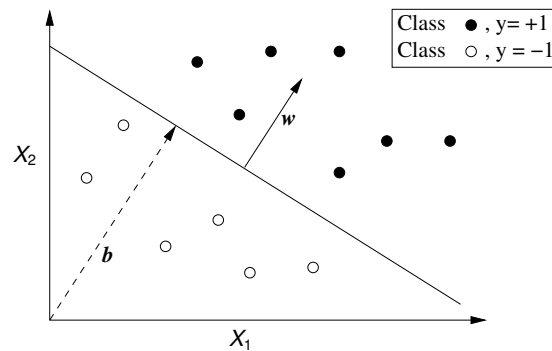
Support Vector Machines

Support Vector Machines (SVMs)

A Bit of History...

- SVM's were introduced in 1992 at the COLT-92 conference
- They gave origin to a new class of algorithms named *kernel machines*
- Since then there has been a growing interest on these methods
- More information may be obtained at www.kernel-machines.org
- A good reference on SVMs:
N. Cristianini and J. Shawe-Taylor: An introduction to Support Vector Machines. Cambridge University Press, 2000.
- SVMs have been applied with success in a wide range of areas like: bio-informatics, text mining, hand-written character recognition, etc.

Two Linearly Separable Classes

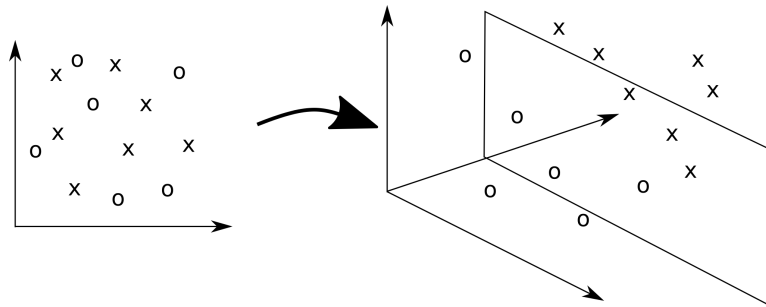


- Obtain a linear separation of the cases (binary classification problems)
- Very simple and effective for linearly separable problems
- Most real-world problems are not linearly separable!

The Basic Idea of SVMs

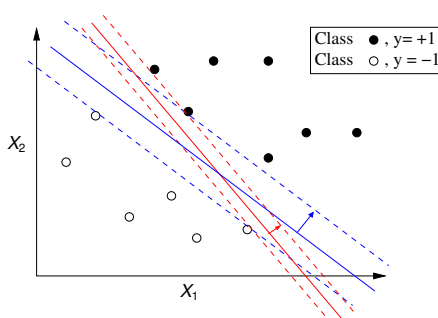
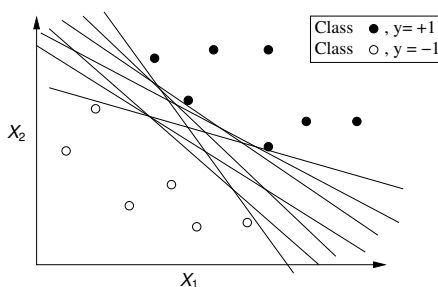
- Map the original data into a new space of variables with very high dimension.
- Use a linear approximation on this new input space.

The Idea in a Figure



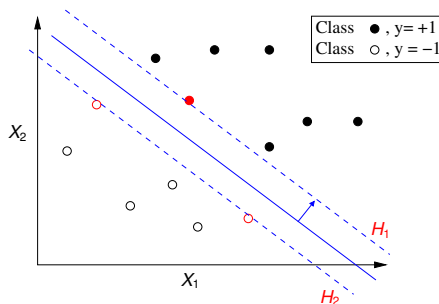
Map the original data into a new (higher dimension) coordinates system where the classes are linearly separable

Maximum Margin Hyperplane



- There is an infinite number of hyperplanes separating the two classes!
- Which one should we choose?!
- We want the one that ensures a better classification accuracy on unseen data
- SVMs approach this problem by searching for the **maximum margin hyperplane**

The Support Vectors



- All cases that fall on the hyperplanes H_1 and H_2 are called the **support vectors**.
- Removing all other cases would not change the solution!

The Optimal Hyperplane

- SVMs use quadratic optimization algorithms to find the optimal hyperplane that maximizes the margin that separates the cases from the 2 classes
- Namely, these methods are used to find a solution to the following equation,

$$L_D = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

Subject to :

$$\alpha_i \geq 0$$

$$\sum_i \alpha_i y_i = 0$$

- In the found solution, the α_i 's > 0 correspond to the support vectors that represent the optimal solution

Recap

- Most real world problems are not linearly separable
- SVMs solve this by “moving” into a extended input space where classes are already linearly separable
- This means the maximum margin hyperplane needs to be found on this new very high dimension space



The Kernel trick

- The solution to the optimization equation involves dot products that are computationally heavy on high-dimensional spaces
- It was demonstrated that the result of these complex calculations is equivalent to the result of applying certain functions (the kernel functions) in the space of the original variables.

The Kernel Trick

Instead of calculating the dot products in a high dimensional space, take advantage of the proof that $K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z})$ and simply replace the complex dot products by these simpler and efficient calculations



Summary of the SVMs Method

- As problems are usually non-linear on the original feature space, move into a high-dimension space where linear separability is possible
- Find the optimal separating hyperplane on this new space using quadratic optimization algorithms
- Avoid the heavy computational costs of the dot products using the kernel trick



How to handle more than 2 classes?

- Solve several binary classification tasks
- Essentially find the support vectors that separate each class from all others

The Algorithm

- Given a m classes task
- Obtain m SVM classifiers, one for each class
- Given a test case assign it to the class whose separating hyperplane is more distant from the test case



Obtaining an SVM in R

The package **e1071**

```
library(e1071)
data(Glass, package='mlbench')
tr <- Glass[1:200,]
ts <- Glass[201:214,]
s <- svm(Type ~ ., tr)
predict(s, ts)

## 201 202 203 204 205 206 207 208 209 210 211 212 213 214
##    7    2    7    7    7    7    7    2    7    7    7    7    7    7
## Levels: 1 2 3 5 6 7
```



Obtaining an SVM in R (2)

The package **e1071**

```
ps <- predict(s, ts)
table(ps, ts$Type)

##
## ps    1    2    3    5    6    7
##    1    0    0    0    0    0    0
##    2    0    0    0    0    0    2
##    3    0    0    0    0    0    0
##    5    0    0    0    0    0    0
##    6    0    0    0    0    0    0
##    7    0    0    0    0    0   12

mc <- table(ps, ts$Type)
error <- 100*(1-sum(diag(mc)))/sum(mc)
error

## [1] 14.28571
```



ε -SV Regression

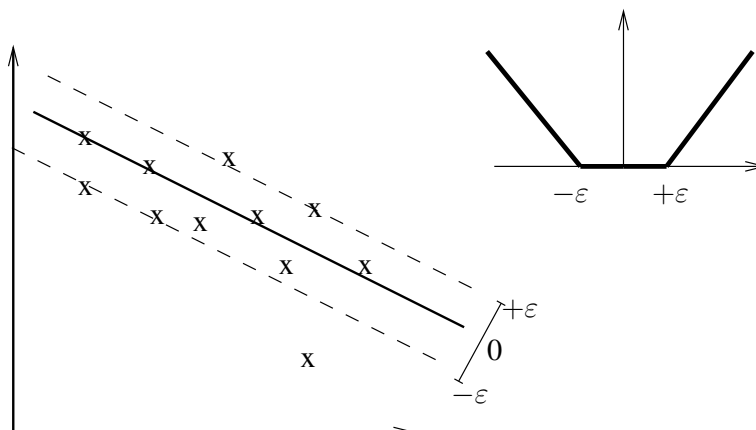
- Vapnik (1995) proposed the notion of ε support vector regression
- The goal in ε -SV Regression is to find a function $f(x)$ that has at most ε deviation from the given training cases
- In other words we do not care about errors smaller than ε

V. Vapnik (1995). The Nature of Statistical Learning Theory. Springer.

ε -SV Regression (cont.)

- ε -SV Regression uses the following error metric,

$$|\xi|_{\varepsilon} = \begin{cases} 0 & \text{if } |\xi| \leq \varepsilon \\ |\xi| - \varepsilon & \text{otherwise} \end{cases}$$



ε -SV Regression (cont.)

- The theoretical development of this idea leads to the following optimization problem,

$$\begin{aligned} \text{Minimize : } & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^l (\xi_i + \xi_i^*) \\ \text{Subject to : } & \begin{cases} y_i - \mathbf{w} \cdot \mathbf{x} - b \leq \varepsilon + \xi_i \\ \mathbf{w} \cdot \mathbf{x} + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases} \end{aligned}$$

where C corresponds to the cost to pay for each violation of the error limit ε



ε -SV Regression (cont.)

- As within classification we use the kernel trick to map a non-linear problem into a high dimensional space where we solve the same quadratic optimization problem as in the linear case
- In summary, by the use of the $|\xi|_\varepsilon$ loss function we reach a very similar optimization problem to find the support vectors of any non-linear regression problem.

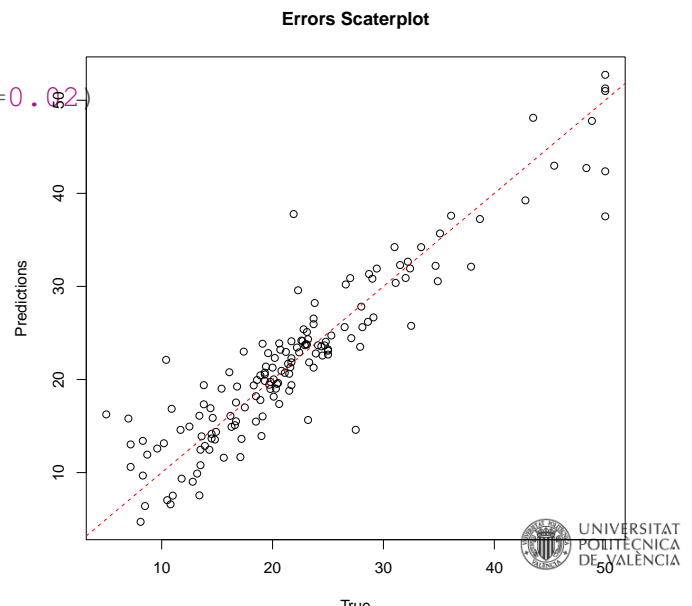


SVMs for regression in R

```
library(e1071)
data(Boston, package='MASS')
set.seed(1234)
sp <- sample(1:nrow(Boston), 354)
tr <- Boston[sp, ]
ts <- Boston[-sp, ]
s <- svm(medv ~ ., tr, cost=10, epsilon=0.02)
preds <- predict(s, ts)
mean((ts$medv-preds)^2)

## [1] 13.82111
```

```
plot(ts$medv, preds, main='Errors Scaterplot',
      ylab='Predictions', xlab='True')
abline(0, 1, col='red', lty=2)
```



Hands On SMVs

Hands on SVMs

The file `Wine.Rdata` contains 2 data frames with data about the quality of “green” wines: i) `redWine` and ii) `whiteWine`. Each of these data sets has information on a series of wine tasting sessions to “green” wines (both red and white). For each wine sample several physico-chemical properties of the wine sample together with a quality score assigned by a committee of wine experts (variable `quality`).

- 1 Obtain and SVM for forecasting the quality of the red variant of “green” wines
- 2 Split the data set in two parts: one with 70% of the samples and the other with the remaining 30%. Obtain an SVM with the first part and apply it to the second. What was the resulting mean absolute error?
- 3 Using the `round()` function, round the predictions obtained in the previous question to the nearest integer. Calculate the error rate of the resulting integers when compared to the true values

k -Nearest Neighbors

k-Nearest Neighbors

k Nearest Neighbors

- The k -nearest neighbor method was first described in the early 1950s.
- This method is computationally intensive with large data sets and it did not enjoy lots of popularity because of this.
- With the advent of cheap computing power its popularity has increased a lot because it is a very simple and effective method that can easily handle both classification and regression problems.
- k -nearest neighbors can be seen as methods that learn by analogy - i.e. they are based on the notion of similarity between cases.

k Nearest Neighbors (cont.)

The Basic Idea

- If we are given a new test case \mathbf{x} for which we want a prediction
 - 1 search in the training set for the most similar cases (the nearest neighbors) to \mathbf{x}
 - 2 use the outcomes of these cases to obtain the prediction for \mathbf{x}

k Nearest Neighbors

Main Characteristics

- The k-nearest neighbors are known as lazy learners as they do not learn any model of the data
- Learning in k-nearest neighbors consists simply in storing the training data
 - Variants here include storing in data structures that provide efficient querying of the nearest neighbors
- They do not make any assumption on the unknown functional form we are trying to approximate, which means that with sufficient data they are applicable to any problem
- They usually achieve good results but...
 - They require a proper distance metric to be defined - issues like normalization, irrelevant variables, unknown values, etc., may have a strong impact on their performance
- They have fast training time, but slow prediction time

The Notion of Similarity

- The key issue on kNN is the notion of similarity
- This notion is strongly related with the notion of **distance between observations**
- Distances among observations in a data set can be used to find the neighbors of a test case



How to Calculate the Distance between 2 Cases?

- The notion of distance is related to the differences between the values on the variables describing the cases

ID	Income	Sex	Position	Age
1	2500	f	manager	35
2	2750	f	manager	30
3	4550	m	director	50

Case 1 is “closer” to case 2 than to 3



The Euclidean Distance Function

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$$

where x_i is the value of case \mathbf{x} on variable i

Example

Given two cases $\mathbf{x} = \langle 3, 5, 1 \rangle$ and $\mathbf{y} = \langle 12, 5.4, -3 \rangle$ their Euclidean distance is given by

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(3 - 12)^2 + (5 - 5.4)^2 + (1 - (-3))^2} = 9.85697$$



A Generalization - the Minkowski distance

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^p |x_i - y_i|^r \right)^{1/r}$$

where if

- $r = 1$ we have what is known as the Manhattan distance (or L_1 -norm)
- $r = 2$ we have the Euclidean distance
- etc.



Potential Problems with Distance Calculation

In domains where cases are described by many variables several problems may arise that may distort the notion of distance between any two cases.

- Different scales of variables
- Different importance of variables
- Different types of data (e.g. both numeric and nominal variables, etc.)
- etc.



Heterogeneous Distance Functions

How to calculate the distance between two cases described by variables with different type (e.g. numeric and nominal variables)?
A possible solution,

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^p \delta_i(x_i, y_i)$$

emque,

$$\delta_i(v_1, v_2) = \begin{cases} 1 & \text{if } i \text{ is nominal e } v_1 \neq v_2 \\ 0 & \text{if } i \text{ is nominal e } v_1 = v_2 \\ \frac{|v_1 - v_2|}{\text{range}(i)} & \text{if } i \text{ is numeric} \end{cases}$$

The distance between $\langle 2500, f, \text{director}, 35 \rangle$ and $\langle 2750, f, \text{director}, 30 \rangle$ would be given by $\frac{|2500-2750|}{\text{range}(\text{Salary})} + 0 + 0 + \frac{|35-30|}{\text{range}(\text{Age})}$



1-Nearest Neighbor Classifier

Method

- Search for the training case most similar to the test case
 - Predict for the test case the class of this nearest neighbor
-
- Very simple method
 - May suffer with the presence of outliers
 - Frequently achieves good results



k-Nearest Neighbor Classifier

- Use the k nearest neighbors to obtain the classification of the test case
- Namely, the majority class on the k neighbors is the prediction of the method
- What should be the value of k ?
 - Frequent values are 3, 5 and 7
 - Odd numbers to avoid draws!
 - It can be estimated experimentally
 - Global estimation searches for the ideal k for a given data set
 - Local estimation methods try to estimate the ideal k for each test case (computationally very demanding!)



k-nearest neighbors in R

- Package **class** contains function `knn()`

```
library(class)
set.seed(1234)
sp <- sample(1:150, 100)
tr <- iris[sp,]
ts <- iris[-sp,]
nn3 <- knn(tr[, -5], ts[, -5], tr[, 5], k=3)
(mtrx <- table(nn3, ts$Species))

##
## nn3          setosa versicolor virginica
## setosa          12           0           0
## versicolor       0           20          1
## virginica        0           1          16

(err <- 1-sum(diag(mtrx))/sum(mtrx))

## [1] 0.04
```



k-nearest neighbors in R - 2

- Package **DMwR2** has a wrapper function with a “standard” interface,

```
library(class)
library(DMwR2)
set.seed(1234)
sp <- sample(1:150, 100)
tr <- iris[sp,]
ts <- iris[-sp,]
nn3 <- kNN(Species ~ ., tr, ts, k=3, stand=TRUE)
(mtrx <- table(nn3, ts$Species))

##
## nn3          setosa versicolor virginica
## setosa          12           0           0
## versicolor       0           20          3
## virginica        0           1          14

(err <- 1-sum(diag(mtrx))/sum(mtrx))

## [1] 0.08
```



Trying to find the “ideal” value of k in R

```

trials <- c(1,3,5,7,11,13,15)
nreps <- 10
res <- matrix(NA,nrow=length(trials),ncol=2)
for(k in seq_along(trials)) {
  errs <- rep(0,nreps)
  for(r in 1:nreps) {
    sp <- sample(1:150,100)
    tr <- iris[sp,]
    ts <- iris[-sp,]
    nn3 <- kNN(Species ~ .,tr,ts,k=trials[k],norm=TRUE)
    mtrx <- table(nn3,ts$Species)
    errs[r] <- 1-sum(diag(mtrx))/sum(mtrx)
  }
  res[k,] <- c(mean(errs),sd(errs))
}
dimnames(res) <- list(paste('k',trials,sep=' '),c('avg','std'))
res

##      avg      std
## k=1  0.082 0.02741
## k=3  0.050 0.02708
## k=5  0.056 0.01265
## k=7  0.052 0.03676
## k=11 0.040 0.02108
## k=13 0.046 0.02503
## k=15 0.070 0.02539

```

k-Nearest Neighbor Regression

Method

- Search for the training case most similar to the test case
- Predict for the test case the **average** of the target variable values of the neighbors

k-nearest neighbors regression in R

- Package **caret** has a function that to obtain these models

```
library(caret)
data(Boston, package="MASS")
set.seed(1234)
sp <- sample(1:506, 354)
tr <- Boston[sp, ]
ts <- Boston[-sp, ]
tgt <- which(colnames(Boston) == "medv")
nn3 <- knnreg(tr[, -tgt], tr[, tgt], k=3)
pnn3 <- predict(nn3, ts[, -tgt])
(mse <- mean((pnn3-ts[, tgt])^2))

## [1] 44.93
```

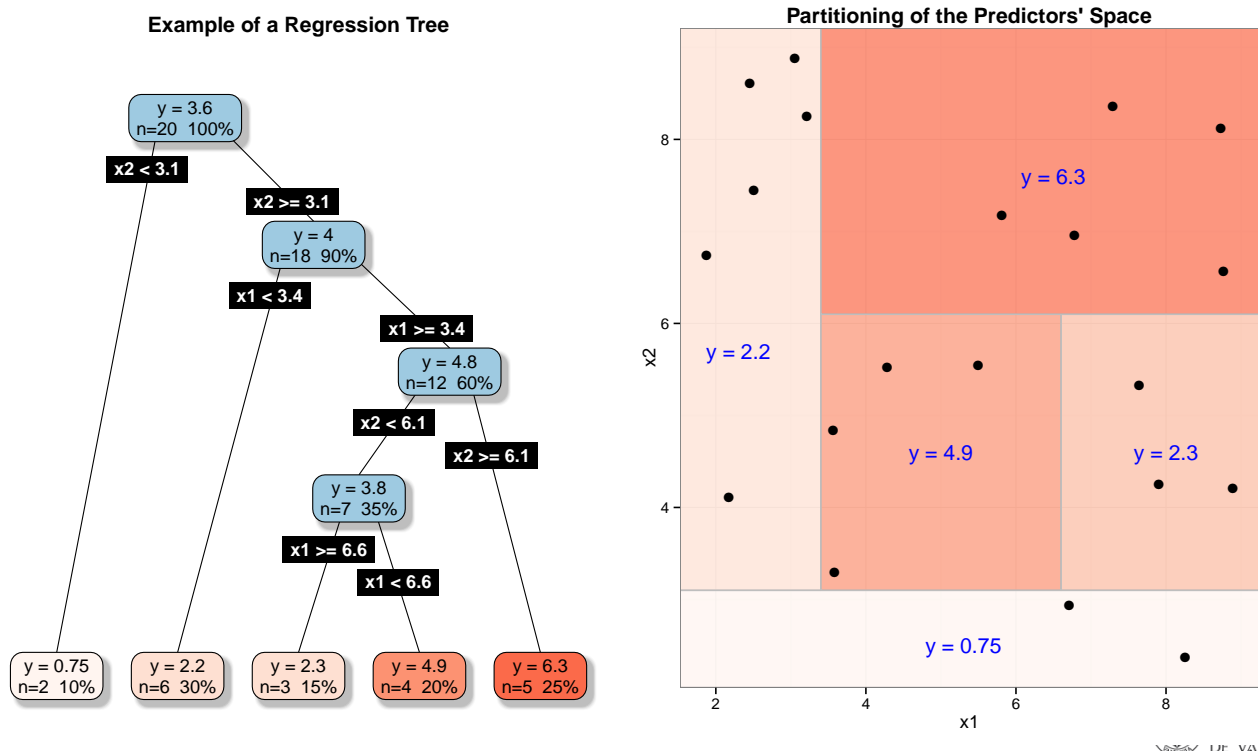
Tree-based Models

Tree-based Models

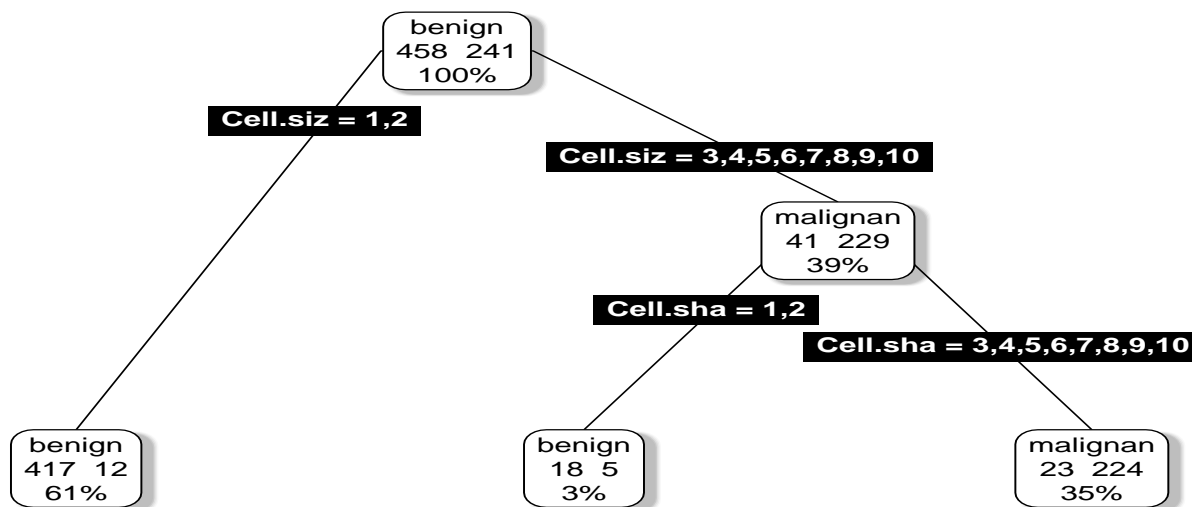
- Tree-based models (both classification and regression trees) are models that provide as result a model based on logical tests on the input variables
- These models can be seen as a partitioning of the input space defined by the input variables
- This partitioning is defined based on carefully chosen logical tests on these variables
- Within each partition all cases are assigned the same prediction (either a class label or a numeric value)
- Tree-based models are known by their (i) computational efficiency; (ii) interpretable models; (iii) embedded variable selection; (iv) embedded handling of unknown variable values and (v) few assumptions on the unknown function being approximated



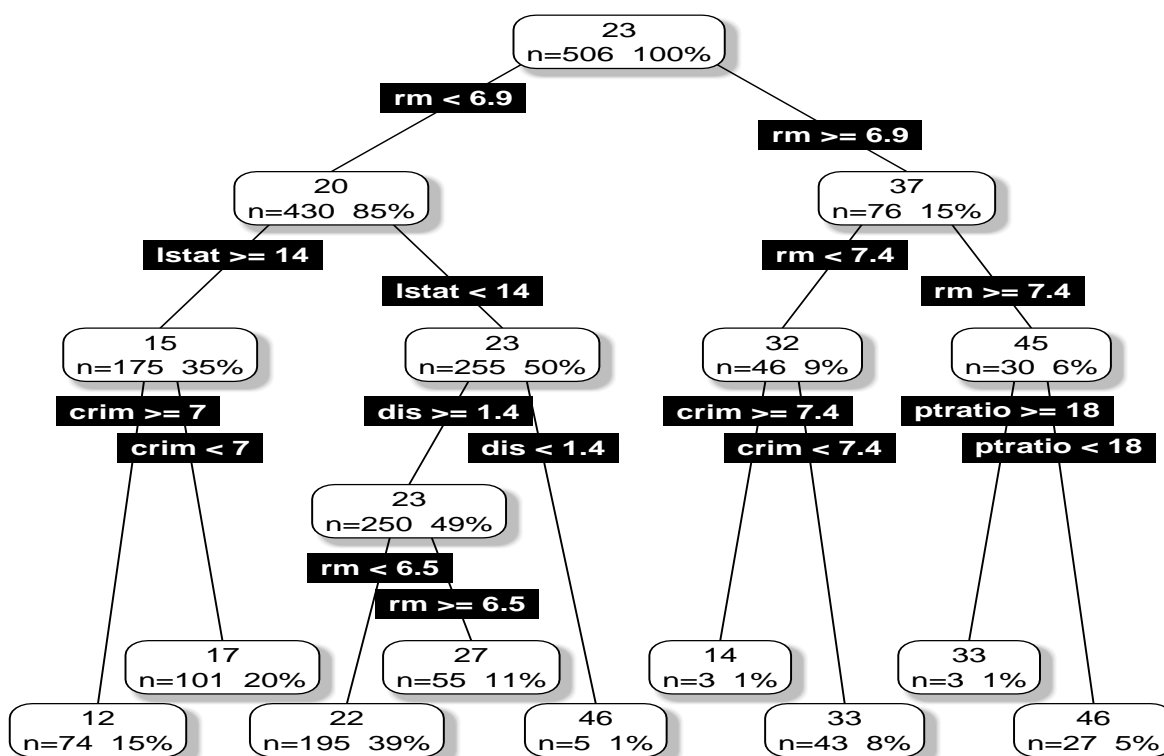
An Example of Trees Partitioning



An Example of a Classification Tree



An Example of a Regression Tree



Tree-based Models

- Most tree-based models are binary trees with logical tests on each node
- Tests on numerical predictors take the form $x_i < \alpha$, with $\alpha \in \mathcal{R}$
- Tests on nominal predictors take the form $x_j \in \{v_1, \dots, v_m\}$
- Each path from the top (root) node till a leaf can be seen as a logical condition defining a region of the predictors space.
- All observations “falling” on a leaf will get the same prediction
 - the majority class of the training cases in that leaf for classification trees
 - the average value of the target variable for regression trees
- The prediction for a new test case is easily obtained by following a path from the root till a leaf according to the case predictors values



The Recursive Partitioning Algorithm

```

1: function RECURSIVEPARTITIONING( $D$ )
   Input :  $D$ , a sample of cases,  $\{\langle x_{i,1}, \dots, x_{i,p}, y_i \rangle\}_{i=1}^{N_{train}}$ 
   Output :  $t$ , a tree node

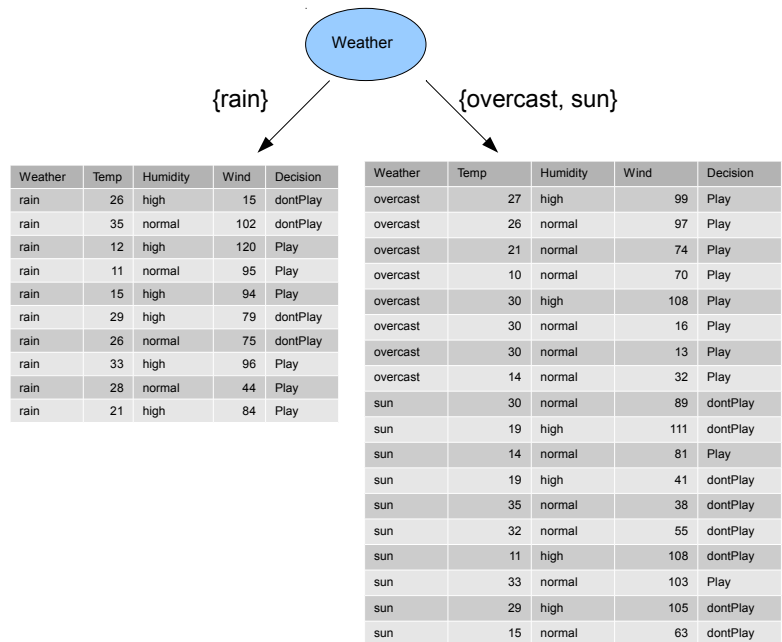
2:   if <TERMINATION CRITERION> then
3:     Return a leaf node with the majority class in  $D$ 
4:   else
5:      $t \leftarrow$  new tree node
6:      $t.split \leftarrow$  <FIND THE BEST PREDICTORS TEST>
7:      $t.leftNode \leftarrow$  RecursivePartitioning( $\mathbf{x} \in D : \mathbf{x} \models t.split$ )
8:      $t.rightNode \leftarrow$  RecursivePartitioning( $\mathbf{x} \in D : \mathbf{x} \not\models t.split$ )
9:     Return the node  $t$ 
10:  end if
11: end function

```



The Recursive Partitioning Algorithm - an example

Weather	Temp.	Humidity	Wind	Decision
rain	26	high	15	dontPlay
rain	35	normal	102	dontPlay
overcast	27	high	99	Play
overcast	26	normal	97	Play
rain	12	high	120	Play
overcast	21	normal	74	Play
sun	30	normal	89	dontPlay
sun	19	high	111	dontPlay
sun	14	normal	81	Play
overcast	10	normal	70	Play
rain	11	normal	95	Play
rain	15	high	94	Play
sun	19	high	41	dontPlay
sun	35	normal	38	dontPlay
rain	29	high	79	dontPlay
rain	26	normal	75	dontPlay
overcast	30	high	108	Play
overcast	30	normal	16	Play
rain	33	high	96	Play
overcast	30	normal	13	Play
sun	32	normal	55	dontPlay
sun	11	high	108	dontPlay
sun	33	normal	103	Play
overcast	14	normal	32	Play
rain	28	normal	44	Play
rain	21	high	84	Play
sun	29	high	105	dontPlay
sun	15	normal	63	dontPlay



The Recursive Partitioning Algorithm (cont.)

Key Issues of the RP Algorithm

- When to stop growing the tree - termination criterion
- Which value to put on the leaves
- How to find the best split test



The Recursive Partitioning Algorithm (cont.)

When to Stop?

Too large trees tend to **overfit** the training data and will perform badly on new data - a question of reliability of error estimates

Which value?

Should be the value that better represents the cases in the leaves

What are the good tests?

A test is good if it is able to **split** the cases of sample in such a way that they form partitions that are “purer” than the parent node



Classification vs Regression Trees

- They are both grown using the Recursive Partitioning algorithm
- The main difference lies on the used preference criterion
- This criterion has impact on:
 - The way the best test for each node is selected
 - The way the tree avoids over fitting the training sample
- Classification trees typically use criteria related to error rate (e.g. the Gini index, the Gain ratio, entropy, etc.)
- Regression trees typically use the least squares error criterion



How to Evaluate a Test in Classification Trees?

Gini Impurity

- The Gini index of a data set D where each example belongs to one of c classes is given by,

$$Gini(D) = 1 - \sum_{i=1}^c p_i^2$$

where p_i is the probability of class i usually estimated with the observed frequency on the training data

- If the data set is split on a logical test T then the resulting Gini index is given by,

$$Gini_T(D) = \frac{|D_T|}{|D|} Gini(D_T) + \frac{|D_{\neg T}|}{|D|} Gini(D_{\neg T})$$

- In this context the reduction in impurity given by T is,

$$\Delta Gini_T(D) = Gini(D) - Gini_T(D)$$



Gini Impurity - an example

Weather	Temp	Humidity	Wind	Decision
rain	26	high	15	dontPlay
rain	35	normal	102	dontPlay
rain	12	high	120	jogar
rain	11	normal	95	jogar
rain	15	high	94	jogar
rain	29	high	79	dontPlay
rain	26	normal	75	dontPlay
rain	33	high	96	jogar
rain	28	normal	44	jogar
rain	21	high	84	jogar
overcast	27	high	99	jogar
overcast	26	normal	97	jogar
overcast	21	normal	74	jogar
overcast	10	normal	70	jogar
overcast	30	high	108	jogar
overcast	30	normal	16	jogar
overcast	30	normal	13	jogar
overcast	14	normal	32	jogar
sunny	30	normal	89	dontPlay
sunny	19	high	111	dontPlay
sunny	14	normal	81	jogar
sunny	19	high	41	dontPlay
sunny	35	normal	38	dontPlay
sunny	32	normal	55	dontPlay
sunny	11	high	108	dontPlay
sunny	33	normal	103	jogar
sunny	29	high	105	dontPlay
sunny	15	normal	63	dontPlay

$$Gini(D) = 1 - \left(\left(\frac{16}{16+12} \right)^2 + \left(\frac{12}{16+12} \right)^2 \right) = 0.49$$

$$Gini_{Weather \in \{rain\}}(D) = \frac{10}{28} \cdot Gini(D_{Weather \in \{rain\}}) + \frac{18}{28} \cdot Gini(D_{Weather \notin \{rain\}}) = 0.489$$

$$Gini(D_{Weather \in \{rain\}}) = 1 - \left(\left(\frac{4}{4+6} \right)^2 + \left(\frac{6}{4+6} \right)^2 \right) = 0.48$$

$$Gini(D_{Weather \notin \{rain\}}) = 1 - \left(\left(\frac{8}{8+10} \right)^2 + \left(\frac{10}{8+10} \right)^2 \right) = 0.49$$

$$\Delta Gini_{Weather \in \{rain\}}(D) = 0.49 - 0.489 = 0.001$$

Calculate the value of $\Delta Gini_{Weather \in \{overcast\}}(D)$



Which Tests are Tried?

Numeric Predictors

- Given a set of data D and a continuous variable A let $V_{A,D}$ be the set of values of A occurring in D
- Start by ordering the set $V_{A,D}$
- Evaluate all tests $A < x$ where x takes as values all mid-points between every successive value in the ordered set



Numeric Predictors - an example

- Given the unsorted values of $Temp$:
26 35 27 26 12 21 30 19 14 10 11 15 19 35 29 26 30 30 33 30 32
11 33 14 28 21 29 15
- Start by ordering them:
10 11 11 12 14 14 15 15 19 19 21 21 26 26 26 27 28 29 29 30 30
30 30 32 33 33 35 35
- Then try (i.e. evaluate) all tests in between each value:
 - $Temp < 10.5$
 - $Temp < 11.5$
 - $Temp < 13$
 - $Temp < 14.5$
 - etc.
- Choose the test with the best score to be the best test in variable $Temp$



Which Tests are Tried?

Nominal Predictors

- Given a set of data D and a nominal variable A let $V_{A,D}$ be the set of values of A occurring in D
- Evaluate all possible combinations of subset of values in $V_{A,D}$
- Note that there are some optimizations that reduce the computational complexity of this search



Nominal Predictors - an example

- Given the values of `Weather`:
rain, overcast, sunny
- Try (i.e. evaluate) all subsets of these values:
 - $Weather \in \{rain\}$
 - $Weather \in \{overcast\}$
 - $Weather \in \{sunny\}$
- Choose the test with the best score to be the best test in variable *Weather*



How to Evaluate a Test in Regression Trees?

Least Squares Regression Trees

- Regression trees are usually grown by trying to **minimize the sum of the squared errors**, leading to *Least Squares Regression Trees*
- According to the LS (Least Squares) criterion the error of the cases D in a node of the tree is given by,

$$Err(D) = \frac{1}{|D|} \sum_{\langle \mathbf{x}_i, y_i \rangle \in D} (y_i - k)^2$$

where D is the sample of cases in a node t , $|D|$ is the cardinality of this set and k is a constant used to represent the cases in the node

- It can be easily proven that the constant k that minimizes this error is the **average target variable value** of the cases in the node



Least Squares Regression Trees

- Any logical test T divides the cases D of a node in two partitions, D_T and D_{-T} . The resulting pooled error is given by,

$$Err_T(D) = \frac{|D_T|}{|D|} \times Err(D_T) + \frac{|D_{-T}|}{|D|} \times Err(D_{-T})$$

where $|D_T|/|D|$ ($|D_{-T}|/|D|$) is the proportion of cases going to the left (right) branch of the node containing the cases D

- We can estimate the value of the split test T by the respective error reduction,

$$\Delta Err_T(D) = Err(D) - Err_T(D)$$

- Finding the best split test for a node involves evaluating all possible tests for this node using the above equations



Least Squares Regression Trees (cont.)

- For continuous variables this requires a sorting operation on the values of this variable occurring in the node
- After this sorting, a fast incremental algorithm (Torgo, 1999) can be used to find the best cut-point value for the test
- With respect to nominal variables, Breiman and colleagues (1984) have proved a theorem that avoids trying all possible combinations of values, reducing the computational complexity of this task from $O(2^v - 1)$ to $O(v - 1)$, where v is the number of values of the nominal variable

Breiman et al. (1984): Classification and Regression Trees

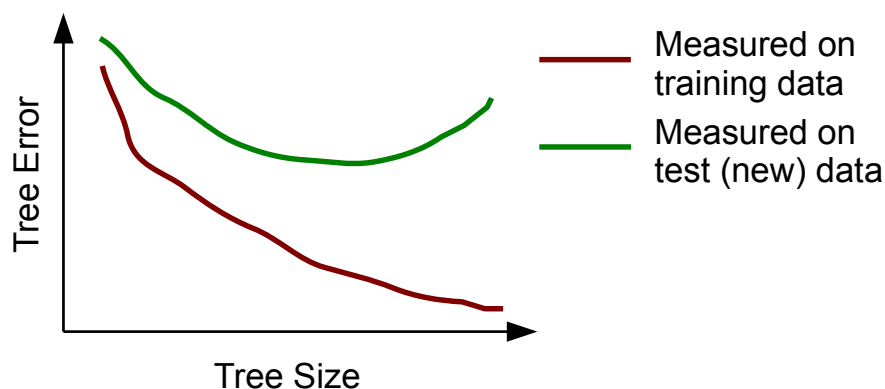
Torgo L. (1999): Inductive learning of tree-based regression models. PhD thesis, Department of Computer Science, Faculty of Sciences, University of Porto.

Torgo,L. (2011) : Regression Trees. In Encyclopedia of Machine Learning, C.Sammut and G.I.Webb (Eds.). Pages 842–845, Springer.



Deciding when to stop growing the trees

- The scores discussed before keep improving as we grow the tree
- At an extreme, an overly large tree, will perfectly fit the given training data (i.e. all cases are correctly predicted by the tree)
- Such huge trees are said to be overfitting the training data and will most probably perform badly on a new set of data (a test set), as they have captured spurious characteristics of the training data



Deciding when to stop growing the trees - 2

- As we go down in the tree the decisions on the tests are made on smaller and smaller sets, and thus potentially less reliable decisions are made
- The standard procedure in tree learning is to grow an overly large tree and then use some statistical procedure to **prune** unreliable branches from this tree. The goal of this procedure is to try to obtain reliable estimates of the error of the tree. This procedure is usually called **post-pruning** a tree.
- An alternative procedure (not so frequently used) is to decide during tree growth when to stop. This is usually called **pre-pruning**.



(Post-)Pruning a Tree

Cost-complexity and Error-complexity Pruning

- Grown and overly large tree
- Generate a sequence of sub-trees
 - Error-complexity criterion for regression trees
 - Cost-complexity criterion for classification trees
- Use cross validation to estimate the error of these trees
- Use the x -SE rule to select the best sub-tree



Classification and Regression Trees in R

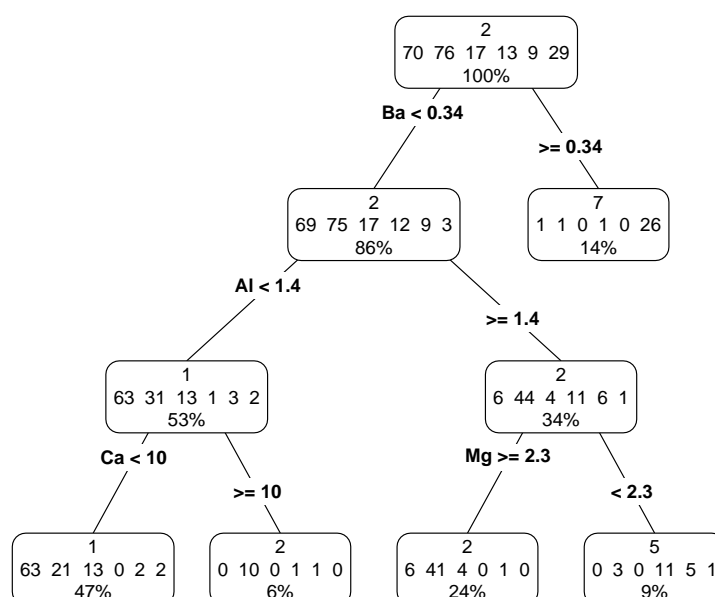
The package `rpart`

- Package `rpart` implements most of the ideas of the system CART that was described in the book “Classification and Regression Trees” by Breiman and colleagues
- This system is able to obtain classification and regression trees.
- For classification trees it uses the Gini score to grow the trees and it uses Cost-Complexity post-pruning to avoid over fitting
- For regression trees it uses the least squares error criterion and it uses Error-Complexity post-pruning to avoid over fitting
- On package `DMwR2` you may find function `rpartXse()` that grows and prunes a tree in a way similar to CART using the above infra-structure



Illustration using a classification task - *Glass*

```
library(DMwR2)
library(rpart.plot)
data(Glass, package='mlbench')
ac <- rpartXse(Type ~ ., Glass)
prp(ac, type=4, extra=101)
```



How to use the trees for Predicting?

```
tr <- Glass[1:200,]
ts <- Glass[201:214,]
ac <- rpartXse(Type ~ .,tr)
predict(ac,ts)

##      1      2      3      5 6      7
## 201 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 202 0 0.3636364 0.6363636 0.00000000 0 0.0000000
## 203 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 204 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 205 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 206 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 207 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 208 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 209 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 210 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 211 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 212 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 213 0 0.0000000 0.0000000 0.09090909 0 0.9090909
## 214 0 0.0000000 0.0000000 0.09090909 0 0.9090909
```



How to use the trees for Predicting? (cont.)

```
predict(ac,ts,type='class')

## 201 202 203 204 205 206 207 208 209 210 211 212 213 214
## 7 3 7 7 7 7 7 7 7 7 7 7 7 7
## Levels: 1 2 3 5 6 7

ps <- predict(ac,ts,type='class')
table(ps,ts$Type)

##
## ps 1 2 3 5 6 7
## 1 0 0 0 0 0 0
## 2 0 0 0 0 0 0
## 3 0 0 0 0 0 1
## 5 0 0 0 0 0 0
## 6 0 0 0 0 0 0
## 7 0 0 0 0 0 13

mc <- table(ps,ts$Type)
err <- 100*(1-sum(diag(mc))/sum(mc))
err

## [1] 7.142857
```

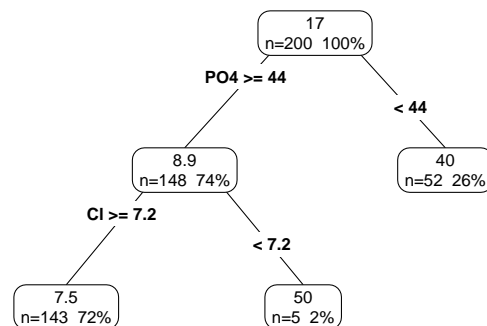


Illustration using a regression task

Forecasting Algae a1

```
prp(ar, type=4, extra=101)
```

```
library(DMwR2)
library(rpart.plot)
data(algae)
d <- algae[, 1:12]
ar <- rpartXse(a1 ~ ., d)
```



How to use the trees for Predicting?

```
tr <- d[1:150,]
ts <- d[151:200,]
ar <- rpartXse(a1 ~ ., tr)
preds <- predict(ar, ts)
mae <- mean(abs(preds - ts$a1))
mae

## [1] 12.27911

cr <- cor(preds, ts$a1)
cr

## [1] 0.512407
```



Hands on Tree-based Models - the Wines data

File `Wine.Rdata` contains two data frames with data on green wine quality: (i) `redWine` and (ii) `whiteWine`. Each of these data sets contains a series of tests with green wines (red and white). For each of these tests the values of several physicochemical variables together with a quality score assigned by wine experts (column `quality`).

- 1 Build a regression tree for the white wines data set
- 2 Obtain a graph of the obtained regression tree
- 3 Apply the tree to the data used to obtain the model and calculate the mean squared error of the predictions
- 4 Split the data set in two parts: 70% of the tests and the remaining 30%. Using the larger part to obtain a regression tree and apply it to the other part. Calculate again the mean squared error. Compare with the previous scores and comment.



Model Ensembles and Random Forests

Model Ensembles

What?

- Ensembles are collections of models that are used together to address a certain prediction problem

Why? (Diettrich, 2002)

- For complex problems it is hard to find a model that “explains” all observed data.
- Averaging over a set of models typically leads to significantly better results.

Diettrich, T. G. (2002). Ensemble Learning. In *The Handbook of Brain Theory and Neural Networks*, Second edition, (M.A. Arbib, Ed.), Cambridge, MA: The MIT Press, 2002. 405-408.



The Bias-Variance Decomposition of Prediction Error

- The prediction error of a model can be split in two main components: the bias and the variance components
- The bias component is the part of the error that is due to the poor ability of the model to fit the seen data
- The variance component has to do with the sensibility of the model to the given training data



The Bias-Variance Decomposition of Prediction Error

- Decreasing the bias by adjusting more to the training sample will most probably lead to a higher variance - the over-fitting phenomenon
- Decreasing the variance by being less sensitive to the given training data will most probably have as consequence a higher bias
- In summary: there is a well-known bias-variance trade-off in learning a prediction model

Ensembles are able to reduce both components of the error

Their approach consist on applying the same algorithm to different samples of the data and use the resulting models in a voting schema to obtain predictions for new cases



Random Forests (Breiman, 2001)

- Random Forests put the ideas of sampling the cases and sampling the predictors, together in a single method
 - Random Forests combine the ideas of bagging together with the idea of random selection of predictors
- Random Forests consist of sets of tree-based models where each tree is obtained from a bootstrap sample of the original data and uses some form of random selection of variables during tree growth

Breiman, L. (2001): "Random Forests". Machine Learning 45 (1): 5—32.



Random Forests - the algorithm

- For each of the k models
 - Draw a random sample with replacement to obtain the training set
 - Grow a classification or regression tree
 - On each node of the tree choose the best split from a randomly selected subset m of the predictors
- The trees are fully grown, i.e. no pruning is carried out



Random Forests in R

The package `randomForest`

```
library(randomForest)
data(Boston, package="MASS")
samp <- sample(1:nrow(Boston), 354)
tr <- Boston[samp, ]
ts <- Boston[-samp, ]
m <- randomForest(medv ~ ., tr)
ps <- predict(m, ts)
mean(abs(ts$medv-ps))

## [1] 2.378855
```



A classification example

```
data(Glass, package='mlbench')
set.seed(1234)
sp <- sample(1:nrow(Glass), 150)
tr <- Glass[sp,]
ts <- Glass[-sp,]
m <- randomForest(Type ~ ., tr, ntree=3000)
ps <- predict(m, ts)
table(ps, ts$Type)

##
## ps  1  2  3  5  6  7
##  1 13  5  3  0  0  1
##  2  2 18  0  3  0  2
##  3  0  0  1  0  0  0
##  5  0  0  0  4  0  0
##  6  0  1  0  0  3  0
##  7  0  0  0  0  0  8

mc <- table(ps, ts$Type)
err <- 100*(1-sum(diag(mc))/sum(mc))
err

## [1] 26.5625
```



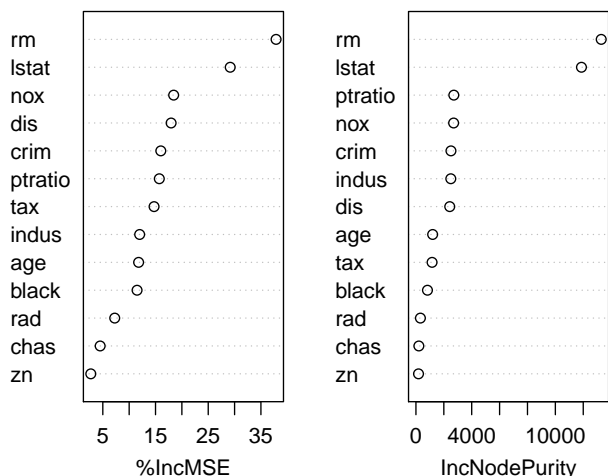
Other Uses of Random Forests

Variable Importance

```
data(Boston, package='MASS')
library(randomForest)
m <- randomForest(medv ~ ., Boston,
                  importance=T)
importance(m)

##           %IncMSE  IncNodePurity
## crim    16.001604    2511.3914
## zn       2.719681     184.4274
## indus   11.992644    2501.0269
## chas     4.496731     208.3667
## nox     18.440180    2702.4705
## rm      37.873226   13288.7533
## age     11.793865    1198.7370
## dis     17.957678    2423.8487
## rad      7.259293     320.4829
## tax     14.721102    1157.0856
## ptratio 15.715445    2716.8744
## black   11.498495     826.2531
## lstat   29.172401   11871.6578
```

Feature Relevance Scores



```
varImpPlot(m, main="Feature Relevance Scores")
```



Hands on Linear Regression and Random Forests

the Algae data set

Load in the data set `algae` from package **DMwR2** and answer the following questions:

- 1 How would you obtain a random forest to forecast the value of alga *a4*
- 2 Repeat the previous exercise but now using a linear regression model. Try to simplify the model using the `step()` function.
- 3 Obtain the predictions of the two previous models for the data used to obtain them. Draw a scatterplot comparing these predictions
- 4 The data frame named `test.algae` contains a test set with some extra 140 water samples for which we want predictions. Use the previous two models to obtain predictions for *a4* on these new samples. Check what happened to the test cases with NA's. Fill-in the NA's on the test set and repeat the experiment.



Evaluation Methodologies and Comparison of Models

Performance Estimation

The setting

- Predictive task: unknown function $Y = f(\mathbf{x})$ that maps the values of a set of predictors into a target variable value (can be a classification or a regression problem)
- A (training) data set $\{ \langle \mathbf{x}_i, y_i \rangle \}_{i=1}^N$, with known values of this mapping
- Performance evaluation criterion(a) - metric(s) of predictive performance (e.g. error rate or mean squared error)
- **How to obtain a reliable estimates of the predictive performance** of any solutions we consider to solve the task using the available data set?



Reliability of Estimates

Resubstitution estimates

- Given that we have a data set one possible way to obtain an estimate of the performance of a model is to evaluate it on this data set
- This leads to what is known as a **resubstitution estimate** of the prediction error
- **These estimates are unreliable and should not be used as they tend to be over-optimistic!**



Reliability of Estimates

Resubstitution estimates (2)

- Why are they unreliable?
 - Models are obtained with the goal of optimizing the selected prediction error statistic on the given data set
 - In this context it is expected that they get good scores!
 - The given data set is just a **sample of the unknown distribution** of the problem being tackled
 - What we would like is to have the performance of the model on this distribution
 - As this is usually impossible the best we can do is to evaluate the model on **new samples** of this distribution



Goal of Performance Estimation

Main Goal of Performance Estimation

Obtain a **reliable estimate** of the **expected prediction error** of a model on the unknown data distribution

- In order to be reliable it should be based on evaluation on unseen cases - *a test set*



Goal of Performance Estimation (2)

- Ideally we want to repeat the testing several times
- This way we can collect a series of scores and provide as our estimate the average of these scores, together with the standard error of this estimate
- In summary:
 - calculate the sample mean prediction error on the repetitions as an estimate of the true population mean prediction error
 - complement this sample mean with the standard error of this estimate



Goal of Performance Estimation (3)

- The golden rule of Performance Estimation:
The data used for evaluating (or comparing) any models cannot be seen during model development.



Goal of Performance Estimation (4)

- An experimental methodology should:
 - Allow obtaining several prediction error scores of a model, E_1, E_2, \dots, E_k
 - Such that we can calculate a sample mean prediction error

$$\bar{E} = \frac{1}{k} \sum_{i=1}^k E_i$$

- And also the respective standard error of this estimate

$$SE(\bar{E}) = \frac{s_E}{\sqrt{k}}$$

where s_E is the sample standard deviation of E measured as

$$\sqrt{\frac{1}{k-1} \sum_{i=1}^k (E_i - \bar{E})^2}$$


The Holdout Method

The Holdout Method and Random Subsampling

- The holdout method consists on **randomly dividing the available data sample in two sub-sets** - one used for training the model; and the other for testing/evaluating it
 - A frequently used proportion is 70% for training and 30% for testing



The Holdout Method (2)

- If we have a small data sample there is the danger of either having a too small test set (unreliable estimates as a consequence), or removing too much data from the training set (worse model than what could be obtained with the available data)
- We only get one prediction error score - no average score nor standard error
- If we have a very large data sample this is actually the preferred evaluation method

Random Subsampling

- The Random Subsampling method is a variation of holdout method and it simply consists of repeating the holdout process several times by randomly selecting the train and test partitions
- Has the same problems as the holdout with the exception that we already get several scores and thus can calculate means and standard errors
- If the available data sample is too large the repetitions may be too demanding in computation terms

The Holdout method in R

```

library(DMwR2)
set.seed(1234)
data(Boston, package='MASS')
## random selection of the holdout
trPerc <- 0.7
sp <- sample(1:nrow(Boston), as.integer(trPerc*nrow(Boston)))
## division in two samples
tr <- Boston[sp,]
ts <- Boston[-sp,]
## obtaining the model and respective predictions on the test set
m <- rpartXse(medv ~., tr)
p <- predict(m, ts)
## evaluation
mean((ts$medv-p)^2)

## [1] 22.1313

```



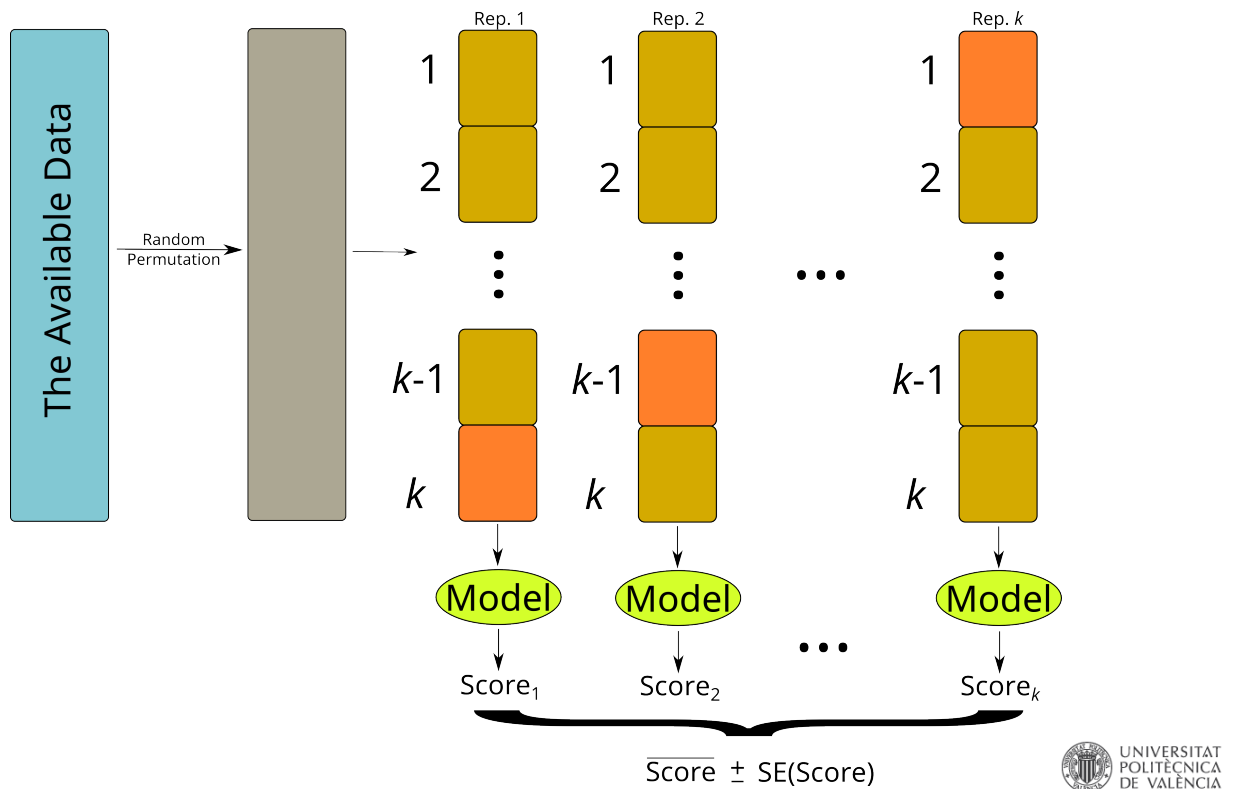
Cross Validation

The k-fold Cross Validation Method

- The idea of k-fold Cross Validation (CV) is similar to random subsampling
- It essentially consists of k repetitions of training on part of the data and then test on the remaining
- The difference lies on the way the partitions are obtained



The k-fold Cross Validation Method (cont.)



Leave One Out Cross Validation Method (LOOCV)

- Similar idea to k-fold Cross Validation (CV) but in this case on each iteration a single case is left out of the training set
- This means it is essentially equivalent to n -fold CV, where n is the size of the available data set

The Bootstrap Method

- Train a model on a random sample of size n with replacement from the original data set (of size n)
 - Sampling with replacement means that after a case is randomly drawn from the data set, it is “put back on the sampling bag”
 - This means that several cases will appear more than once on the training data
 - On average only 63.2% of all cases will be on the training set
- Test the model on the cases that were not used on the training set
- Repeat this process many times (typically around 200)
- The average of the scores on these repetitions is known as the ϵ_0 bootstrap estimate
- The .632 bootstrap estimate is obtained by $.368 \times \epsilon_r + .632 \times \epsilon_0$, where ϵ_r is the resubstitution estimate



Bootstrap

Bootstrap in R

```

data(Boston, package='MASS')
nreps <- 200
scores <- vector("numeric", length=nreps)
n <- nrow(Boston)
set.seed(1234)
for(i in 1:nreps) {
  # random sample with replacement
  sp <- sample(n, n, replace=TRUE)
  # data splitting
  tr <- Boston[sp,]
  ts <- Boston[-sp,]
  # model learning and prediction
  m <- lm(medv ~., tr)
  p <- predict(m, ts)
  # evaluation
  scores[i] <- mean((ts$medv-p)^2)
}
# calculating means and standard errors
summary(scores)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 16.37   21.70   24.20   24.56   26.47   48.82

```



The Infra-Structure of package performanceEstimation

- The package **performanceEstimation** provides a set of functions that can be used to carry out comparative experiments of different models on different predictive tasks
- This infra-structure can be applied to any model/task/evaluation metric
- Installation:
 - Official release (from CRAN repositories):

```
install.packages("performanceEstimation")
```

- Development release (from Github):

```
library(devtools) # You need to install this package before!
install_github("ltorgo/performanceEstimation", ref="develop")
```



The Infra-Structure of package performanceEstimation

- The main function of the package is `performanceEstimation()`
 - It has 3 arguments:
 - 1 The predictive tasks to use in the comparison
 - 2 The models to be compared
 - 3 The estimation task to be carried out
- The function implements a wide range of experimental methodologies including all we have discussed



A Simple Example

- Suppose we want to estimate the mean squared error of regression trees in a certain regression task using cross validation

```
library(performanceEstimation)
library(DMwR2)
data(Boston, package='MASS')
res <- performanceEstimation(
  PredTask(medv ~ ., Boston),
  Workflow("standardWF", learner="rpartXse"),
  EstimationTask(metrics="mse", method=CV(nReps=1, nFolds=10)))
```



A Simple Example (2)

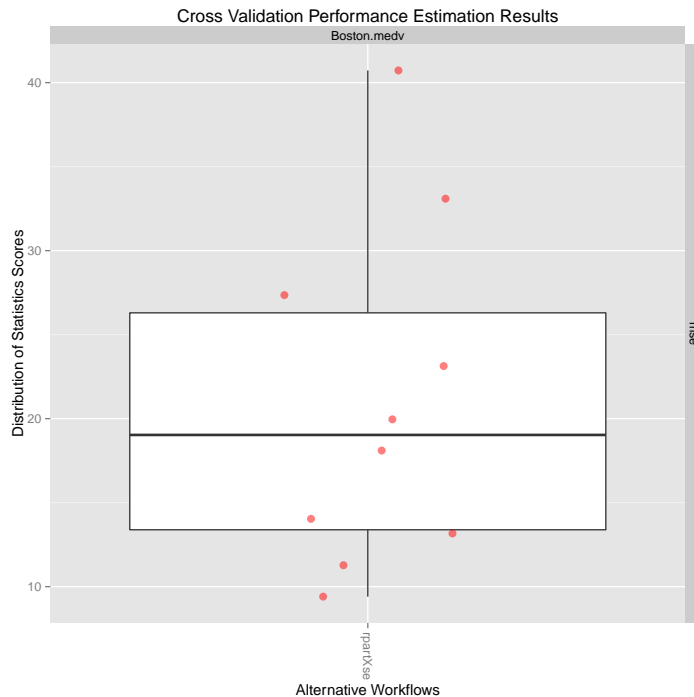
```
summary(res)

##
## == Summary of a Cross Validation Performance Estimation Experiment ==
##
## Task for estimating mse using
## 1 x 10 - Fold Cross Validation
## Run with seed = 1234
##
## * Predictive Tasks :: Boston.medv
## * Workflows :: rpartXse
##
## -> Task: Boston.medv
## *Workflow: rpartXse
##      mse
## avg      21.02393
## std      10.15683
## med      19.02955
## iqr      12.91203
## min       9.40574
## max      40.72403
## invalid  0.00000
```



A Simple Example (3)

```
plot(res)
```



Predictive Tasks

- Objects of class **PredTask** describing a predictive task
 - Classification
 - Regression
 - Time series forecasting
- Created with the constructor with the same name

```
data(iris)
PredTask(Species ~ ., iris)

## Prediction Task Object:
## Task Name      :: iris.Species
## Task Type      :: classification
## Target Feature  :: Species
## Formula        :: Species ~ .
## Task Data Source :: iris

PredTask(Species ~ ., iris, "IrisDS", copy=TRUE)

## Prediction Task Object:
## Task Name      :: IrisDS
## Task Type      :: classification
## Target Feature  :: Species
## Formula        :: Species ~ .
## Task Data Source :: internal 150x5 data frame.
```



Workflows

- Objects of class **Workflow** describing an approach to a predictive task
 - Standard Workflows
 - Function `standardWF` for classification and regression
 - Function `timeseriesWF` for time series forecasting
 - User-defined Workflows



Standard Workflows for Classification and Regression Tasks

```
library(e1071)
Workflow("standardWF", learner="svm", learner.pars=list(cost=10, gamma=0.1))

## Workflow Object:
## Workflow ID      :: svm
## Workflow Function :: standardWF
## Parameter values:
## learner -> svm
## learner.pars -> cost=10 gamma=0.1
```

“standardWF” can be omitted ...

```
Workflow(learner="svm", learner.pars=list(cost=5))

## Workflow Object:
## Workflow ID      :: svm
## Workflow Function :: standardWF
## Parameter values:
## learner -> svm
## learner.pars -> cost=5
```

Standard Workflows for Classification and Regression Tasks (cont.)

- Main parameters of the constructor:
 - Learning stage
 - `learner` - which function is used to obtain the model for the training data
 - `learner.pars` - list with the parameter settings to pass to the learner
 - Prediction stage
 - `predictor` - function used to obtain the predictions (defaults to `predict()`)
 - `predictor.pars` - list with the parameter settings to pass to the predictor



Standard Workflows for Classification and Regression Tasks (cont.)

- Main parameters of the constructor (cont.):
 - Data pre-processing
 - `pre` - vector with function names to be applied to the training and test sets before learning
 - `pre.pars` - list with the parameter settings to pass to the functions
 - Predictions post-processing
 - `post` - vector with function names to be applied to the predictions
 - `post.pars` - list with the parameter settings to pass to the functions



Standard Workflows for Classification and Regression Tasks (cont.)

```
data(algae, package="DMwR2")
res <- performanceEstimation(
  PredTask(a1 ~ ., algae[, 1:12], "A1"),
  Workflow(learner="lm", pre="centralImp", post="onlyPos"),
  EstimationTask("mse", method=CV()) # defaults to 1x10-fold CV
)

##
##
## ##### PERFORMANCE ESTIMATION USING CROSS VALIDATION #####
##
## ** PREDICTIVE TASK :: A1
##
## ++ MODEL/WORKFLOW :: lm
## Task for estimating mse using
## 1 x 10 - Fold Cross Validation
## Run with seed = 1234
## Iteration :*****
```



Evaluating Variants of Workflows

Function `workflowVariants()`

Sometimes you want to evaluate different parameter variants of the same workflow - that is the goal of function `workflowVariants()`. It produces a vector of **Workflow** objects without having to specify all of them.

```
library(e1071)
data(Boston, package="MASS")
res2 <- performanceEstimation(
  PredTask(medv ~ ., Boston),
  workflowVariants(learner="svm",
    learner.pars=list(cost=1:5, gamma=c(0.1, 0.01))),
  EstimationTask(metrics="mse", method=CV())
```



Evaluating Variants of Workflows (cont.)

```
summary(res2)

##
## == Summary of a Cross Validation Performance Estimation Experiment ==
##
## Task for estimating mse using
## 1 x 10 - Fold Cross Validation
## Run with seed = 1234
##
## * Predictive Tasks :: Boston.medv
## * Workflows :: svm.v1, svm.v2, svm.v3, svm.v4, svm.v5, svm.v6, svm.v7, svm.v8, svm.v9, svm.v10
##
## -> Task: Boston.medv
## *Workflow: svm.v1
##      mse
## avg    14.80685
## std    10.15295
## med    12.27015
## iqr    11.87737
## min     5.35198
## max    38.39681
## invalid 0.00000
##
## *Workflow: svm.v2
##      mse
## avg    11.995178
## std     7.908371
## med     8.359433
## iqr    11.626306
## min     4.842848
```

AT
CA
IA

```
## invalid 0.00000
##
## *Workflow: svm.v3
```

Exploring the Results

```
getWorkflow("svm.v1", res2)

## Workflow Object:
## Workflow ID      :: svm.v1
## Workflow Function :: standardWF
## Parameter values:
## learner.pars -> cost=1 gamma=0.1
## learner -> svm

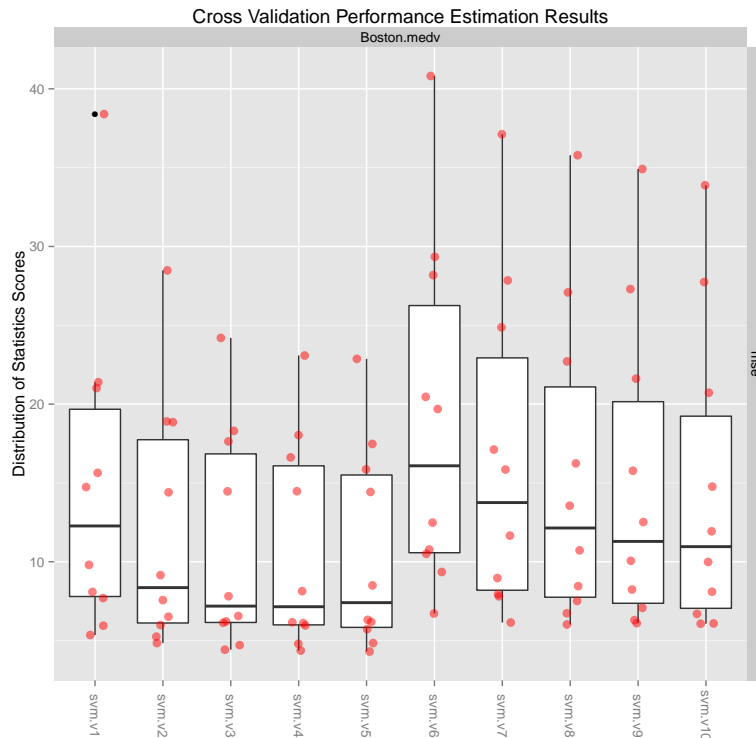
topPerformers(res2)

## $Boston.medv
## Workflow Estimate
## mse svm.v5 10.65
```

```
## *Workflow: svm.v/
##      mse
```

Visualizing the Results

```
plot(res2)
```



Estimation Tasks

- Objects of class **EstimationTask** describing the estimation task
 - Main parameters of the constructor
 - **metrics** - vector with names of performance metrics
 - **method** - object of class **EstimationMethod** describing the method used to obtain the estimates

```
EstimationTask(metrics=c("F", "rec", "prec"), method=Bootstrap(nReps=100))

## Task for estimating F,rec,prec using
## 100 repetitions of e0 Bootstrap experiment
## Run with seed = 1234
```



Performance Metrics

- Many classification and regression metrics are available
 - Check the help page of functions `classificationMetrics` and `regressionMetrics`
- User can provide a function that implements any other metric she/he wishes to use
 - Parameters **evaluator** and **evaluator.pars** of the `EstimationTask` constructor



Comparing Different Algorithms on the Same Task

```
library(randomForest)
library(e1071)
res3 <- performanceEstimation(
  PredTask(medv ~ ., Boston),
  workflowVariants("standardWF",
    learner=c("rpartXse", "svm", "randomForest")),
  EstimationTask(metrics="mse", method=CV(nReps=2, nFolds=5)))
```



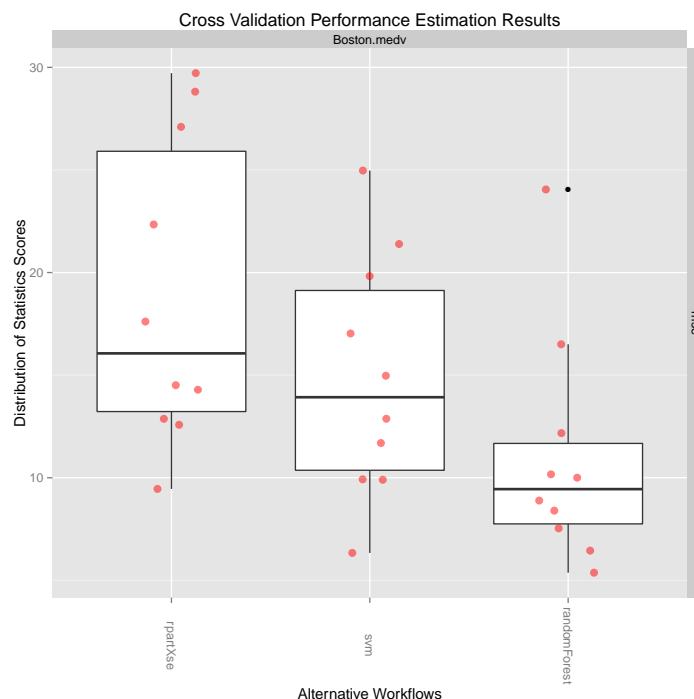
Some auxiliary functions

```
rankWorkflows(res3, 3)
```

```
## $Boston.medv
## $Boston.medv$mse
##      Workflow Estimate
## 1 randomForest 10.95412
## 2          svm 14.89183
## 3      rpartXse 18.92990
```

The Results

```
plot(res3)
```

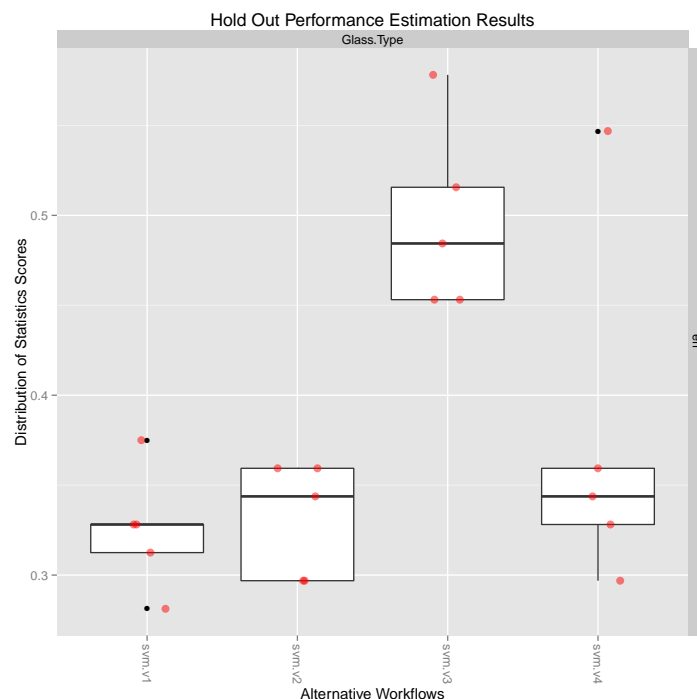


An example using Holdout and a classification task

```
data(Glass, package='mlbench')
res4 <- performanceEstimation(
  PredTask(Type ~ ., Glass),
  workflowVariants(learner="svm", # You may omit "standardWF" !
                  learner.pars=list(cost=c(1, 10),
                                    gamma=c(0.1, 0.01))),
  EstimationTask(metrics="err", method=Holdout(nReps=5, hldSz=0.3)))
```

The Results

```
plot(res4)
```



An example involving more than one task

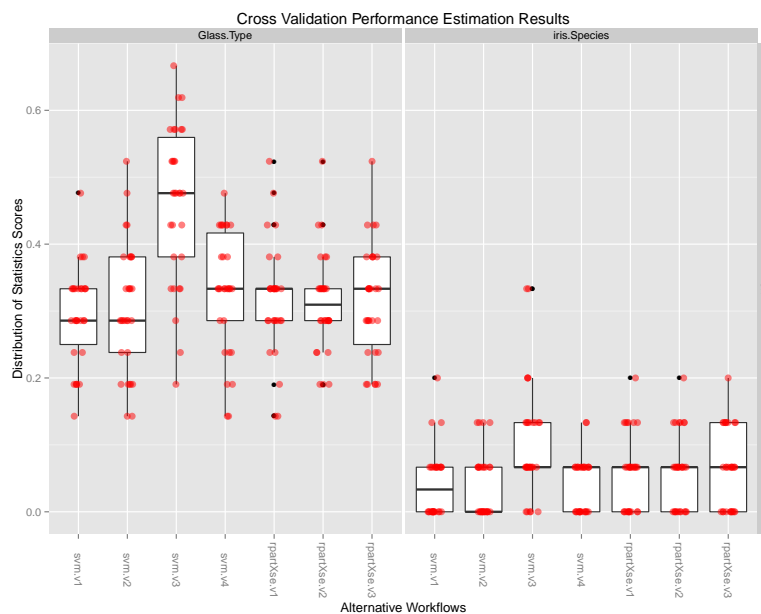
```

data(Glass, package='mlbench')
data(iris)
res5 <- performanceEstimation(
  c(PredTask(Type ~ ., Glass), PredTask(Species ~ ., iris)),
  c(workflowVariants(learner="svm",
                    learner.pars=list(cost=c(1,10),
                                       gamma=c(0.1,0.01))),
    workflowVariants(learner="rpartXse",
                    learner.pars=list(se=c(0,0.5,1)),
                    predictor.pars=list(type="class"))),
  EstimationTask(metrics="err", method=CV(nReps=3))

```

The Results

```
plot(res5)
```



The Results (2)

```
topPerformers(res5)
```

```
## $Glass.Type
##   Workflow Estimate
## err  svm.v1      0.294
##
## $iris.Species
##   Workflow Estimate
## err  svm.v2      0.04
```

```
topPerformer(res5, "err", "Glass.Type")
```

```
## Workflow Object:
## Workflow ID      :: svm.v1
## Workflow Function :: standardWF
##   Parameter values:
## learner.pars  -> cost=1 gamma=0.1
## learner      -> svm
```



Hands on Performance Estimation

the Algae data set

Load in the data set `algae` and answer the following questions:

- 1 Estimate the MSE of a regression tree for forecasting alga *a1* using 10-fold Cross validation.
- 2 Repeat the previous exercise this time trying some variants of random forests. Check what are the characteristics of the best performing variant.
- 3 Compare the results in terms of mean absolute error of the default variants of a regression tree, a linear regression model and a random forest, in the task of predicting alga *a3*. Use 2 repetitions of a 5-fold Cross Validation experiment.
- 4 Carry out an experiment designed to select what are the best models for each of the seven harmful algae. Use 10-fold Cross Validation. For illustrative purposes consider only the default variants of regression trees, linear regression and random forests.



Are the Observed Differences Statistically Significant?

Statistical Hypothesis Testing

- Test if some result is unlikely to have occurred by chance
- The null hypothesis: there is no difference among a set of alternative workflows
- This hypothesis is rejected if the result of the test has a p-value less than a certain selected threshold (typically 0.01 or 0.05, i.e. 99% or 95% confidence)

- There are many statistical tests that could be used
- The work by Demsar (2006) includes what are the current recommendations for different experimental setups

J. Demsar. Statistical comparisons of classifiers over multiple data sets. Journal of Machine Learning Research, 7:1–30, 2006.



Paired Comparisons on a Task

Wilcoxon Signed Rank Test

- The null hypothesis: the difference between the two workflows is zero
- This hypothesis is rejected if the result of the test has a p-value less than a certain selected threshold (typically 0.01 or 0.05, i.e. 99% or 95% confidence)



A Simple Example

```
library(performanceEstimation)
library(DMwR2) # because of rpartXse
data(Boston, package="MASS")
res <- performanceEstimation(
  PredTask(medv ~ ., Boston),
  workflowVariants(learner="rpartXse", learner.pars=list(se=c(0, 0.5, 1))),
  EstimationTask(metrics="mse", method=CV(nReps=3, nFolds=10))
)
```



A Simple Example

```
pres <- pairedComparisons(res)

## Warning in pairedComparisons(res): With less 2 tasks the Friedman,
## Nemenyi and Bonferroni-Dunn tests are not calculated.

pres$mse$WilcoxonSignedRank.test

## , , Boston.medv
##
##           MedScore DiffMedScores    p.value
## rpartXse.v1 18.18101             NA         NA
## rpartXse.v2 19.54956    -1.36855309 0.5837571
## rpartXse.v3 18.21299    -0.03198033 0.5027610
```



Which ones are significant at some level?

```

signifDiffs (pres,p.limit=0.05)

## $mse
## $mse$WilcoxonSignedRank.test
## $mse$WilcoxonSignedRank.test$Boston.medv
##      MedScore DiffMedScores      p.value
##      18.18101          NA          NA
##
##
## $mse$t.test
## $mse$t.test$Boston.medv
##      AvgScore DiffAvgScores      p.value
##      19.65952          NA          NA

```

Paired Comparisons on a Multiple Tasks

Demsar (2006) recommended procedure

- Step 1: Friedman test
 - Null hypothesis: all workflows are equivalent and so their rankings across the tasks are equal
- If this hypothesis is rejected then we can move to the second step
 - Paired comparisons among all pairs of workflows
 - Nemenyi post-hoc test
 - Null hypothesis: there is no significant difference among the ranks of a certain pair of workflows
 - Paired comparisons against a baseline
 - Bonferroni-Dunn post-hoc test
 - Null hypothesis: there is no significant difference among the ranks of a certain workflow and the baseline

An Example with Several Tasks

```

library(performanceEstimation)
library(e1071)
library(randomForest)
tgts <- 12:18
tasks <- c()
for(t in tgts)
  tasks <- c(tasks,
             PredTask(as.formula(paste(colnames(algae)[t], '~ .')),
                      algae[, c(1:11, t)],
                      paste0("algaA", t-11),
                      copy=TRUE))
res.algae <- performanceEstimation(
  tasks,
  workflowVariants(learner=c("svm", "lm", "randomForest"),
                   pre="knnImp"),
  EstimationTask("mae", method=CV())
)

```



An Example with Several Tasks (cont.)

Can we reject the hypothesis that the workflows have the same ranking across all tasks?

```

pres <- pairedComparisons(res.algae)
pres$mae$F.test

## $chi
## [1] 12.28571
##
## $FF
## [1] 43
##
## $critVal
## [1] 0.3574087
##
## $rejNull
## [1] TRUE

```



An Example with Several Tasks (cont.)

Are there any significant differences among the workflows?

```
pres$mae$Nemenyi.test

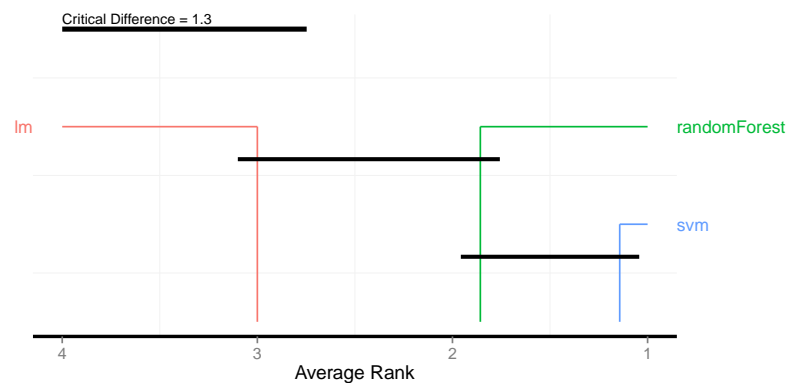
## $critDif
## [1] 1.252761
##
## $rkDifs
##           svm      lm randomForest
## svm      0.0000000 1.857143    0.7142857
## lm       1.8571429 0.000000    1.1428571
## randomForest 0.7142857 1.142857    0.0000000
##
## $signifDifs
##           svm      lm randomForest
## svm      FALSE  TRUE     FALSE
## lm       TRUE  FALSE   FALSE
## randomForest FALSE FALSE   FALSE
```



CD diagrams for the Nemenyi test

Average rank differences that are not statistically significant are connected

```
CDdiagram.Nemenyi(pres)
```



An Example with Several Tasks (cont.)

Suppose “lm” was our baseline system and we wanted to check if the other alternatives were able to improve over it on these tasks

```
pres <- pairedComparisons(res.algae,baseline="lm")
pres$mae$BonferroniDunn.test

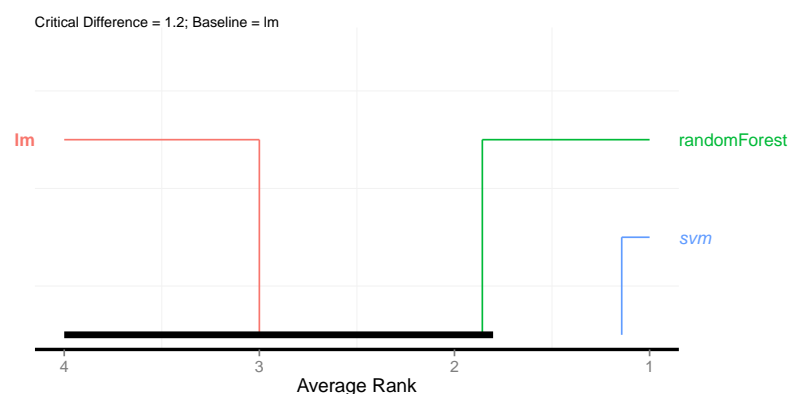
## $critDif
## [1] 1.19808
##
## $baseline
## [1] "lm"
##
## $rkDifs
##          svm randomForest
## 1.857143  1.142857
##
## $signifDifs
##          svm randomForest
##          TRUE          FALSE
```



CD diagrams for the Bonferroni Dunn test

Are the average ranks of the other systems significantly better than the one of “lm”?

```
CDdiagram.BD(pres)
```



Hands on Statistical Significance

Using the *algae* data set from package **DMwR2** answer the following questions

- 1 For the 7 different algae, choose a reasonable set of SVM variants and estimate their MSE error.
- 2 Check if these alternatives are significantly better than the SVM with the default parameter settings
- 3 Present the results of the previous question visually