

Machine Learning com Linguagem R

Deivison Venicio Souza

Universidade Federal Pará - UFPA
Engenheiro Florestal, Me. Ciências Florestais
Programa de Pós-graduação em Engenharia Florestal - UFPR
Especialização Data Science & Big Data (Em andamento) - UFPR
(deivisonvs@ufpa.br)

07/12/2018
Curitiba, PR

Sumário

1 Diferença: AI, ML e DL

2 Machine Learning

- Tipos de Aprendizado
- ML nas Ciências Florestais

3 Machine Learning no R

- A linguagem R
- Pacotes para ML

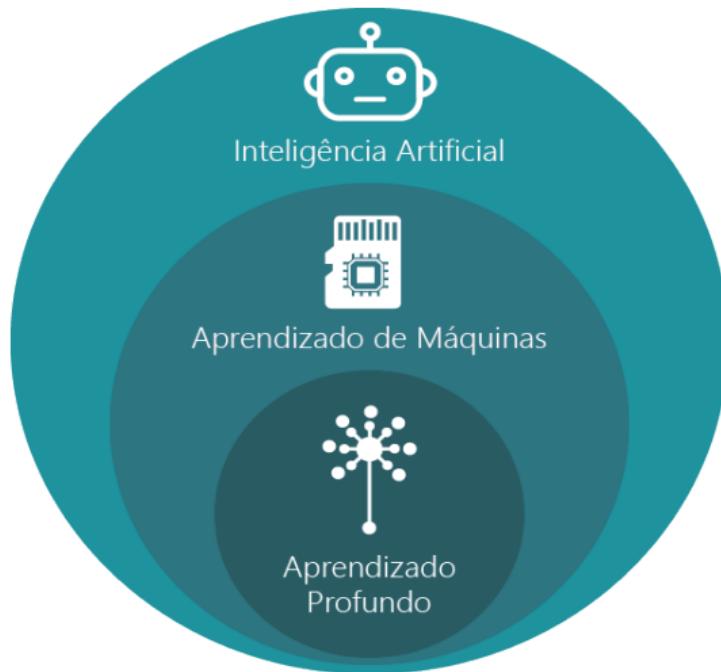
- Popularidade de pacotes

4 Pacote Caret

5 ML supervisionado

- Pré-processamento
- Divisão dos dados
- Treinamento e Ajuste de Hiperparâmetros

Diferença: AI, ML e DL?



AI - Definição

O termo *Artificial Intelligence* foi cunhado pela primeira vez em 1956.

Definição - John McCarthy

A ciência e engenharia de fazer máquinas inteligentes.

AI - Definição

Outra definição

“Programas (agentes inteligentes) com a capacidade de **aprender e raciocinar como seres humanos**.”

“O estudo e projeto de agentes inteligentes.”

“Um agente inteligente é um sistema que percebe seu ambiente e toma atitudes que maximizam suas chances de sucesso.” ([Wikipedia](#))

A AI está no cotidiano!

① Reconhecimento Facial (*Face Recognition*)

- Como o [Facebook](#) é capaz de reconhecer uma pessoa em uma foto?

② Sistemas de Recomendação (*Recommendation Systems*)

- Como o [Google Maps](#) é capaz de traçar a melhor rota para um destino?
- Como a [NetFlix](#) faz para recomendar séries e filmes (e outros) para os usuários?
([SR NetFlix](#))

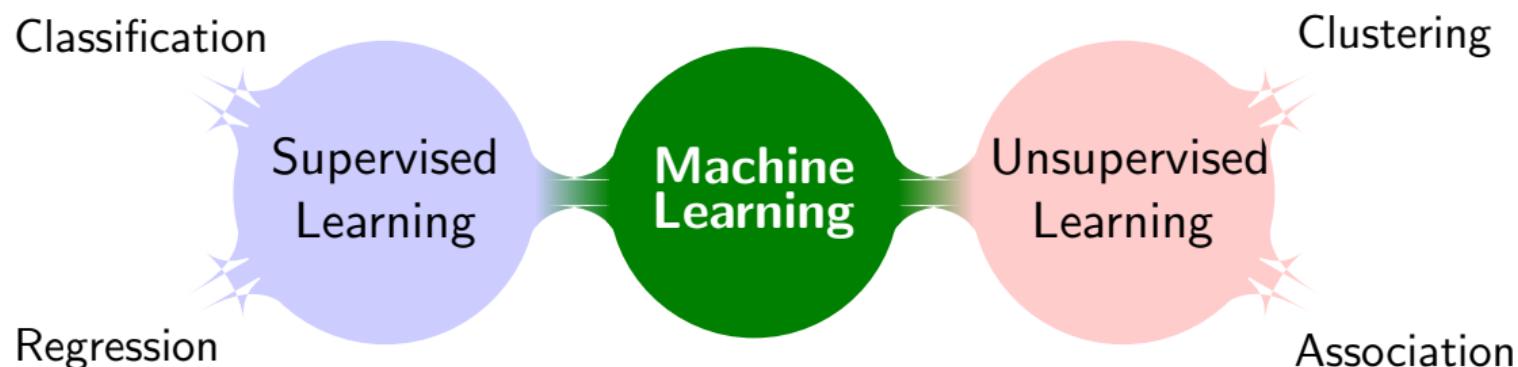
Machine Learning

Arthur Lee Samuel: Termo “Machine Learning” em 1959.

Definição - Parafraseando Samuel (1959)

É um **subconjunto de inteligência artificial** que frequentemente usa técnicas estatísticas para dar aos computadores a capacidade de “aprender” com **dados**, sem serem explicitamente programados.

ML - Tipos de Aprendizado



Supervised Learning

Classificação

Uma tarefa de classificação é quando a variável resposta é uma categoria.
Ex.: Reconhecimento de espécies florestais a partir de imagens.

Regressão

Um tarefa de regressão é quando a variável resposta é um valor real (numérica).
Ex.: Biomassa, Volume, Altura de árvores, entre outros.

SL - Uma rápida intuição

Iris Data Set (Fisher, 1936) - UCI Machine Learning Repository

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
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
...				
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.9	3.1	4.9	1.5	versicolor
5.5	2.3	4.0	1.3	versicolor
6.5	2.8	4.6	1.5	versicolor
...				
6.3	3.3	6.0	2.5	virginica
5.8	2.7	5.1	1.9	virginica
7.1	3.0	5.9	2.1	virginica
6.3	2.9	5.6	1.8	virginica
6.5	3.0	5.8	2.2	virginica



Iris Setosa



Iris Virginica



Iris Versicolor

SL - Uma rápida intuição

Meta: Modelar uma árvore de decisão e resolver o problema de reconhecimento das espécies de Íris.

Possibilidade: Árvore de Classificação e Regressão (CART) (BREIMAN et al., 1984)

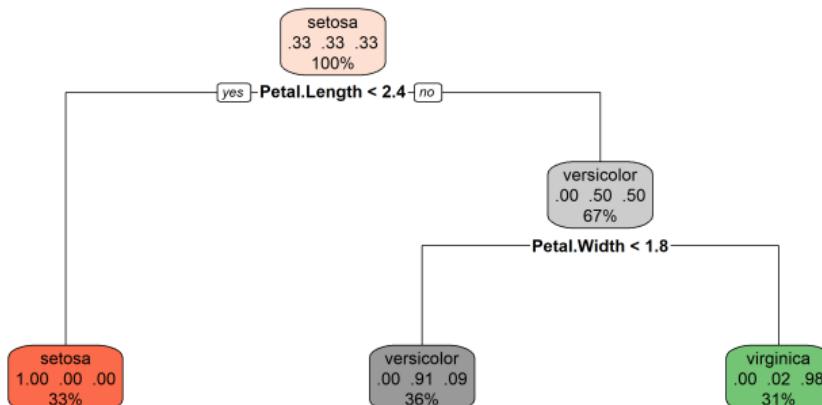
Linguagem R: Package **Rpart** (Recursive Partitioning and Regression Trees)
(THERNEAU et al., 2017)

Qual a ideia do CART?

“A ideia é dividir o espaço de covariável em várias partições e ajustar um modelo constante da variável resposta em cada partição.”

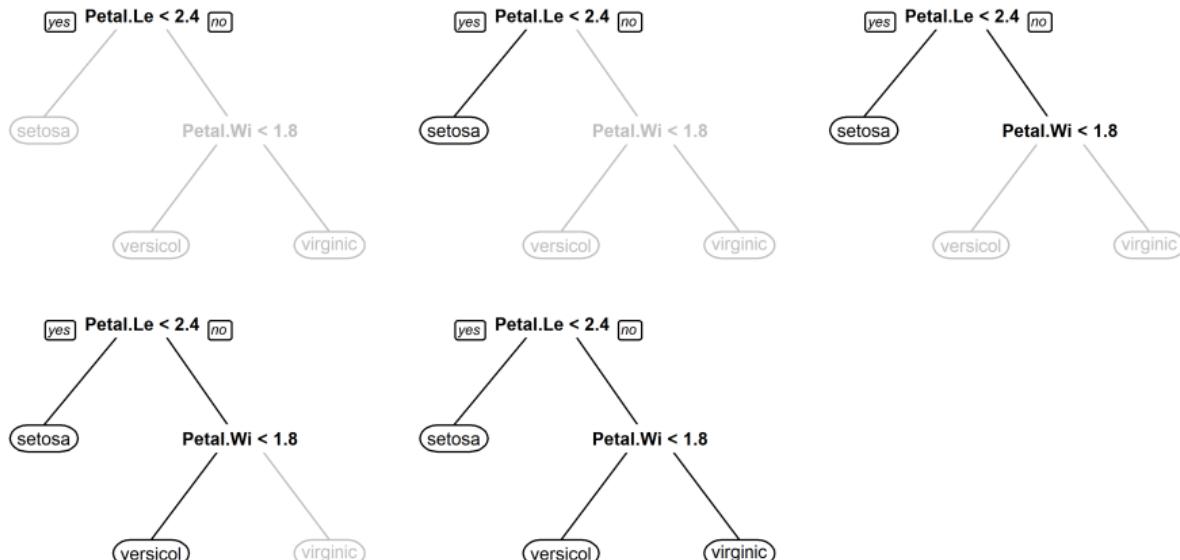
SL - Uma rápida intuição

Árvore de Classificação



SL - Uma rápida intuição

Regras da Árvore de Classificação



Aplicações de ML nas Ciências Florestais

ML nas Ciências Florestais - Regressão

Sanquetta et al. Carbon Balance and Management 2013, 8:6
<http://www.cbmjournal.com/content/8/1/6>



METHODOLOGY

Open Access

On the use of data mining for estimating carbon storage in the trees

Carlos Roberto Sanquetta^{1†}, Jaime Wojciechowski^{2†}, Ana Paula Dalla Corte^{1*†}, Aurélio Lourenço Rodrigues^{3†} and Greyce Charllyne Benedet Maas^{3†}

Sanquetta et al. (2013)

ML nas Ciências Florestais - Regressão

Artificial Intelligence Models to Estimate Biomass of Tropical Forest Trees

Razer Anthom Nizer Rojas Montaño, Carlos Roberto Sanquetta, Jaime Wojciechowski, Eduardo Mattar, Ana Paula Dalla Corte, and Eduardo Todt

Montaño et al. (2018)

ML nas Ciências Florestais - Classificação

Machine Vision and Applications (2015) 26:279–293
DOI 10.1007/s00138-015-0659-0

ORIGINAL PAPER

Forest species recognition based on dynamic classifier selection and dissimilarity feature vector representation

J. G. Martins · L. S. Oliveira · A. S. Britto Jr. ·
R. Sabourin

Martins et al. (2015)

ML nas Ciências Florestais - Classificação

[Ecological Informatics 46 \(2018\) 1–7](#)



Contents lists available at [ScienceDirect](#)

Ecological Informatics

journal homepage: www.elsevier.com/locate/ecolinf



Automatic classification of native wood charcoal

T.M. Maruyama^a, L.S. Oliveira^a, A.S. Britto Jr^{b,c,*}, S. Nisgoski^a



Maruyama et al. (2018)

Machine Learning e Linguagem R

A Linguagem R

R é uma linguagem e ambiente para computação estatística e gráficos.

- **Sistemas Operacionais:** Compila e roda em uma ampla variedade de plataformas (Linux, Windows e MacOS);
 - **Software Livre:** R está disponível como Software Livre, em forma de código fonte;
 - **CRAN:** O repositório oficial de pacotes do R. ([CRAN Link](#))



CRAN Task View

Atualmente, são 102 pacotes sobre ML publicados no (CRAN Task View).

CRAN Task View: Machine Learning & Statistical Learning
 (Mantenedor: Torsten Hothorn, Versão: 21/07/2018)

ahaz	e1071	GMMBoost	LiblineaR	partykit	relaxo	spa
arules	earth	gradDescent	LogicReg	pdp	rgenoud	stabs
BART	effects	grf	LTRCTrees	penalized	RLT	SuperLearner
bartMachine	elasticnet	grplasso	maptree	penalizedLDA	Rmalschains	svmpath
BayesTree	ElemStatLearn	grpreg	mboost	picasso	rminer	tensorflow
biglasso	evclass	h2o	mlr	plotmo	rnn	tgp
bmrmm	evtree	hda	model4you	quantregForest	ROCR	tree
Boruta	FCNN4R	hdi	MXM	randomForest	RoughSets	trtf
bst	frbs	hdm	ncvreg	randomForestSRC	rpart	varSelRF
C50	GAMBoost	ICEbox	nnet	ranger	RPMM	vcrpart
caret	gamboostLSS	ipred	oem	rattle	RSNNS	wsrf
CORElearn	gbm	kernlab	OneR	Rborist	RWeka	xgboost
CoxBoost	ggRandomForests	klaR	opusminer	RcppDL	RXshrink	-
Cubist	glmnet	lars	pamr	rdetools	sda	-
deepnet	glmpath	lasso2	party	REEMtree	SIS	-



Qual pacote usar?



Qual pacote usar?

- ① **Objetivo:** Qual(is) algoritmo(s) será(ão) testado(s)?;
- ② **Popularidade:** Pode-se medir a popularidade dos pacotes de AM;
- ③ **Número de downloads:** Observar o número de downloads do pacote no repositório CRAN (`pacote cranlogs → cran_downloads()`); e
- ④ **Tempo de execução:** Maior velocidade de execução da tarefa → **Big Data**.

Qual pacote usar?

- ① **Objetivo:** Qual(is) algoritmo(s) será(ão) testado(s)?;
- ② **Popularidade:** Pode-se medir a popularidade dos pacotes de AM;
- ③ **Número de downloads:** Observar o número de downloads do pacote no repositório CRAN (`pacote cranlogs → cran_downloads()`); e
- ④ **Tempo de execução:** Maior velocidade de execução da tarefa → **Big Data**.

Qual pacote usar?

- ① **Objetivo:** Qual(is) algoritmo(s) será(ão) testado(s)?;
- ② **Popularidade:** Pode-se medir a popularidade dos pacotes de AM;
- ③ **Número de downloads:** Observar o número de downloads do pacote no repositório CRAN (`pacote cranlogs → cran_downloads()`); e
- ④ **Tempo de execução:** Maior velocidade de execução da tarefa → **Big Data**.

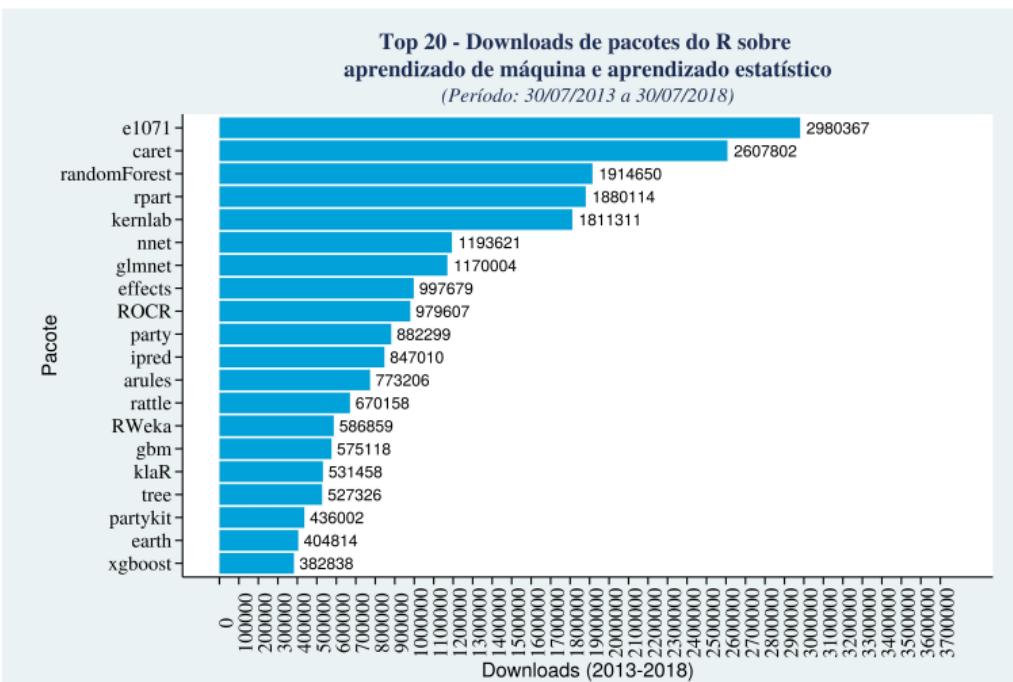
Qual pacote usar?

- ① **Objetivo:** Qual(is) algoritmo(s) será(ão) testado(s)?;
- ② **Popularidade:** Pode-se medir a popularidade dos pacotes de AM;
- ③ **Número de downloads:** Observar o número de downloads do pacote no repositório CRAN ([pacote cranlogs → cran_downloads\(\)](#)); e
- ④ **Tempo de execução:** Maior velocidade de execução da tarefa → [Big Data](#).

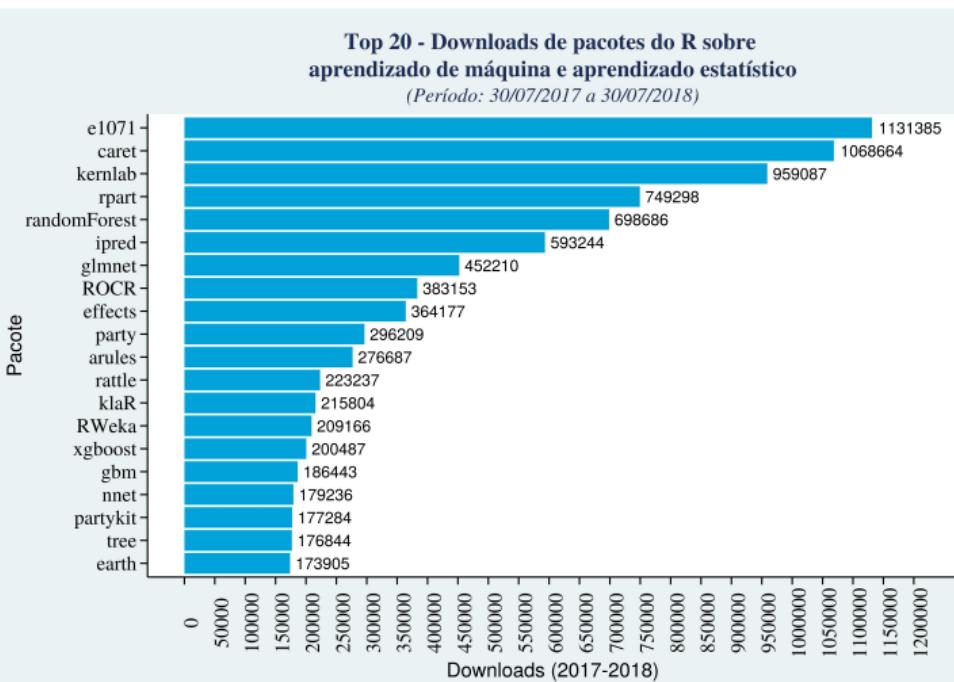
Qual pacote usar?

- ① **Objetivo:** Qual(is) algoritmo(s) será(ão) testado(s)?;
- ② **Popularidade:** Pode-se medir a popularidade dos pacotes de AM;
- ③ **Número de downloads:** Observar o número de downloads do pacote no repositório CRAN ([pacote cranlogs → cran_downloads\(\)](#)); e
- ④ **Tempo de execução:** Maior velocidade de execução da tarefa → **Big Data**.

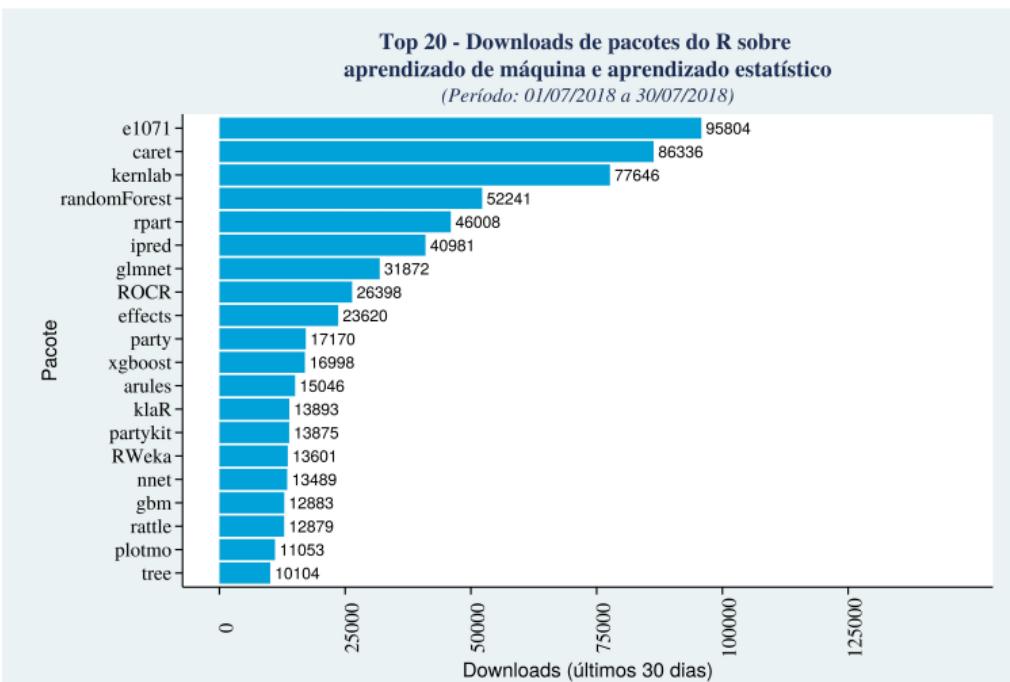
TOP 20 - Downloads (últimos 5 anos)



TOP 20 - Downloads (último ano)



TOP 20 - Downloads (30 dias)



Pacote **caret**

(**C**lassification **A**nd **R**egression **T**raining)

(KUHN; JOHNSON, 2013; KUHN, 2018)

"Constitui um conjunto de funções que tentam simplificar o processo de construção de modelos preditivos."

web page: <https://topepo.github.io/caret/index.html>

Pacote Caret



Pacote Caret

O pacote **caret** surge com o objetivo de: (KUHN, 2008)

- **Sintaxe:** Eliminar diferenças sintáticas entre muitas funções de construção de modelos preditivos;
- **Abordagens:** Desenvolver um conjunto de abordagens semi-automatizadas para otimizar os valores de hiperparâmetros de ajuste;
- **Processamento Paralelo:** Criar um pacote que possa ser facilmente estendido para sistemas de processamento paralelo.

Pacote Caret

O pacote **caret** surge com o objetivo de: (KUHN, 2008)

- **Sintaxe:** Eliminar diferenças sintáticas entre muitas funções de construção de modelos preditivos;
- **Abordagens:** Desenvolver um conjunto de abordagens semi-automatizadas para otimizar os valores de hiperparâmetros de ajuste;
- **Processamento Paralelo:** Criar um pacote que possa ser facilmente estendido para sistemas de processamento paralelo.

Pacote Caret

O pacote **caret** surge com o objetivo de: (KUHN, 2008)

- **Sintaxe:** Eliminar diferenças sintáticas entre muitas funções de construção de modelos preditivos;
- **Abordagens:** Desenvolver um conjunto de abordagens semi-automatizadas para otimizar os valores de hiperparâmetros de ajuste;
- **Processamento Paralelo:** Criar um pacote que possa ser facilmente estendido para sistemas de processamento paralelo.

Pacote Caret

O pacote **caret** surge com o objetivo de: (KUHN, 2008)

- **Sintaxe:** Eliminar diferenças sintáticas entre muitas funções de construção de modelos preditivos;
- **Abordagens:** Desenvolver um conjunto de abordagens semi-automatizadas para otimizar os valores de hiperparâmetros de ajuste;
- **Processamento Paralelo:** Criar um pacote que possa ser facilmente estendido para sistemas de processamento paralelo.

Pacote Caret

Por que usar o pacote **caret**?



Pacote Caret

Em termos gerais, o pacote `caret` possui ferramentas para:



Pacote Caret

Função	Tarefa
findCorrelation()	Encontrar variáveis altamente correlacionadas
nearZeroVar()	Identificar preditores com variância próxima de zero
preProcess()	Realizar pré-processamento
createDataPartition()	Dividir aleatoriamente o conjunto de dados (estratificação)
train()	Ajustar modelos preditivos utilizando reamostragem
varImp()	Estimar a importância das variáveis preditoras
resamples()	Agrupar e visualizar os resultados da reamostragem
diff.resamples()	Fazer inferências sobre diferenças de desempenho de modelos
confusionMatrix()	Criar uma matriz de confusão
plotObsVsPred()	Gerar gráfico de valores observados versus preditos

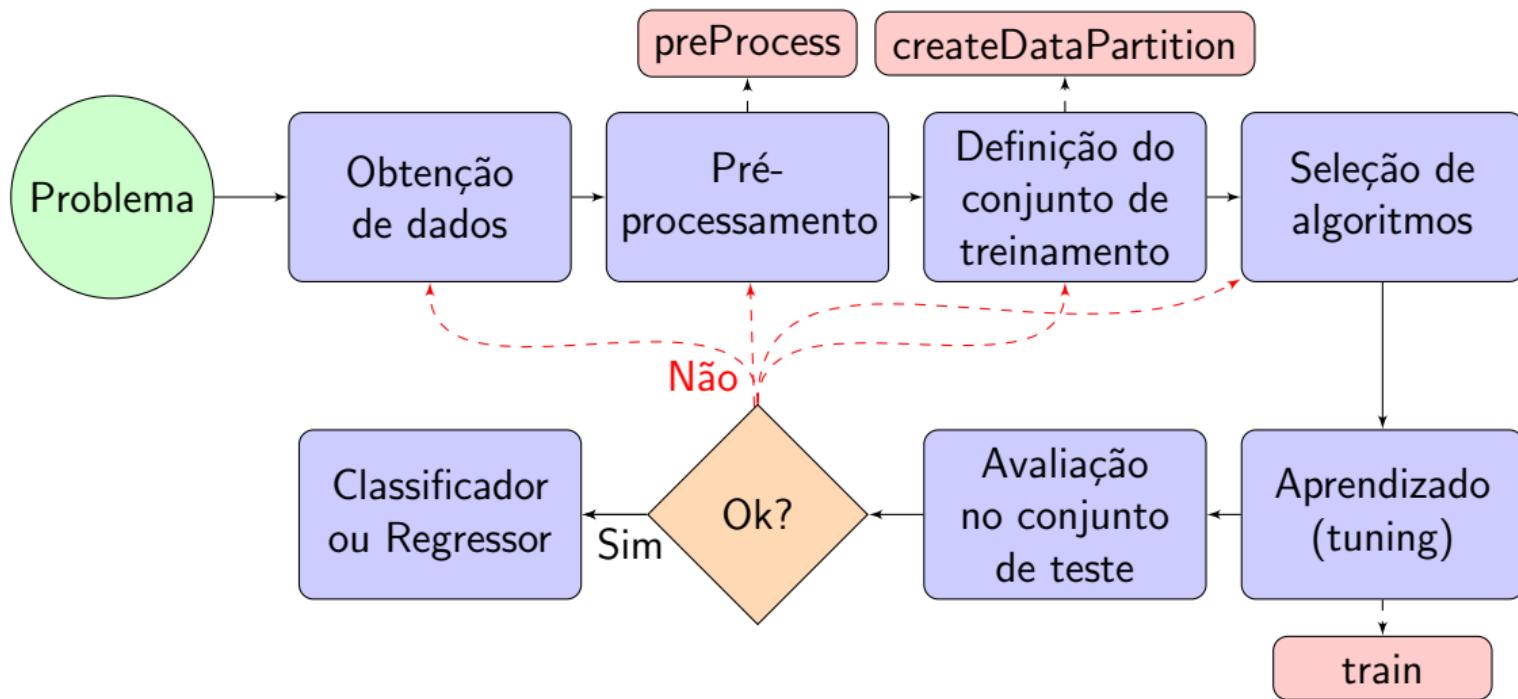
Processo SML

(Supervised Machine Learning)

Uma abordagem usando o pacote [caret](#).

(KUHN; JOHNSON, 2013; KUHN, 2018)

Processo SML - Simplificado



Pré-processamento

Função preProcess()

Pré-processamento

Definição

As técnicas de pré-processamento de dados geralmente se referem à **adição**, **exclusão** ou **transformação** dos dados do conjunto de treinamento (KUHN; JOHNSON, 2013), previamente à construção do modelo preditivo.

Pré-processamento

Por que?

Os diferentes modelos possuem **sensibilidades distintas** para os tipos de preditores. Assim, a forma com que os preditores entram no modelo também é importante (KUHN; JOHNSON, 2013).

A transformação de preditores é estratégica para diminuir os efeitos de variáveis com **maiores escalas** de medidas sobre a determinação das métricas de distância (WITTEN et al., 2017).

Pré-processamento

Por que?

Os diferentes modelos possuem **sensibilidades distintas** para os tipos de preditores. Assim, a forma com que os preditores entram no modelo também é importante (KUHN; JOHNSON, 2013).

A transformação de preditores é estratégica para diminuir os efeitos de variáveis com **maiores escalas** de medidas sobre a determinação das métricas de distância (WITTEN et al., 2017).

Pré-processamento

Por que?

Os diferentes modelos possuem **sensibilidades distintas** para os tipos de preditores. Assim, a forma com que os preditores entram no modelo também é importante (KUHN; JOHNSON, 2013).

A transformação de preditores é estratégica para diminuir os efeitos de variáveis com **maiores escalas** de medidas sobre a determinação das **métricas de distância** (WITTEN et al., 2017).

Pré-processamento

Boa prática?

Testar uma série de transformações dos dados brutos combinadas com vários algoritmos de aprendizado de máquina, pode ajudar na descoberta de boas representações dos dados e algoritmos melhores capazes de explorar a estrutura dessas representações (BROWNLEE, 2017).

Algoritmos: k -nearest neighbors, Support Vector Machines, Artificial Neural Networks.

Pré-processamento

Função preProcess() - métodos

Pré-processamento

Função `preProcess()` - métodos

a) Métodos para padronização ou normalização de variáveis preditoras:

- "center": Subtrai cada valor x_i da média $\bar{x} = (x_i - \bar{x})$.
- "scale": Divide cada valor x_i pelo desvio padrão $sd_{(x)} = (x_i / sd_{(x)})$.
- "center" e "scale": Padroniza os dados ($\bar{x} = 0$ e $sd_{(x)} = 1$).
- "range": Dimensiona os dados no intervalo [0, 1].

b) Métodos para transformação → distribuição mais simétrica:

- "BoxCox": Uma transformação Box-Cox (diferentes de zero e positivos).

Pré-processamento

Função preProcess() - métodos

a) Métodos para padronização ou normalização de variáveis preditoras:

- "center": Subtrai cada valor x_i da média $\bar{x} = (x_i - \bar{x})$.
- "scale": Divide cada valor x_i pelo desvio padrão $sd_{(x)} = (x_i / sd_{(x)})$.
- "center" e "scale": Padroniza os dados ($\bar{x} = 0$ e $sd_{(x)} = 1$).
- "range": Dimensiona os dados no intervalo [0, 1].

b) Métodos para transformação → distribuição mais simétrica:

- "BoxCox": Uma transformação Box-Cox (diferentes de zero e positivos).

Pré-processamento

Função `preProcess()` - métodos

c) Métodos de imputação de dados:

- "knnImpute": Encontra os k vizinhos mais próximos (euclidiana) → média.
- "medianImpute": Imputação via medianas de cada preditor.
- "bagImpute": Imputação via bagging.

d) Outros métodos: "YeoJohnson", "expoTrans", "pca", "ica", "spatialSign", "corr", "zv", "nzv", and "conditionalX".

Obs.: A função `preProcess` apenas estima os parâmetros necessários para cada método. Em seguida, deve-se usar a função `predict.preProcess` para aplicar o(s) método(s) em conjuntos de dados específicos.

Pré-processamento

Função `preProcess()` - métodos

c) Métodos de imputação de dados:

- "knnImpute": Encontra os k vizinhos mais próximos (euclidiana) → média.
- "medianImpute": Imputação via medianas de cada preditor.
- "bagImpute": Imputação via bagging.

d) Outros métodos: "YeoJohnson", "expoTrans", "pca", "ica", "spatialSign", "corr", "zv", "nzv", and "conditionalX".

Obs.: A função `preProcess` apenas estima os parâmetros necessários para cada método. Em seguida, deve-se usar a função `predict.preProcess` para aplicar o(s) método(s) em conjuntos de dados específicos.

Pré-processamento

Função `preProcess()` - métodos

c) Métodos de imputação de dados:

- "knnImpute": Encontra os k vizinhos mais próximos (euclidiana) → média.
- "medianImpute": Imputação via medianas de cada preditor.
- "bagImpute": Imputação via bagging.

d) Outros métodos: "YeoJohnson", "expoTrans", "pca", "ica", "spatialSign", "corr", "zv", "nzv", and "conditionalX".

Obs.: A função `preProcess` apenas estima os parâmetros necessários para cada método. Em seguida, deve-se usar a função `predict.preProcess` para aplicar o(s) método(s) em conjuntos de dados específicos.

Pré-processamento

Função preProcess() - Exemplo intuitivo

```

df <- data.frame(
  especie = c("Acapu", "Araucaria", "Mogno", "Cedro", "Ipe"),
  diâmetro = c(NA, 27.0, 33.6, 42.6, 52.1),
  altura = c(8.4, 8.7, 9.1, NA, 15.4),
  cortar = c("Não", "Não", "Não", "Não", "Sim"),
  stringsAsFactors = TRUE)
df

```

especie	diametro	altura	cortar
Acapu	NA	8.4	Não
Araucaria	27.0	8.7	Não
Mogno	33.6	9.1	Não
Cedro	42.6	NA	Não
Ipe	52.1	15.4	Sim

Pré-processamento

Função preProcess() - "center" → "scale"

```
# Center & Scale  
library(caret)  
CentsCl <- preProcess(x=df[,2:3], method = c("center", "scale"))  
predict(CentsCl, df)
```

especie	diametro	altura	cortar
Acapu	NA	-0.5977922	Não
Araucaria	-1.0830748	-0.5081234	Não
Mogno	-0.4785680	-0.3885650	Não
Cedro	0.3457596	NA	Não
Ipe	1.2158832	1.4944806	Sim

Pré-processamento

Função preProcess() - "center" → "scale" → "knnImpute"

```
#knnImpute  
knnImpute <- preProcess(x=df[,2:3], method = c("knnImpute"), k=3, knnSummary = mean)  
predict(knnImpute, df)
```

especie	diametro	altura	cortar
Acapu	-0.1152532	-0.5977922	Não
Araucaria	-1.0830748	-0.5081234	Não
Mogno	-0.4785680	-0.3885650	Não
Cedro	0.3457596	0.1992641	Não
Ipe	1.2158832	1.4944806	Sim

Divisão dos dados

Função `createDataPartition()`

Divisão dos dados

Em geral, em ML reporta-se a três conjuntos de dados (WITTEN et al., 2017):

- **Conjunto de aprendizado (treino):** Usado para construir os modelos preditivos (classificação ou regressão).
- **Conjunto de validação:** Usado para otimizar os hiperparâmetros dos modelos e escolher a configuração de melhor desempenho preditivo.
- **Conjunto de teste:** Usado para estimar o desempenho preditivo final do modelo otimizado (*tuning model*).

Divisão dos dados

Importante - Conjuntos independentes

"Cada um dos três conjuntos deve ser escolhido de forma independente"
(WITTEN et al., 2017)

- O conjunto de validação deve ser diferente do conjunto de treinamento para obter um bom desempenho no estágio de otimização ou seleção; e
- O conjunto de teste deve ser diferente de ambos para obter uma estimativa confiável da taxa de erro real.

Divisão dos dados

Função `createDataPartition()`

A função `createDataPartition` pode ser usada para criar divisões aleatórias estratificadas (*stratified random split*) de um conjunto de dados (KUHN, 2008).

```
createDataPartition(y, times = 1, p = 0.5, list = TRUE, groups = min(5,length(y)))
```

y = variável que servirá de base para o particionamento.

times = número de partições a serem geradas.

p = porcentagem de dados do conjunto de treinamento.

list = argumento lógico indicando o formato de saída dos resultados.

Divisão dos dados

Função `createDataPartition()`

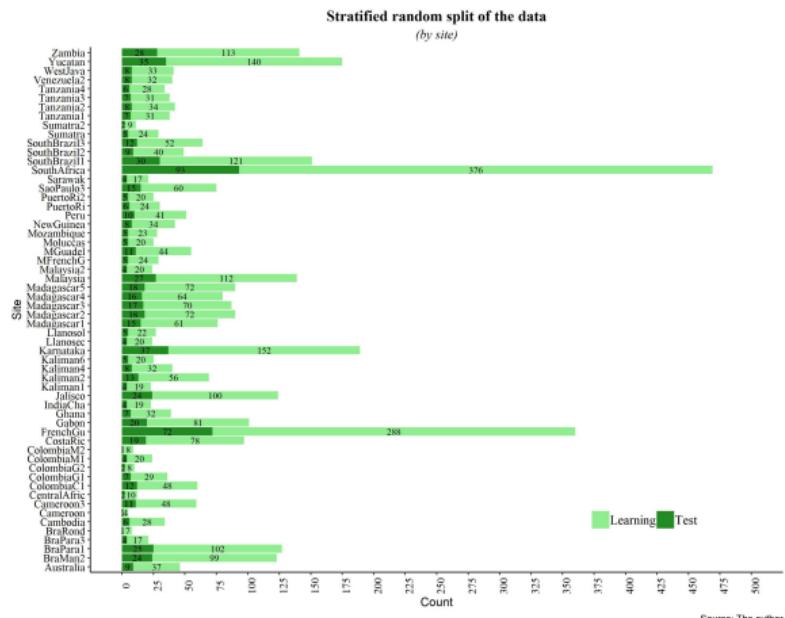
A amostragem aleatória é feita dentro dos níveis de **y**:

Se $y = \text{fator}$ → a amostragem aleatória é feita considerando os níveis de fatores. Assim, busca obter representatividade de todas as classes, bem como equilibrar a quantidade amostrada dentro das classes.

Se $y = \text{numérico}$ → a amostra é dividida em subgrupos com base em percentis. Assim, a amostragem aleatória é feita dentro desses subgrupos (KUHN et al., 2017).

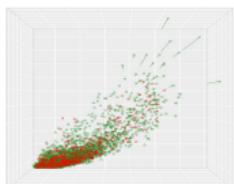
Divisão dos dados

Função createDataPartition() - Split (Biomassa de árvores)



Divisão dos dados

Função `createDataPartition()` - Garantia das propriedades estatísticas



Correlation matrix

Correlation Matrix				
	D	H	WSG	AGB
D	1	0.8351	-0.172	0.7896
H	0.8351	1	-0.2237	0.6151
WSG	-0.172	-0.2237	1	-0.0517
AGB	0.7896	0.6151	-0.0517	1

Descriptive analysis (all data set)

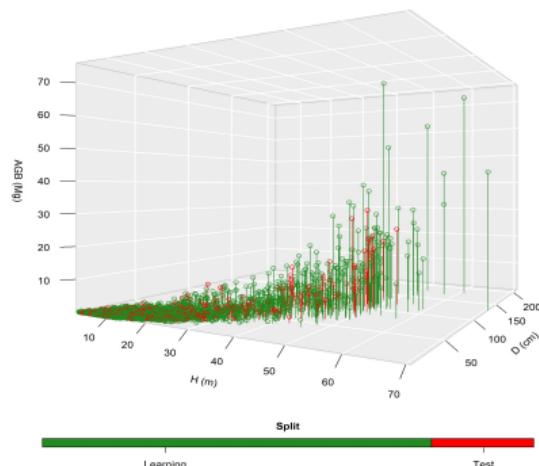
Descriptive statistics (Excel 2010)							
	n	Min	Max	Mean	Stddev	Skewness	Kurtosis
D	4004	5	212	23.9877	24.0853	2.4973	7.6611
H	4004	1.2	70.7	16.0376	10.7719	1.2724	1.4764
WSG	4004	0.09	1.2	0.6335	0.1643	-0.1257	-0.1253
AGR	4004	0.0012	76.0635	1.1341	3.918	8.2436	100.3593

Descriptive analysis (learning data set)

Descriptive Statistics (Running Date: 2023-01-01)							
	n	Min	Max	Mean	Stddev	Skewness	Kurtosis
D	3226	5	212	24.2082	24.5409	2.5326	7.9362
H	3226	1.2	70.7	16.0613	10.8284	1.2842	1.5351
WSG	3226	0.09	1.2	0.6324	0.1635	-0.1177	-0.1508
AGB	3226	0.0012	76.0635	1.1703	4.1065	8.3353	100.2212

Descriptive analysis (test data set)

	n	Min	Max	Mean	Stdev	Skewness	Kurtosis
D	778	5	132	23.0734	22.0875	2.2318	5.2511
H	778	2.1	52.7	15.9395	10.5406	1.2145	1.176
WSG	778	0.1	1.08	0.6382	0.1679	-0.1605	-0.0379
AGR	778	0.0012	30.531	0.9841	3.0107	5.6705	38.1428



Treinamento e Ajuste de Hiperparâmetros

Função train()

Treinamento e Ajuste de Hiperparâmetros

Hyperparameter

Muitos algoritmos de ML possuem "parâmetros" que podem ser ajustados (tuned) para otimizar seu desempenho. Esses parâmetros são denominados **Hiperparâmetros** (WITTEN et al., 2017; PROBST et al., 2018).

Hyperparameter tuning

O termo *hyperparameter tuning* (hyperparameter optimization) pode ser definido como o processo de encontrar boas configurações de hiperparâmetros de um algoritmo para um conjunto de dados específico (PROBST et al., 2018).

Treinamento e Ajuste de Hiperparâmetros



Treinamento e Ajuste de Hiperparâmetros

Alguns métodos e seus hiperparâmetros de ajuste:

Modelo	Método	*Tarefa	Pacote	Hiperparâmetro
k-Nearest Neighbors	knn	C, R	caret	k
SVM with Linear Kernel	svmLinear	C, R	kernlab	C
SVM with Linear Kernel	svmLinear2	C, R	e1071	cost
Multi-Layer Perceptron	mlp	C, R	RSNNS	size
Neural Network	neuralnet	R	neuralnet	layer1,layer2,layer3
Neural Network	nnet	C, R	nnet	size, decay
Ridge Regression	ridge	R	elasticnet	lambda
CART	rpart	C, R	rpart	cp
Random Forest	rf	C, R	randomForest	mtry

*C = Classificação; R = Regressão



Treinamento e Ajuste de Hiperparâmetros

Função train()

A função `train()` possui as seguintes utilidades:

- **Capacidade preditiva:** Obter a estimativa da capacidade preditiva (performance) para diferentes combinações de hiperparâmetros (candidatos) com base na amostra de treinamento, usando técnicas de reamostragem; e
- **Modelo final:** Indicar o melhor modelo preditivo. Isto é, aquele com hiperparâmetros de ajuste ótimo (optimal tuning hyperparameters). Este é o modelo final (tuning model).

Treinamento e Ajuste de Hiperparâmetros

Função train()

Parâmetros:

- **x**: matriz ou dataframe que contém as variáveis preditoras;
- **y**: vector que contém a variável de resposta;
- **form**: formula para indicar as variáveis preditoras e a resposta;
- **data**: dataframe que contém o conjunto de dados;
- **method**: método de construção do modelo preditivo (algoritmos);
- **preProcess**: pré-processamento das variáveis preditoras;
- **metric**: métrica de avaliação da capacidade preditiva dos modelos ("RMSE", "Rsquared", "Precisão", "ROC" ...)

Treinamento e Ajuste de Hiperparâmetros

Função train()

Parâmetros (Cont.):

- **trControl**: controlar a construção do modelo e definir a técnica de reamostragem;
- **tuneGrid**: receber um dataframe com os candidatos a hiperparâmetro de ajuste ótimo; e
- **tuneLength**: número de níveis para cada hiperparâmetro de ajuste (usar apenas caso tuneGrid não esteja especificado).

Treinamento e Ajuste de Hiperparâmetros

Função train() - Parâmetro trControl()

```
trainControl(method="repeatedcv", number=10, repeats=10)
```

- **method:**

- "none"* = ajusta um único modelo para todo conj. treino.
- "cv"* = k-fold cross-validation.
- "repeatedcv"* = repeated k-fold cross-validation.
- "boot"* = bootstrapping.
- "LOOCV"* = leave-one-out cross-validation.

Treinamento e Ajuste de Hiperparâmetros

Função train() - Esquema de treinamento

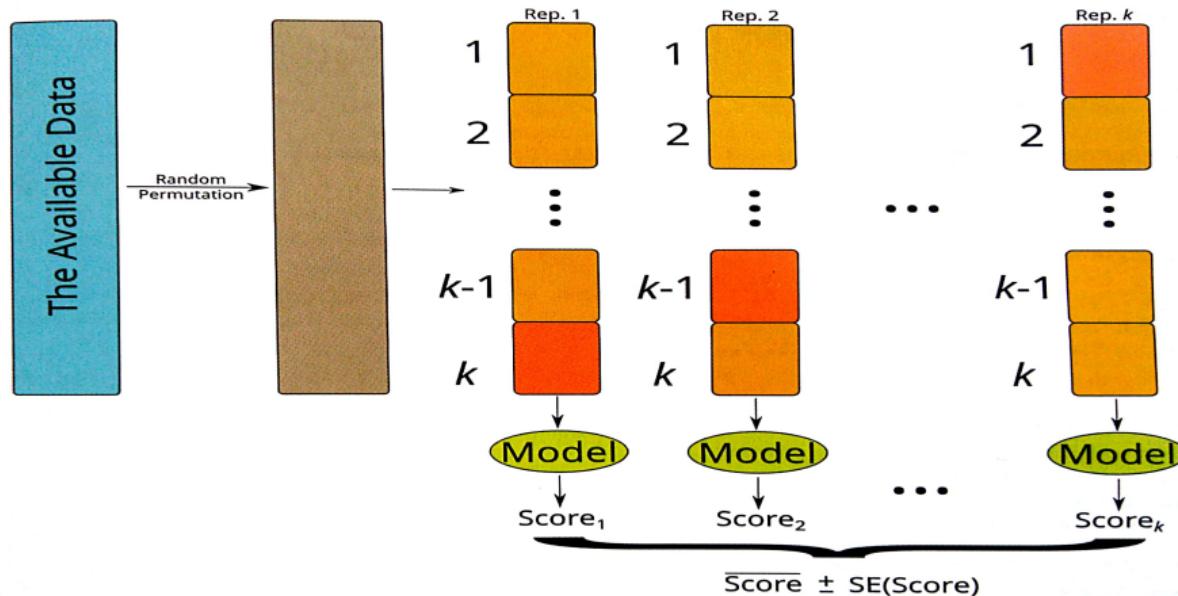
- 1 Define sets of model parameter values to evaluate
- 2 **for** each parameter set **do**
- 3 **for** each resampling iteration **do**
- 4 Hold-out specific samples
- 5 [Optional] Pre-process the data
- 6 Fit the model on the remainder
- 7 Predict the hold-out samples
- 8 **end**
- 9 Calculate the average performance across hold-out predictions
- 10 **end**
- 11 Determine the optimal parameter set
- 12 Fit the final model to all the training data using the optimal parameter set

Treinamento e Ajuste de Hiperparâmetros

Métodos de Reamostragem

- **Hold-out validation** : Separa 2/3 para treinamento e 1/3 para teste.
- **k-fold Cross-Validation:** Particiona o conj. de treinamento em k subconjuntos disjuntos ($k-1$), e testa no conjunto de validação (*hold-out*) (80%-20%).
- **Leave-one-out Cross-Validation:** $k=n$
- **Bootstrap:** Amostragem com reposição.

Treinamento e Ajuste de Hiperparâmetros



Fonte: (TORGÓ, 2017) - (Demonstre!)

Mensagem Final

Recomendação

Aprender uma linguagem de programação e melhorar as bases matemáticas e estatísticas.



Mensagem Final

Avanço

Precisamos **compreender melhor** e avançar na implementação de algoritmos de **Aprendizagem de Máquina** para predição de variáveis nas Ciências Florestais!

Bibliografia I

- BREIMAN, L. et al. **Classification and Regression Trees**. Monterey, CA: Wadsworth and Brooks, 1984.
- KUHN, M. Building predictive models in r using the caret package. **Journal of Statistical Software, Articles**, v. 28, n. 5, p. 1–26, 2008. ISSN 1548-7660. Disponível em: <<https://www.jstatsoft.org/v028/i05>>.
- KUHN, M. **caret: Classification and Regression Training**. [S.I.], 2018. R package version 6.0-79. Disponível em: <<https://CRAN.R-project.org/package=caret>>.
- KUHN, M.; JOHNSON, K. **Applied predictive modeling**. [S.I.]: Springer, 2013. v. 810.

Bibliografia II

- KUHN, M. et al. **caret: Classification and Regression Training**. [S.I.], 2017. R package version 6.0-78. Disponível em: <<https://CRAN.R-project.org/package=caret>>.
- PROBST, P. et al. Tunability: Importance of hyperparameters of machine learning algorithms. **arXiv preprint arXiv:1802.09596**, 2018.
- THERNEAU, T. et al. **rpart: Recursive Partitioning and Regression Trees**. [S.I.], 2017. R package version 4.1-11. Disponível em: <<https://CRAN.R-project.org/package=rpart>>.
- TORGÓ, L. **Data mining with R: learning with case studies**. [S.I.]: CRC press, 2017.
- WITTEN, I. H. et al. **Data Mining: Practical machine learning tools and techniques**. [S.I.]: Morgan Kaufmann, 2017.

OBRIGADO!

Deivison Venicio Souza (UFPA)
Email: deivisonvs@ufpa.br

