



Grupo de Usuarios de R de Sevilla

Selección de variables en Machine Learning

Miguel Pavón Barrera

2 Abril 2019, 19 a 20:30 h

Sala TIC4, CRAI Reina Mercedes

<http://bit.ly/SevillaRmeetup>

INTRODUCCIÓN

- OBJETIVO
- PRINCIPALES MÉTODOS
- IMPLEMENTACIÓN EN R



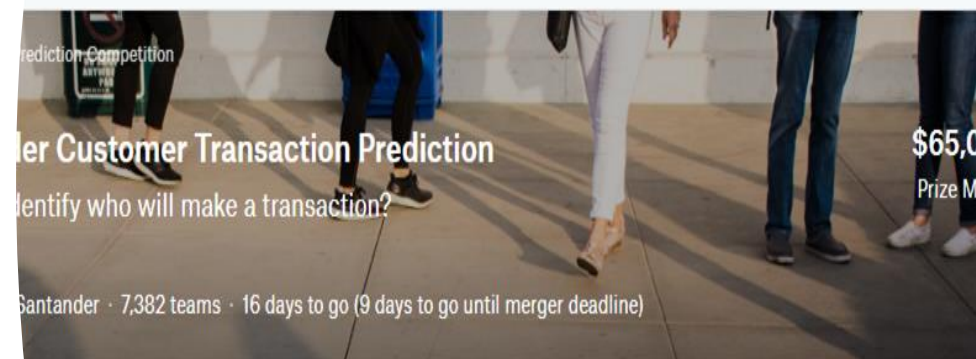
CONJUNTO DE PRUEBAS

- DATASET TRAIN.CSV
- 200 000 OBSERVACIONES
- 200 VARIABLES PREDICTORAS
- PROBLEMA DE CLASIFICACIÓN
- AÑADIMOS 60 VARIABLES RANDOM
- PARTICIONAMOS EN TRAIN Y TEST

SCORE TEST: 0,9630

- DATASET TRAIN.CSV
- 1460 OBSERVACIONES
- 80 VARIABLES PREDICTORAS
- PROBLEMA DE REGRESIÓN
- AÑADIMOS 10 VARIABLES RANDOM

SCORE VAL: 0.1314
SCORE TRAIN: 0.0743



[Data](#) [Kernels](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Team](#) [My Submissions](#) [Submit Prediction](#)

1

with an anonymized dataset containing numeric feature variables, the binary `target` column, and a string `ID_code`.

the value of `target` column in the test set.

The test set contains some rows which are not included in scoring.
Please make sure your submission file is in the correct format.

[API](#) [kaggle competitions download -c santander-custom...](#) [?](#) [Download All](#)

About this file

Description yet

Columns

▲ `ID_code`

`target`

TIPOS DE VARIABLES

- RELEVANTES
- REDUNDANTES
- IRRELEVANTES

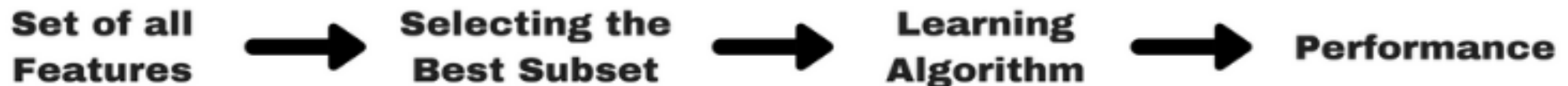


CLASIFICACIÓN MÉTODOS

- **FILTER METHOD – METODOS FILTRO**
- **WRAPPER METHOD – METODOS ENVOLTURA**
- **EMBEDDED METHOD – MÉTODOS EMPOTRADOS**

FILTER METHOD - CARACTERÍSTICAS

- **NO SE BASAN EN ALGORITMOS DE MACHINE LEARNING**
- **MIDEN LA RELACIÓN ENTRE ATRIBUTOS Y VARIABLE RESPUESTA**
- **SIMPLES Y RÁPIDOS COMPUTACIONALMENTE**
- **INDEPENDIENTES DEL CLASIFICADOR**



- **NO DETECTAN ATRIBUTOS REDUNDANTES**
- **IGNORAN LA INTERACCIÓN ENTRE ATRIBUTOS**

FILTER METHOD – PRINCIPALES TÉCNICAS

- CORRELACIÓN
- ANOVA
- CHI.CUADRADO
- ENTROPY METHODS

Feature\Response	Continuous	Categorical
Continuous	Pearson's Correlation	LDA
Categorical	Anova	Chi-Square

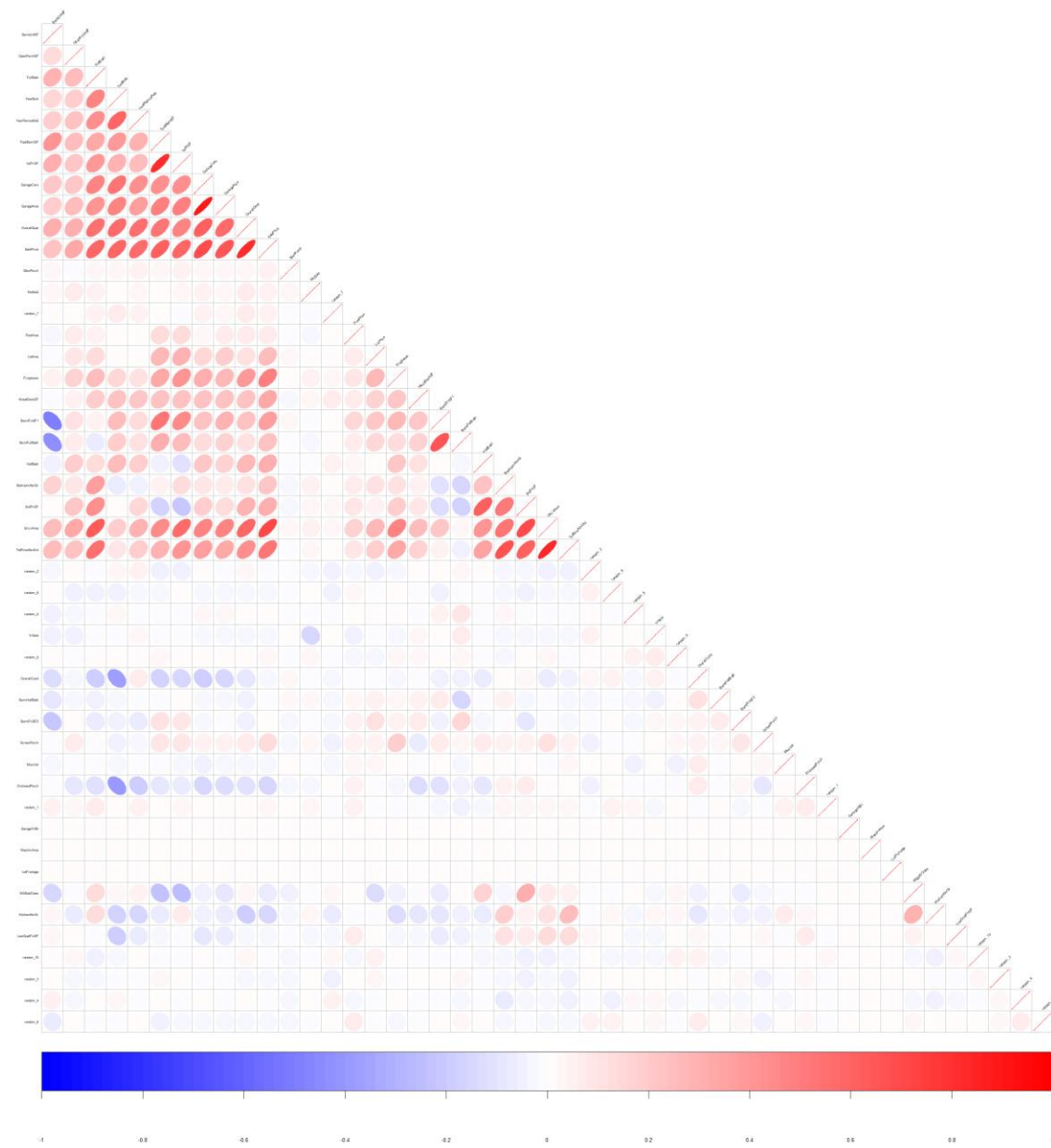
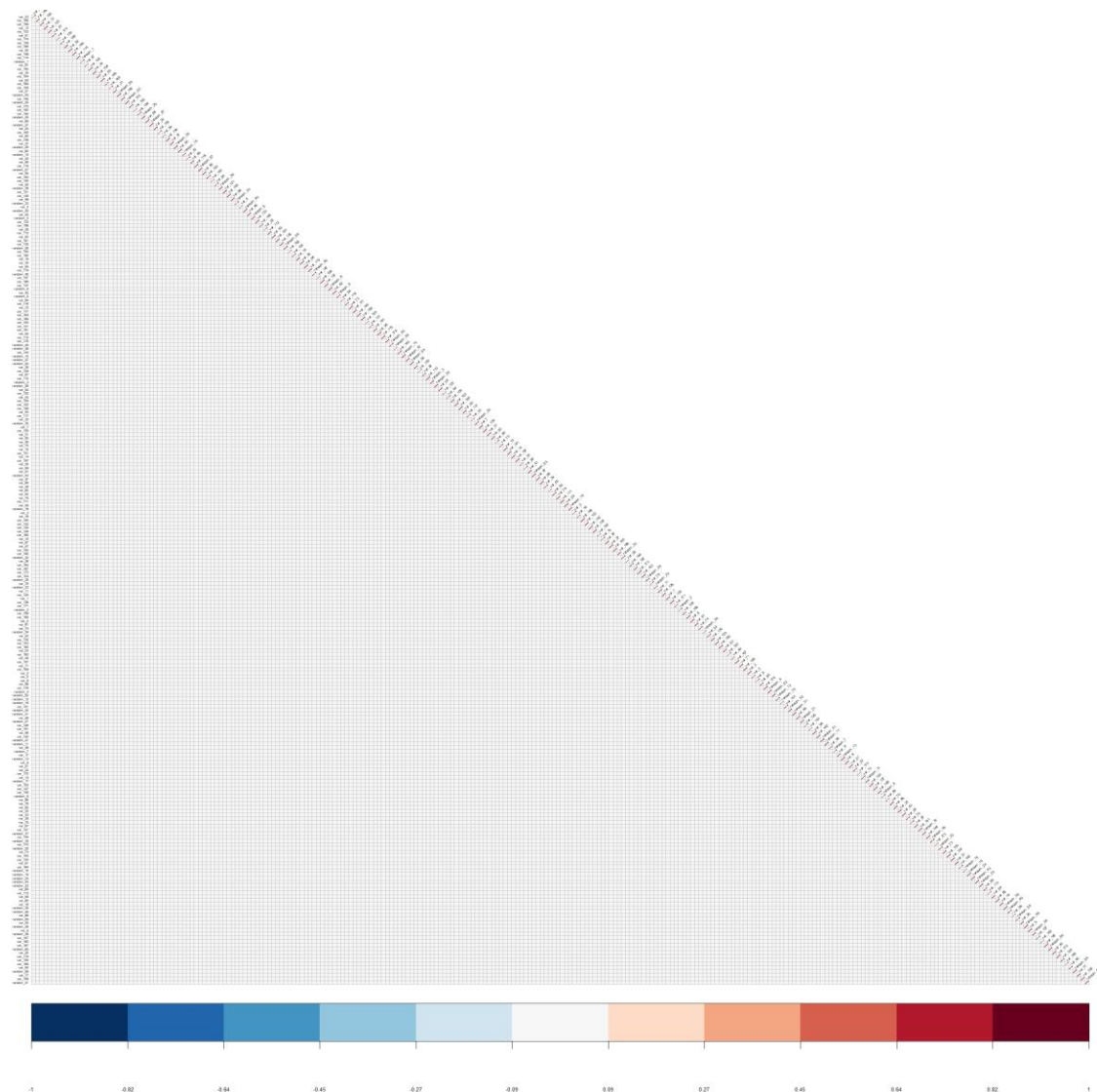
FILTER METHOD – CORRELACIÓN LINEAL

- Mide **dependencia lineal** entre predictores continuos.

Métodos Multivariantes:

- **Método CFS** (CorrelationbasedFeatureSelection)
- **FCBF** (FastCorrelation-BasedFilter)

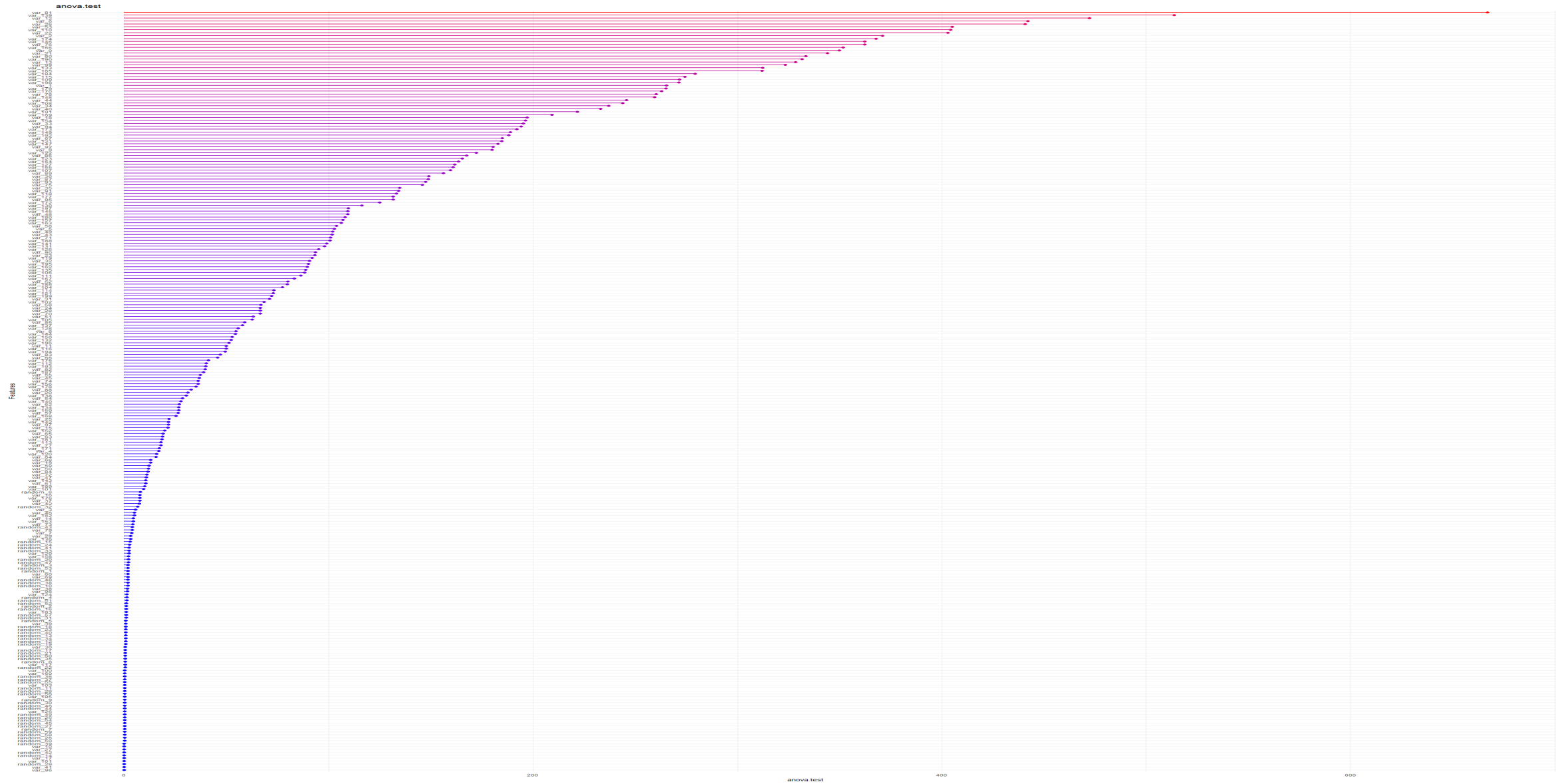
Correlación



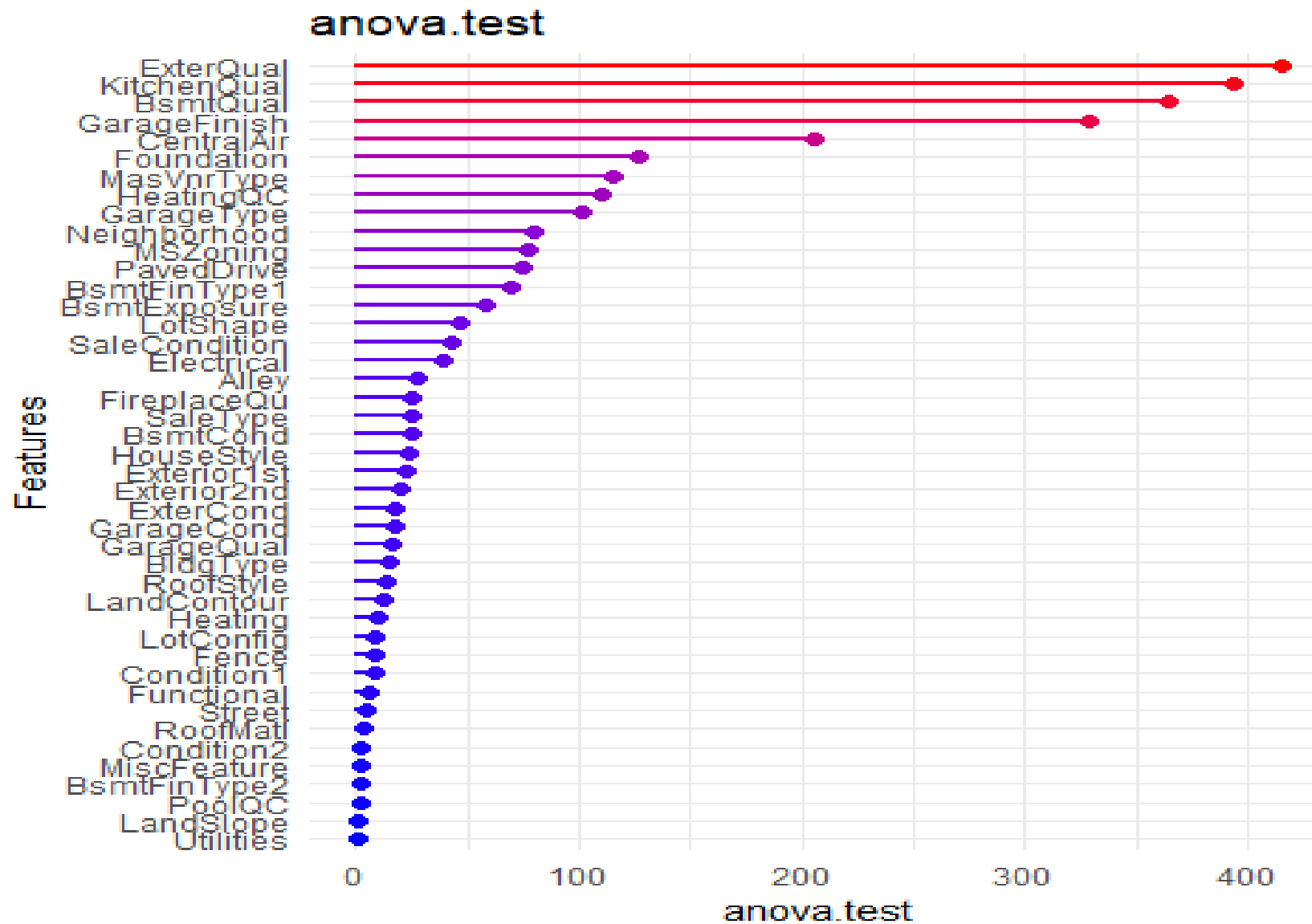
FILTER METHOD - ANOVA

- ESTADÍSTICO F
- REFLEJA EL **GRADO DE PARECIDO** EXISTENTE ENTRE LAS MEDIAS DE UNA **VARIABLE CATEGÓRICA SOBRE UNA CONTINUA**
- SI $F \gg 1$ ENTONCES HAY MAYOR EVIDENCIA
- **SE NECESITA HIPÓTESIS PARA REALIZAR CONTRASTE**

ESTADISTICO F - SANTANDER



ESTADISTICO F - HP



FILTER METHOD - ENTROPÍA

- **MEDIDA DE INCERTIDUMBRE** DE UNA VARIABLE ALEATORIA

$$H(X) = - \sum_i P(x_i) \log P(x_i)$$

Entropía de X después de observar otra variable Y :

$$H(X | Y) = - \sum_j P(y_j) \sum_i P(x_i | y_j) \log P(x_i | y_j)$$

FILTER METHOD - ENTROPÍA

Details

`information.gain` is

$$H(Class) + H(Attribute) - H(Class, Attribute)$$

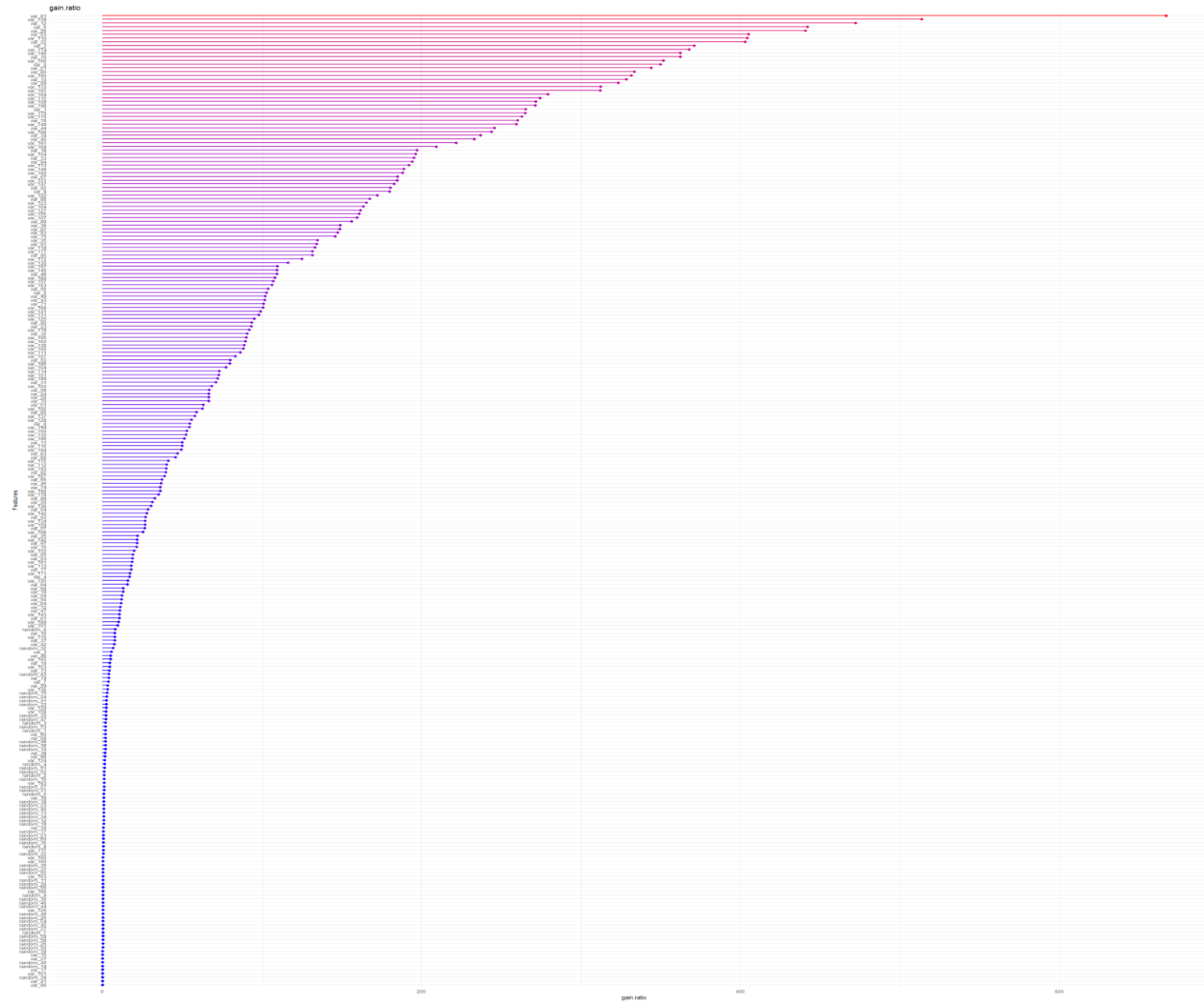
`gain.ratio` is

$$\frac{H(Class) + H(Attribute) - H(Class, Attribute)}{H(Attribute)}$$

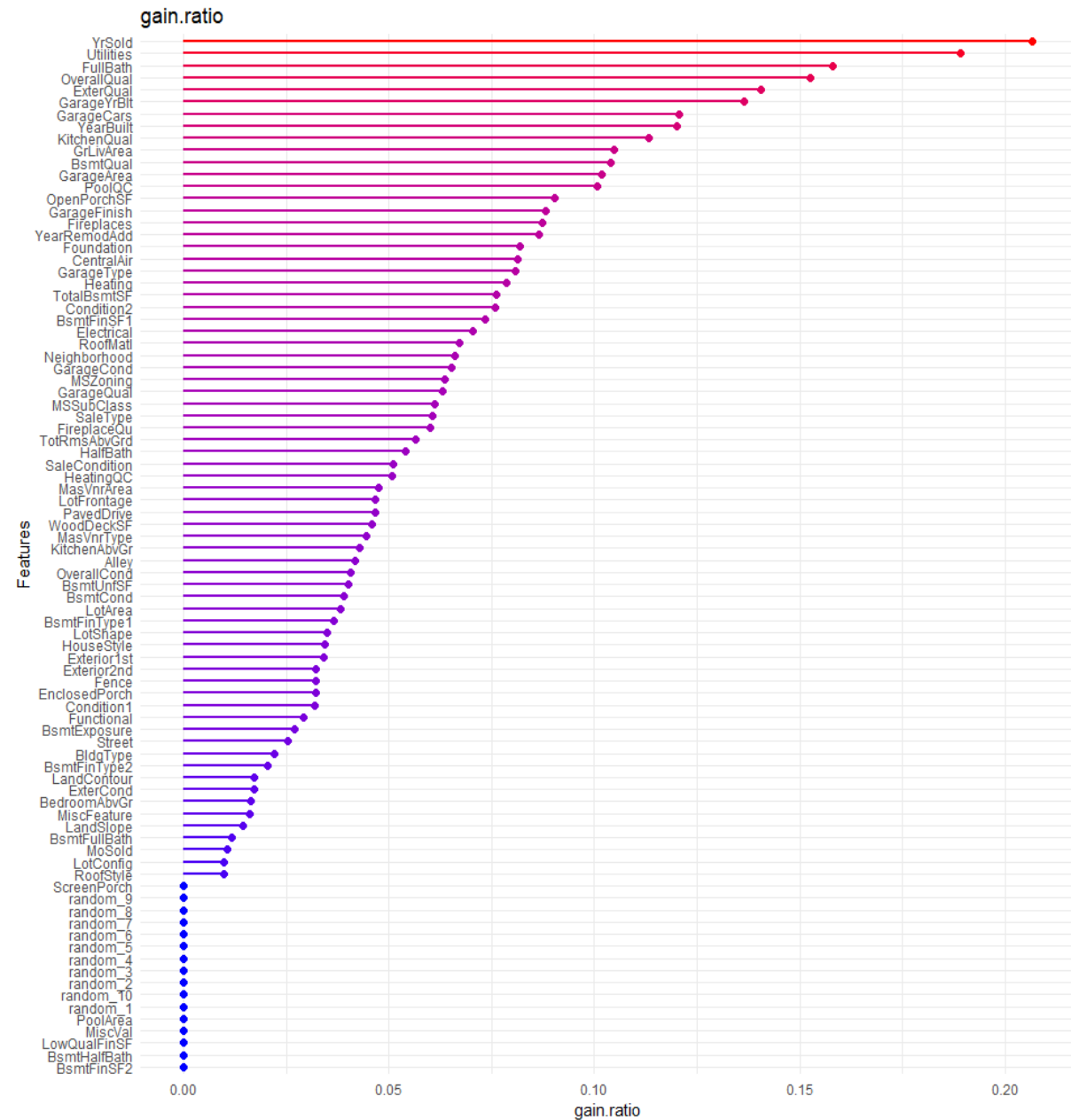
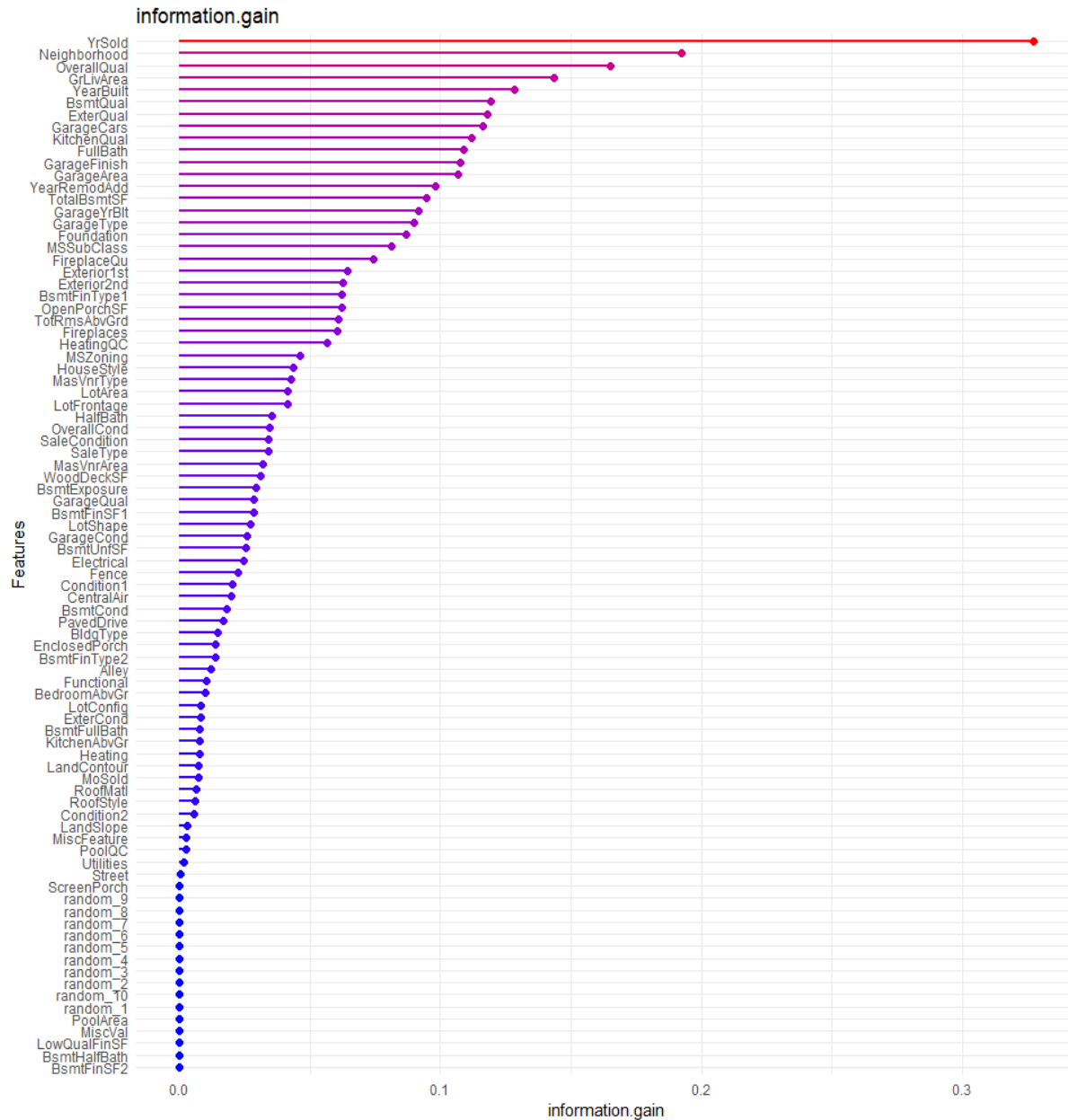
`symmetrical.uncertainty` is

$$2 \frac{H(Class) + H(Attribute) - H(Class, Attribute)}{H(Attribute) + H(Class)}$$

GAIN RATIO - SANTNADER



GAIN - HP



FILTER METHOD – IMPLEMENTACIÓN R

- LIBRERÍA MRL

```
generateFilterValuesData(  
  task      = _____ ,  
  method = _____ )
```

```
makeClassifTask (data, target)
```

```
makeMultilabelTask (data, target)
```

```
makeRegrTask (data, target)
```

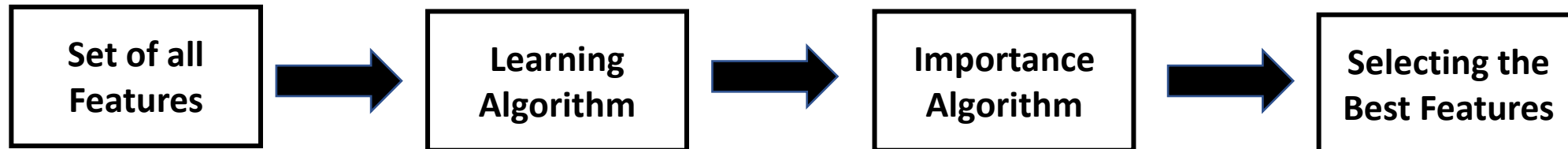
```
listFilterMethods( )
```

id	package	description
1	anovaTest	Anova test for binary and multiclass classification tasks
2	auc	AUC filter for binary classification tasks
3	chiSquare	Chi-squared statistic of independence between feature and target
4	featureImportance	Permutation importance of random forest fitted in package ...
5	fisher	Fisher's discriminant analysis
6	chiSquare	Chi-squared statistic of independence between feature and target
7	gainRatio	Entropy-based gain ratio between feature and target
8	informationGain	Entropy-based information gain between feature and target
9	knnRank	K-nearest neighbor rank
10	linearCorrelation	Linear correlation between feature and target
11	mrmr	Minimum redundancy maximum relevance filter
12	mrml	MRML filter
13	permutationImportance	Permutation importance of random forest fitted in package ...
14	randomForestImportance	Importance of random forest fitted in package ...
15	randomForestVADetect	Minimal depth of 2 variable subset via method varSelect on ...
16	rangerImportance	Variable importance based on ranger permutation importance
17	rangerPermutation	Variable importance based on ranger permutation importance
18	relief	RELIEF algorithm
19	symmetricalUncertainty	Entropy-based symmetrical uncertainty between feature and target
20	univariateModelScore	Resamples an mlr learner for each input feature individually ...
21	variance	A simple variance filter

	id	package	desc
1	anova.test		ANOVA Test for binary and multiclass classification tasks
2	auc		AUC filter for binary classification tasks
3	carscore	care	CAR scores
4	cforest.importance	party	Permutation importance of random forest fitted in package '...
5	chi.squared	FSelector	Chi-squared statistic of independence between feature and t...
6	gain.ratio	FSelector	Entropy-based gain ratio between feature and target
7	information.gain	FSelector	Entropy-based information gain between feature and target
8	kruskal.test		Kruskal Test for binary and multiclass classification tasks
9	linear.correlation		Pearson correlation between feature and target
10	mrmr	mRMRe	Minimum redundancy, maximum relevance filter
11	oneR	FSelector	oneR association rule
12	permutation.importance		Aggregated difference between feature permuted and unpe...
13	randomForest.importance	randomForest	Importance based on OOB-accuracy or node impurity of ran...
14	randomForestSRC.rfsrc	randomForestS...	Importance of random forests fitted in package 'randomFor...
15	randomForestSRC.var.select	randomForestS...	Minimal depth of / variable hunting via method var.select on...
16	ranger.impurity	ranger	Variable importance based on ranger impurity importance
17	ranger.permutation	ranger	Variable importance based on ranger permutation importan...
18	rank.correlation		Spearman's correlation between feature and target
19	relief	FSelector	RELIEF algorithm
20	symmetrical.uncertainty	FSelector	Entropy-based symmetrical uncertainty between feature and...
21	univariate.model.score		Resamples an mlr learner for each input feature individually....
22	variance		A simple variance filter

EMBEDDED METHOD – CARACTERÍSTICAS

- SE BASAN EN ALGORITMOS DE **MACHINE LEARNING**
- ASIGNAN UN GRADO DE IMPORTANCIA A CADA VARIABLE
- RELACIONADO CON LA INTERPRETABILIDAD DE LOS MODELOS



EMBEDDED – PRINCIPALES TÉCNICAS

- REGRESIÓN RIDGE
- REGRESIÓN LASSO
- SVM IMPORTANCE
- RANDOMFOREST IMPORTANCE
- XGBOOST IMPORTANCE
- LIGTHGBM IMPORTANCE
- ETC

EMBEDDED METHOD – REGRESIÓN LASSO

- REGRESIÓN LINEAL **CON RESTRICCIONES**

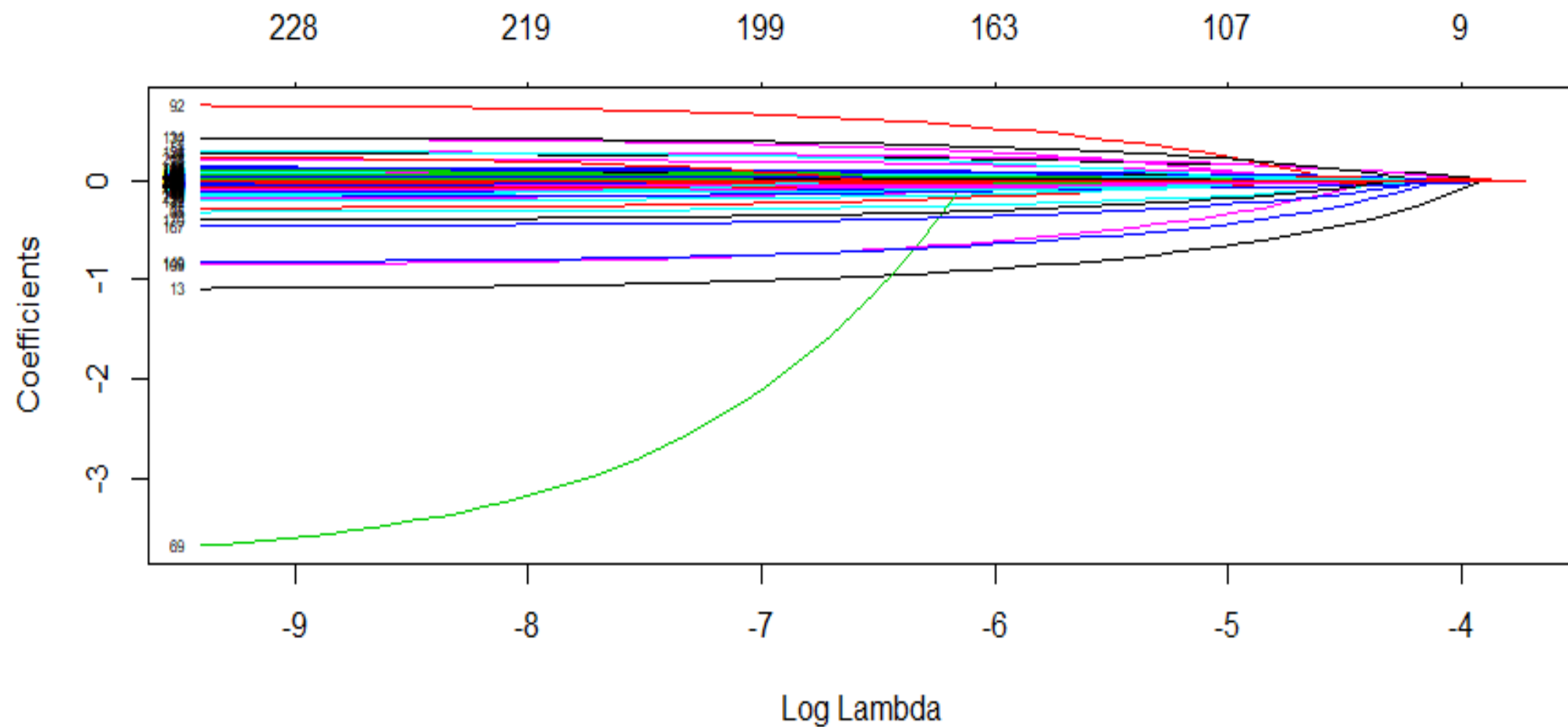
$$\min_{\beta_0, \beta} \left\{ \frac{1}{N} \sum_{i=1}^N (y_i - \beta_0 - x_i^T \beta)^2 \right\} \text{ subject to } \sum_{j=1}^p |\beta_j| \leq t. \quad \longrightarrow \quad \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{N} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

- λ PARÁMETRO REGULARIZACIÓN/PENALIZACIÓN
- **IDEA:** variar λ y comparar la **DEVIANCE**

EMBEDDED METHOD – REGRESIÓN LASSO

```
glmnet( x, y,  
        family = {'binomial', 'Gaussian' ,etc}  
        alpha = 1  
        standardize = {T,F}  
        nlambdas = { 100 , 200 ... } )  
glmnet.cv( )
```

LASSO PLOT



EMBEDDED – RANDOMFOREST IMP.

- Random Forest construye CART de forma paralela
- Cada árbol **MINIMIZA** función de **IMPUREZA**

Impurity	Task	Formula	Description
Gini impurity	Classification	$\sum_{i=1}^C -f_i(1 - f_i)$	f_i is the frequency of label i at a node and C is the number of unique labels.
Entropy	Classification	$\sum_{i=1}^C -f_i \log(f_i)$	f_i is the frequency of label i at a node and C is the number of unique labels.
Variance / Mean Square Error (MSE)	Regression	$\frac{1}{N} \sum_{i=1}^N (y_i - \mu)^2$	y_i is label for an instance, N is the number of instances and μ is the mean given by $\frac{1}{N} \sum_{i=1}^N y_i$

- **OOB** – Conjunto de Validación

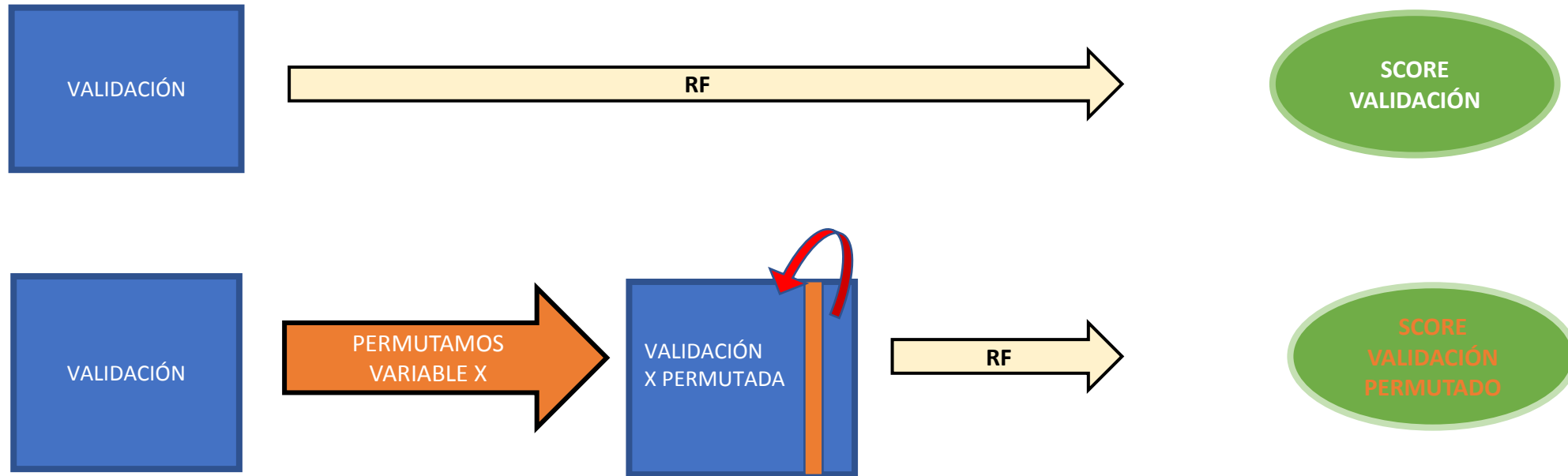
EMBEDDED – RANDOMFOREST IMP.

- **MEDIA DECREASE IMPURITY**

- Mide cuánta “**información**” aporta una variable
- Gini Importance, Information Importance...
- **SESGADO para VARIABLES con ALTA CARDINALIDAD**

EMBEDDED – RANDOMFOREST IMP.

- PERMUTATION IMPORTANCE RF

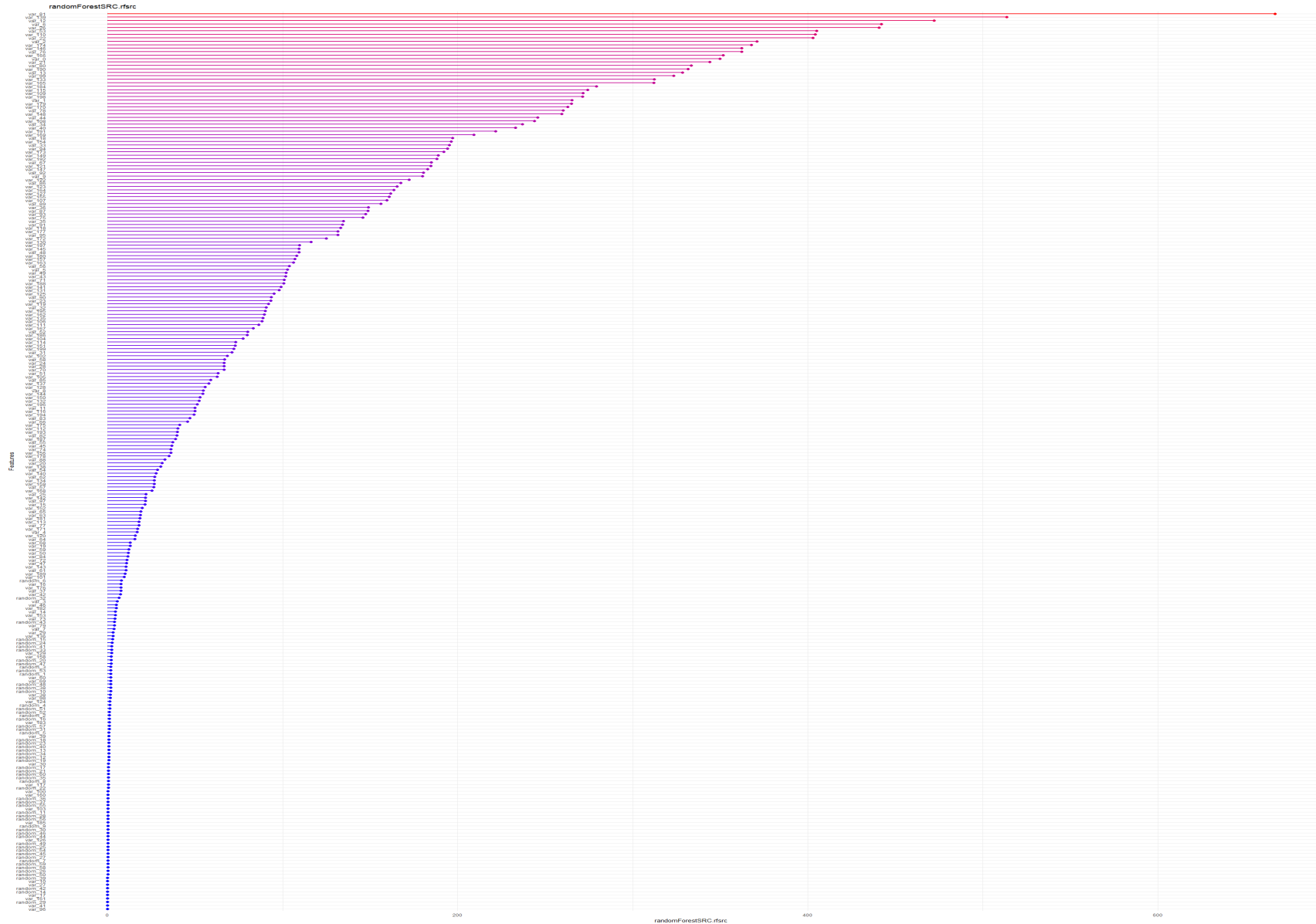


IMPORTANCIA VARIABLE X = SCORE VALIDACION – SCORE VALIDACION PERMUTADO

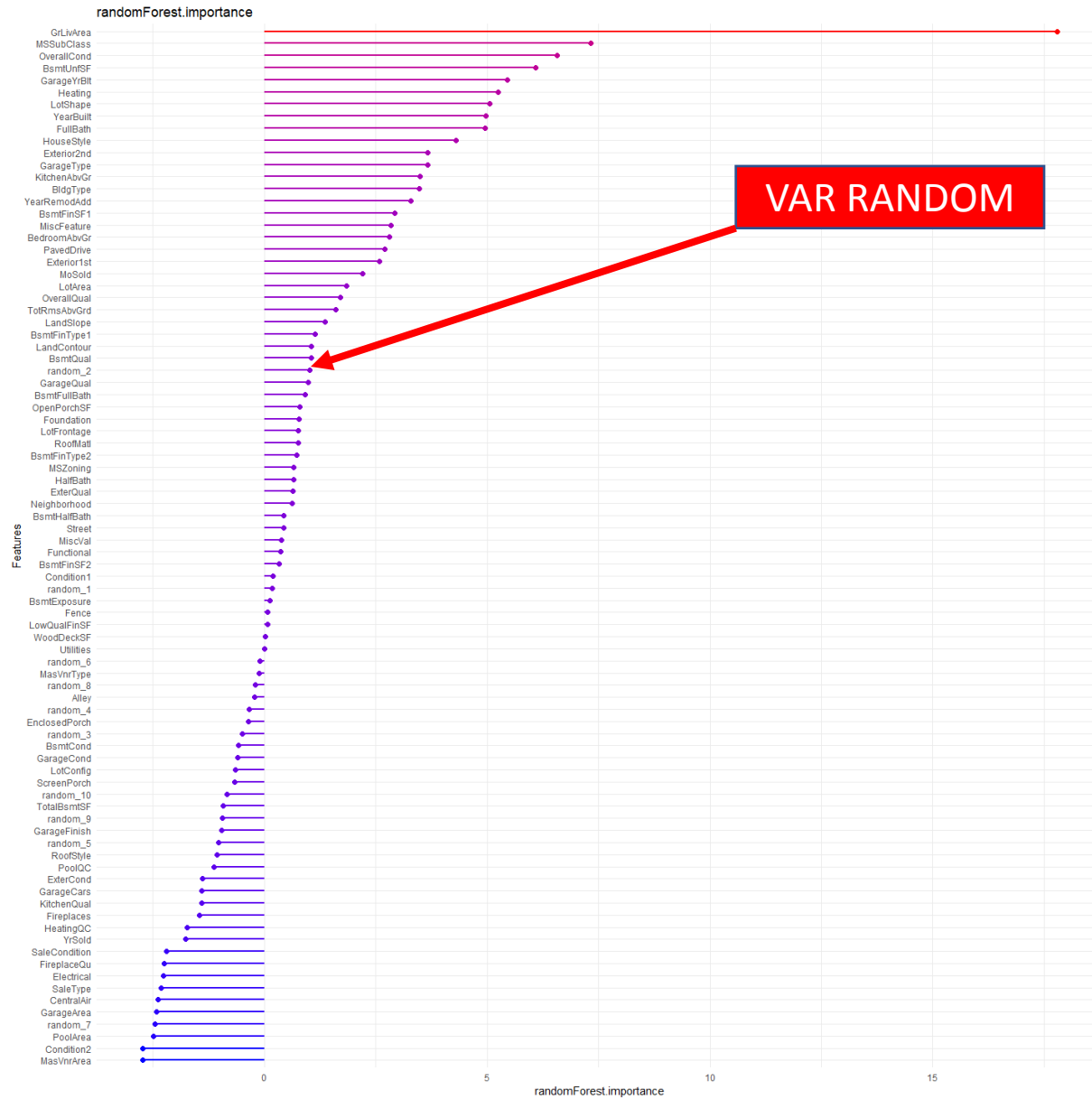
EMBEDDED – RANDOMFOREST IMP.

- PERMUTATION IMPORTANCE
 - Descenso de validación al permutar variable
 - RF validación OOB
 - **Mayor COSTE COMPUTACIONAL**
 - **SESGADO CORRELACIONADAS**

RF PERMUTE IMPORTANCE - SANTANDER



RF PERMUTE IMPORTANCE - HP



RF IMPORTANCIA - IMPLEMENTACIÓN

randomForest library

- 1) `randomForest(formula,
data ,
importance = TRUE , __)`
- 2) `importance(RF , type = {1,2})`

randomForestSRC library

- 1) `rfsrc(formula,
data ,
importance = TRUE , __)`
- 2) `vimp(RF , xvar.names)`

`generateFilterValuesData(task.train , method = {...})`

	id	package	desc
1	cforest.importance	party	Permutation im
2	randomForest.importance	randomForest	Importance ba
3	randomForestSRC.rfsrc	randomForestSRC	Importance of
4	randomForestSRC.var.select	randomForestSRC	Minimal depth

EMBEDDED – XGBOOST IMP.

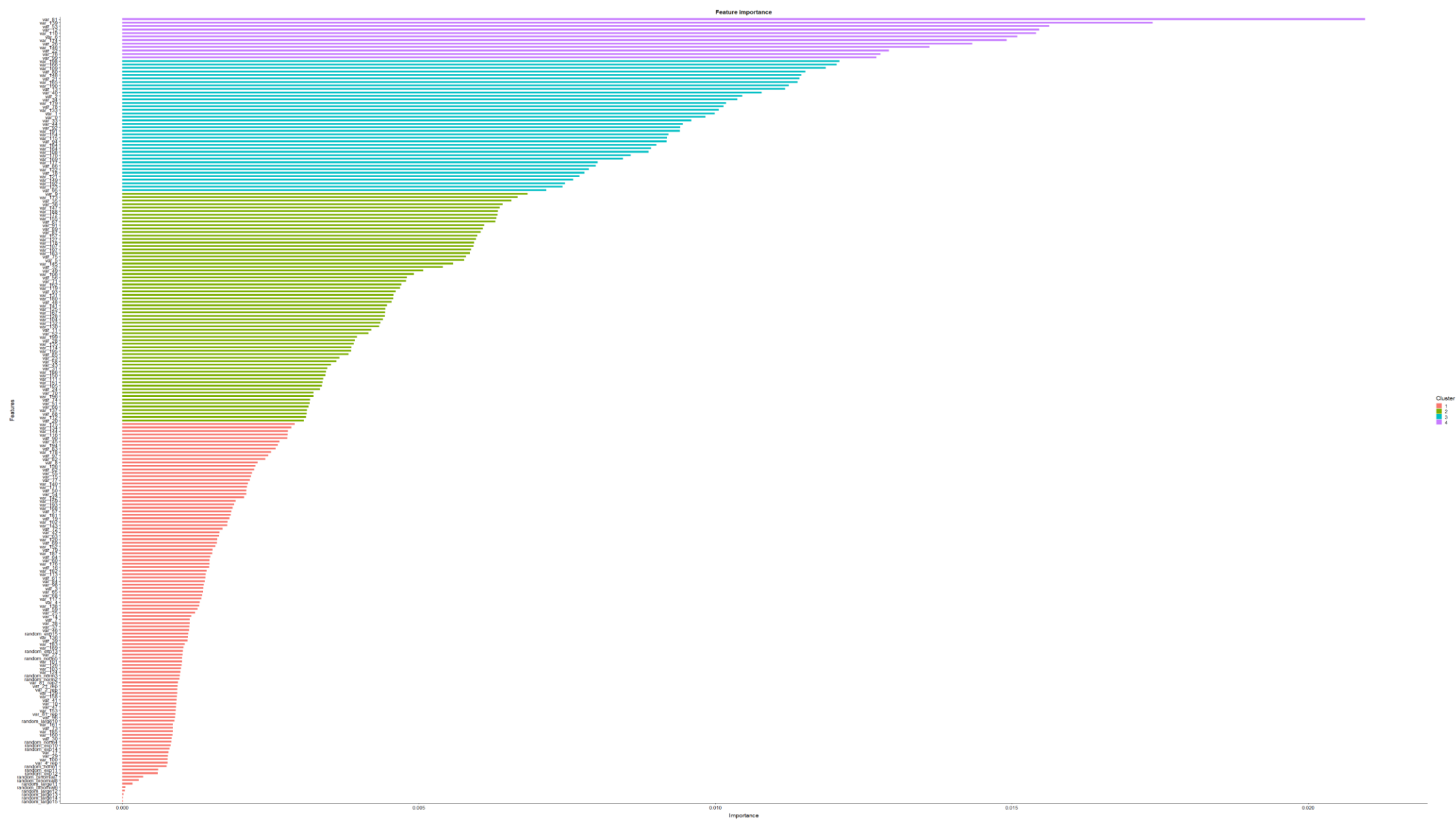
- XGB construye CART de forma secuencial
- **GAIN IMPORTANCE**
- **FREQUENCY IMPORTANCE**
- **COVER IMPORTANCE**

XGBoost

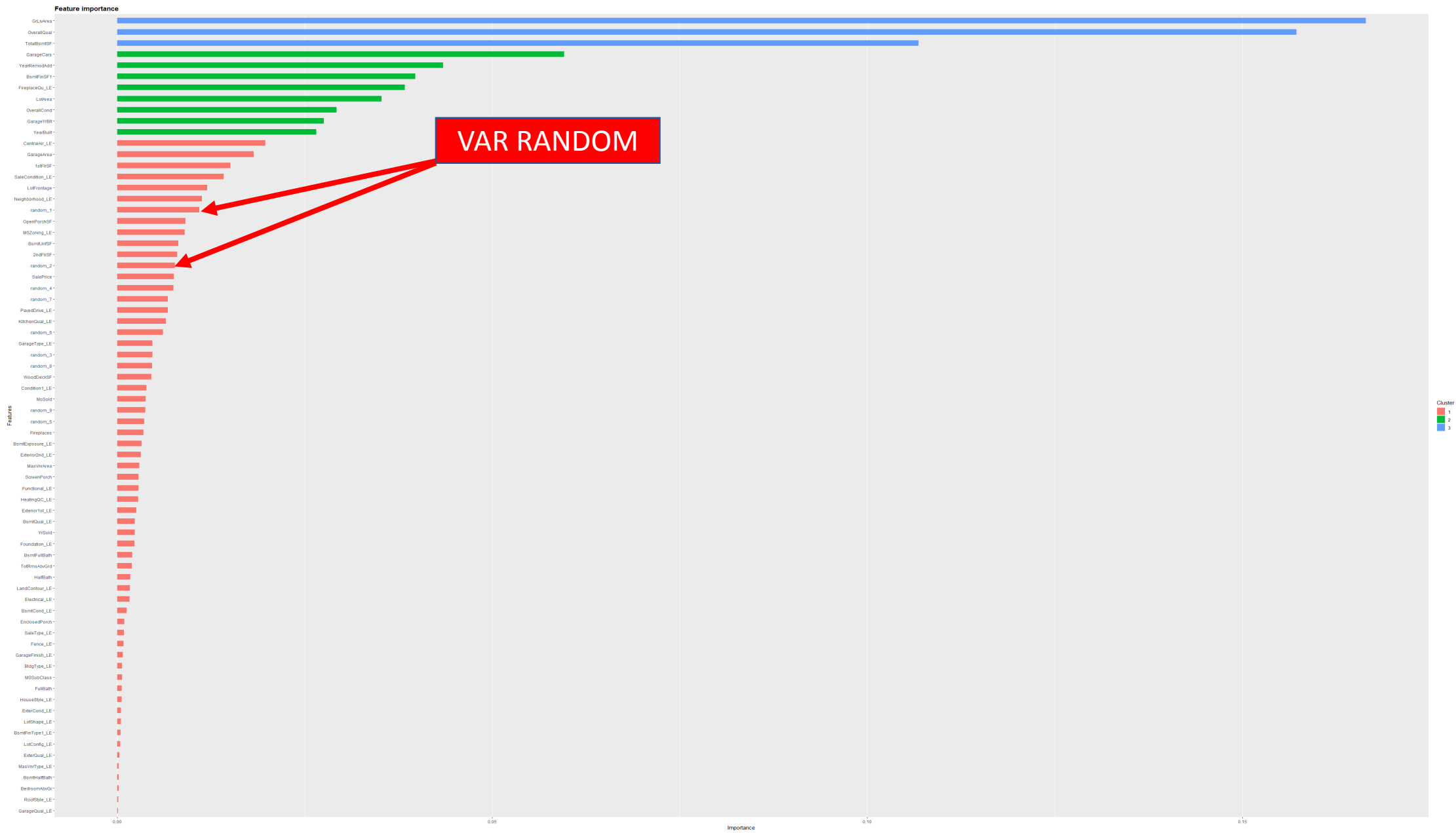
XGB IMPORTANCIA - IMPLEMENTACIÓN

- `xgb.importance(feature_names, xgb.model)`
- `xgb.ggplot.importance(xgb.importance , measure = { 'Gain' ,
'Cover' ,
'Frequency' })`

XGB GAIN - SANTANDER



XGB GAIN - HP



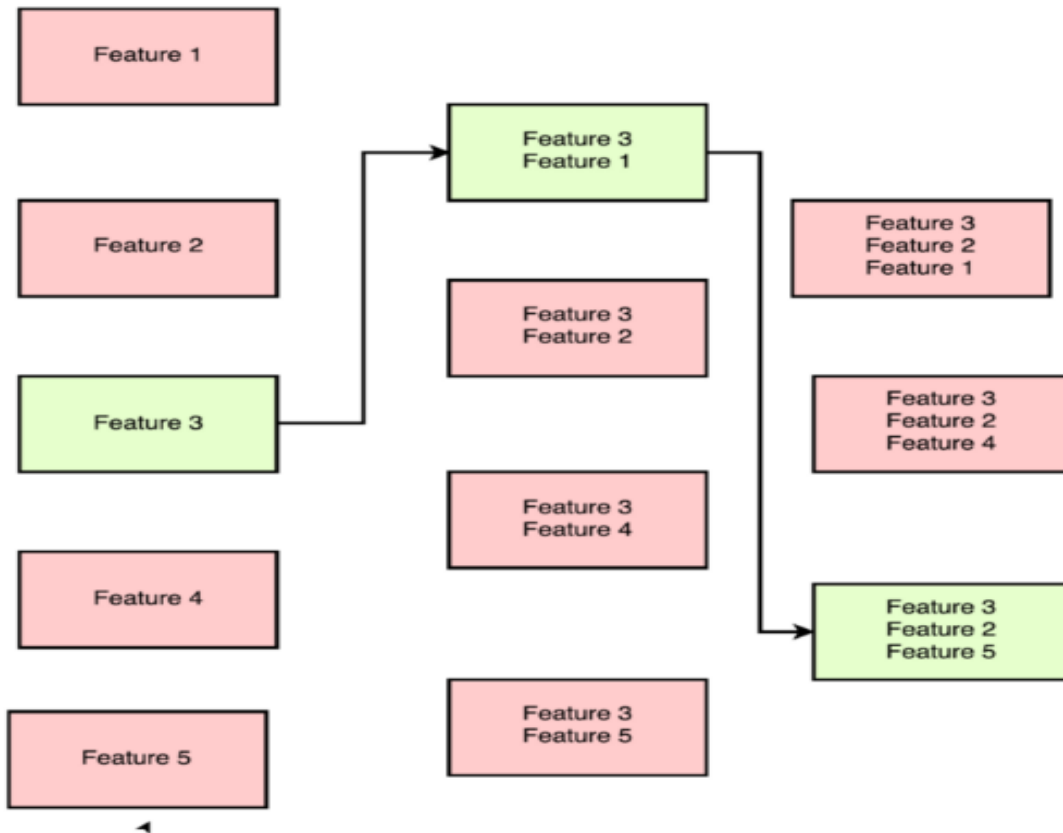
WRAPPER METHOD - CARACTERÍSTICAS

- SE BASAN EN ALGORITMOS DE MACHINE LEARNING
- BUSCAN MEJOR SUBCONJUNTO OPTIMICE MODELO
- COMPUTACIONALMENTE COSTOSOS
- OVERFITTING

WRAPPER METHOD - TÉCNICAS

- **FORDWARD**
- **BACKWARD**
- **GENETIC**
- **BORUTA**

WRAPPER METHOD - FORWARD



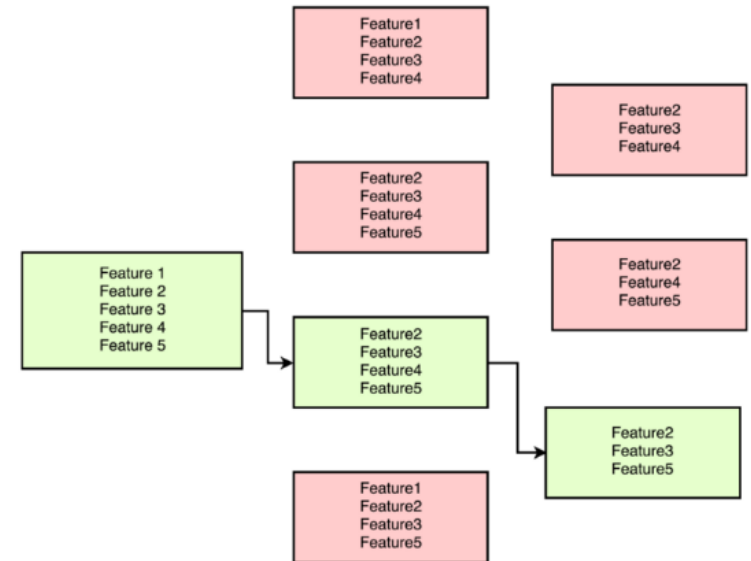
WRAPPER METHOD - BACKWARD

- RECURSIVE ELIMINATION

- CARET \longrightarrow `rfe(x , y , control , sises , metric , params)`

Algoritmo 1: Recursive Feature Elimination

```
1 Entrenar el modelo con todos los atributos;
2 Evaluar el rendimiento del modelo;
3 Ordenar los atributos de acuerdo a una medida de importancia;
4 foreach  $S_i$  ( $i = 1, 2, \dots$ ) do
5   Retener los  $S_i$  atributos más importantes y eliminar los restantes;
6   Entrenar el modelo usando los  $S_i$  atributos;
7   Evaluar el rendimiento del modelo;
8   Recalcular la medida de importancia de los atributos [Opcional];
9 end
10 Ordenar los valores  $S_i$  de acuerdo a su rendimiento;
11 A partir de la lista anterior, elegir un valor  $S_i$  apropiado;
12 Considerar el modelo correspondiente a los  $S_i$  atributos elegidos;
```



WRAPPER METHOD - GENETIC

Algoritmo 3: Algoritmo genético para selección de atributos

```
1 Elegir el criterio de parada, el tamaño de la población ( $m$ ), el número de hijos de cada generación ( $GenSize$ ) y la probabilidad de mutación ( $p_m$ );
2 Hacer  $p$ =Número de atributos;
3 Generar aleatoriamente un conjunto inicial de  $m$  cromosomas binarios, cada uno de longitud  $p$ ;
4 repeat
5   foreach cromosoma do
6     | Entrenar un modelo y calcular el rendimiento del cromosoma;
7   end
8   foreach  $k = 1, \dots, GenSize/2$  do
9     | Seleccionar dos cromosomas basándose en su rendimiento;
10    | Cruzamiento: seleccionar una posición aleatoria en el cromosoma e intercambiar los genes a partir de ese punto;
11    | Mutación: Alterar aleatoriamente, con probabilidad  $p_m$ , los valores binarios de cada gen en cada cromosoma de la descendencia;
12  end
13 until se cumpla el criterio de parada;
```

WRAPPER METHOD – IMPLEMENTACIÓN R

```
selectFeatures( control, learner , task ,resampling , measures = list( ) )
```

```
SelControlExhaustive( max.features )
```

```
FeatSelControlSequential( method = { "sfs", "sbs" } , alpha = 0.02)
```

```
FeatSelControlRandom( )
```

```
FeatSelControlGA( )
```

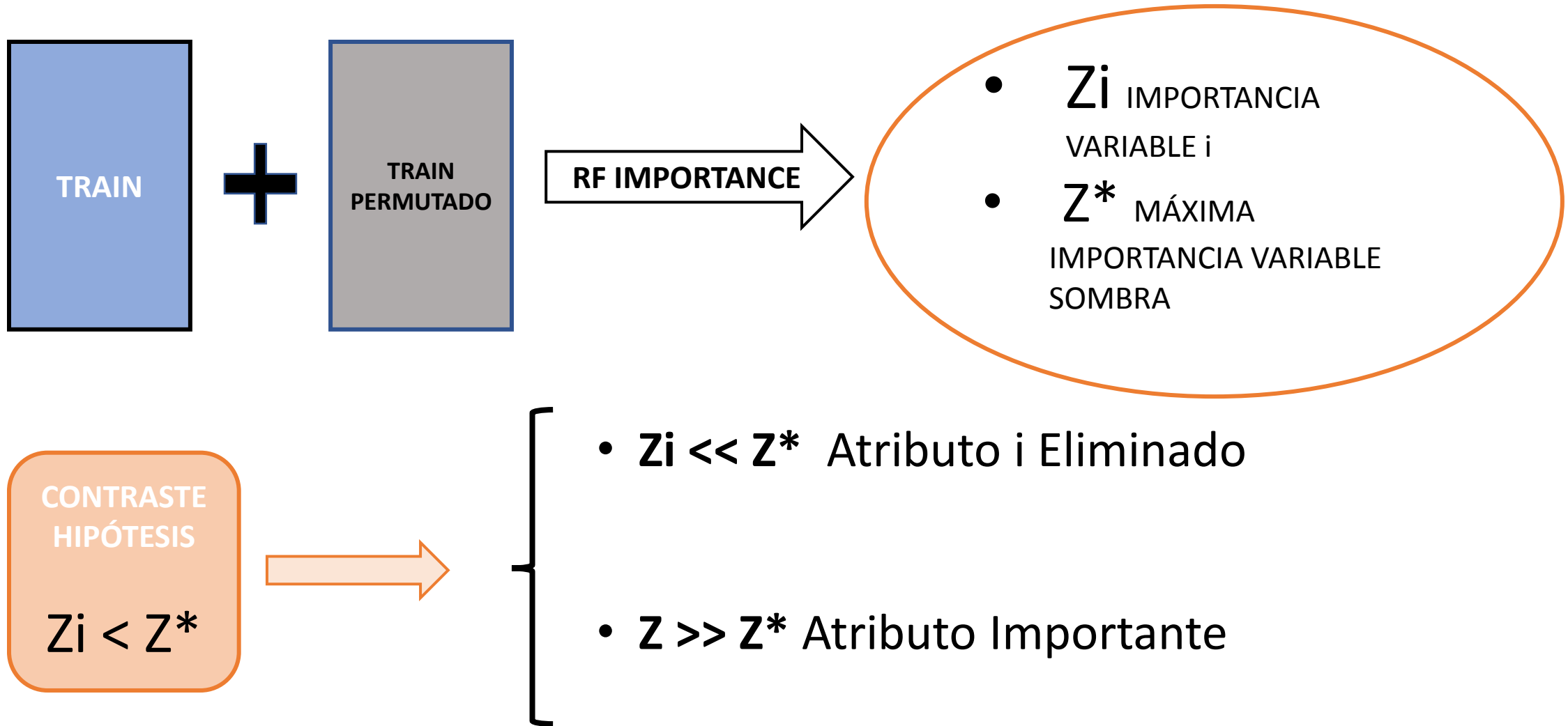
```
makeLearner( cl ="classif.xgboost", parametros )      | listLearners(obj = NA)
```

```
makeResampleDesc( { 'CV' , `LOO` } , iters = 5 , stratify = T)
```

```
makeClassifTask( data , task )
```

[measures: https://mlr.mlr-org.com/articles/tutorial/measures.html](https://mlr.mlr-org.com/articles/tutorial/measures.html)

WRAPPER METHOD - BORUTA



WRAPPER METHOD – BORUTA IMPLEMENTACIÓN

```
Boruta( formula ,  
        data ,  
        pValue = 0.01,  
        mcAdj = TRUE,  
        maxRuns = 100, __ )
```

```
> boruta.trn
```

```
Boruta performed 14 iterations in 17.50179 hours.
```

```
No attributes deemed important.
```

```
No attributes deemed unimportant.
```

```
260 tentative attributes left: random_1, random_10, random_11, random_12, random_13 and 255  
more;
```

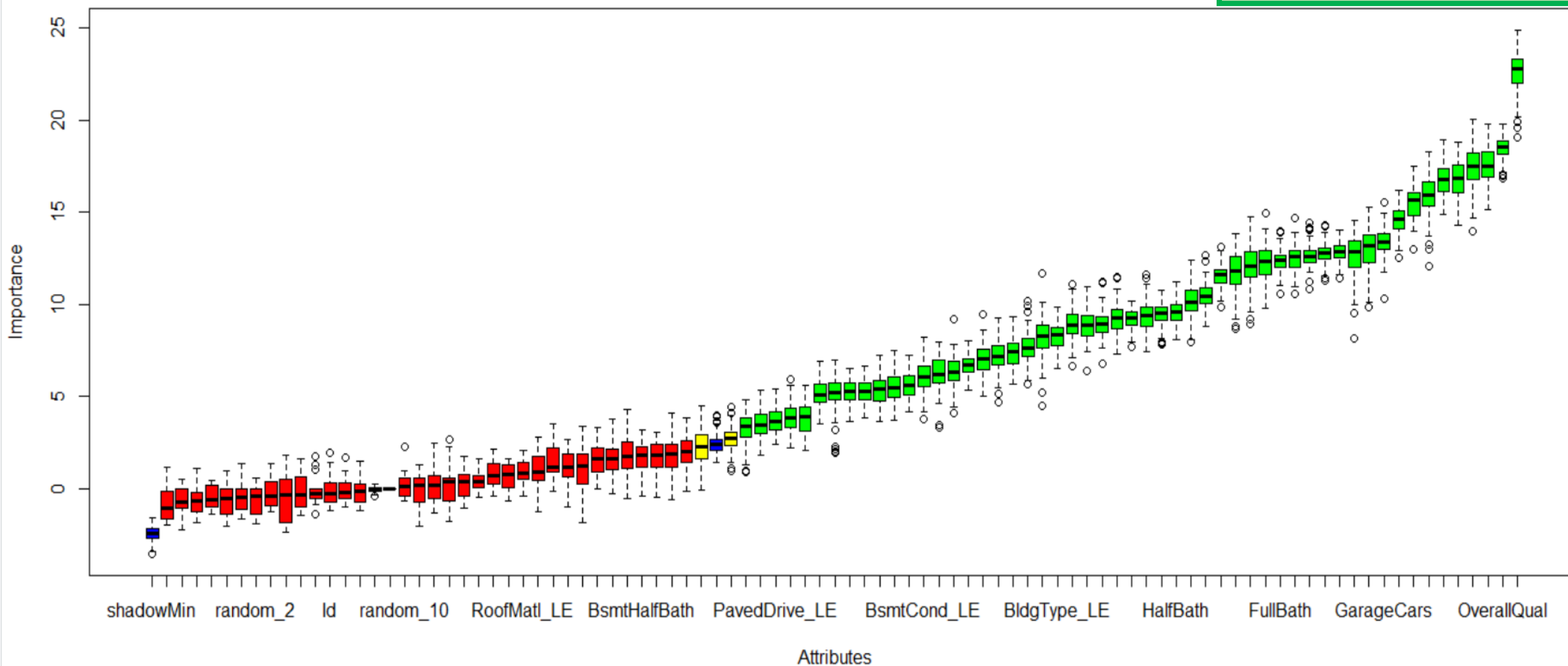
```
> |
```

BORUTA - HOUSEPRICE

55 variables selec.

SCORE VAL: 0.1296

SCORE TRAIN: 0.0785

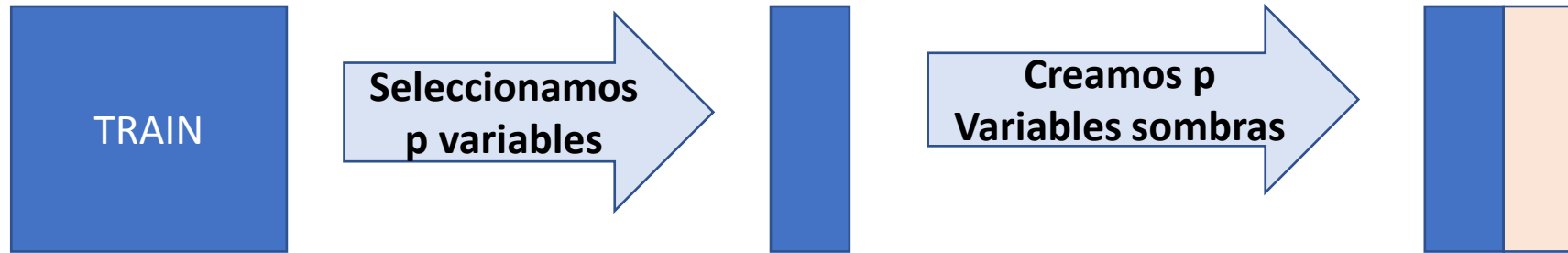


¿Y AHORA...?

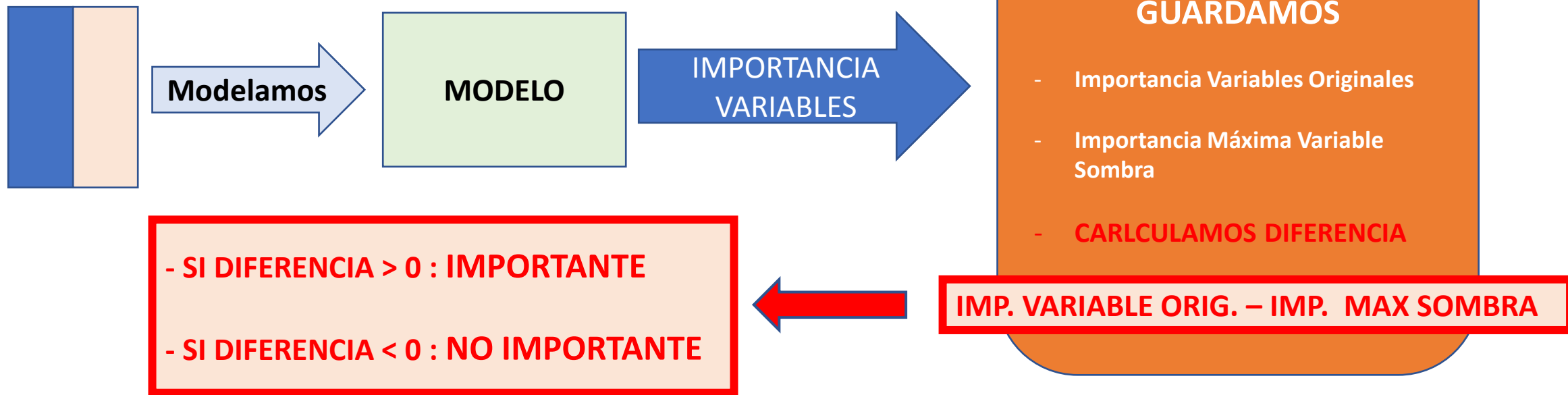


ITERACIÓN i:

1. Seleccionamos p variables y creamos su sombra



2. Creamos modelo simple



FIN:

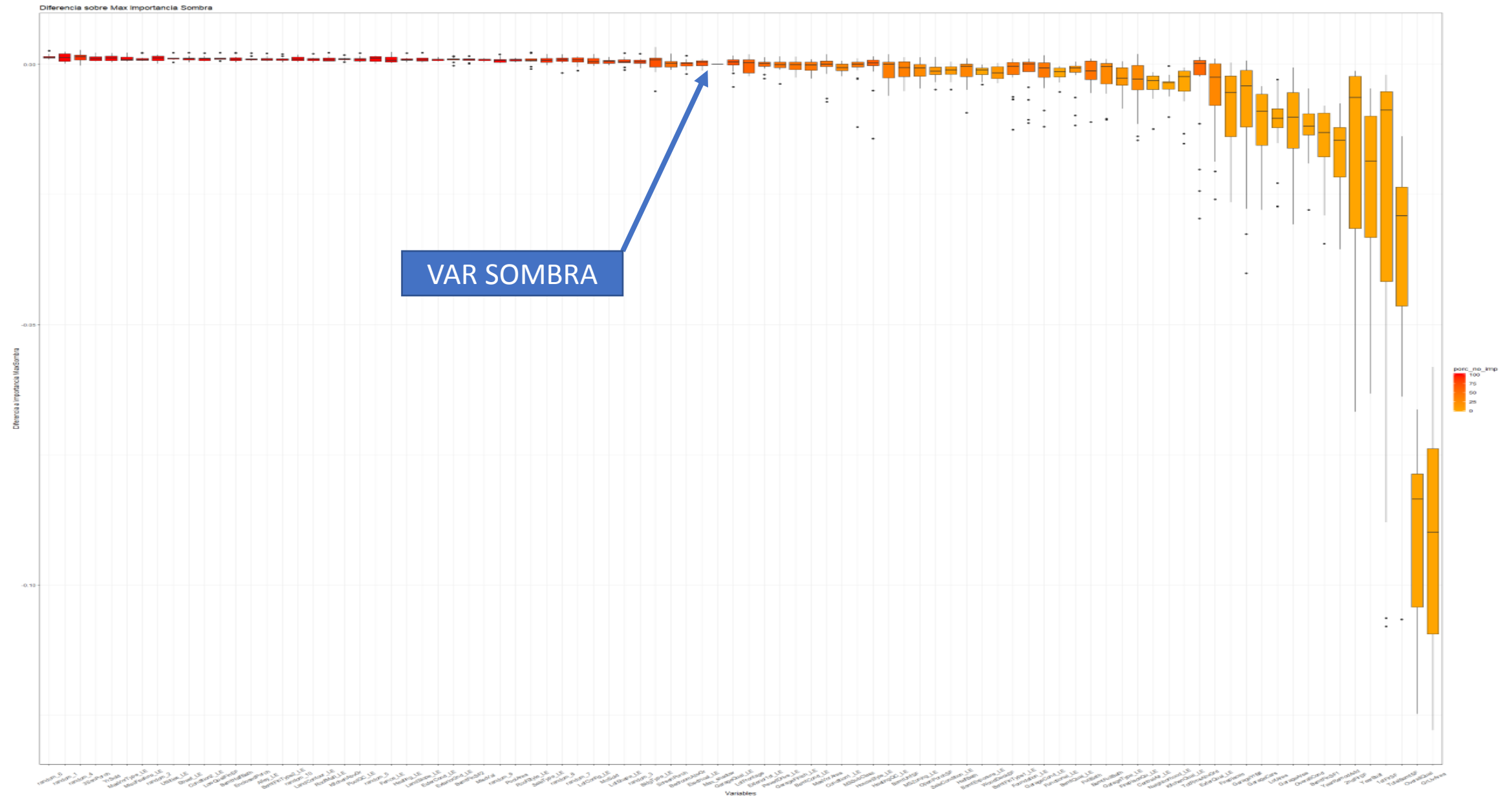
Disponemos de un histórico

FEATURE	NUM ITERACIONES	Num Iteraciones no Importante	% No Importante	Media Diff MAX_SHADOW
random	10	10	100%	-0,04852
Var_81	30	0	0%	3,2455
Ramdom_2	12	10	83%	-0,0254
Noise	30	9	30%	0,05458

TOMAMOS DECISIÓN:

- Elegimos % no importante >> 50%
- Elegimos Media_DIFERENCIA > 0
- ... etc

IMPORTANCIA HP



EJEMPLO NUEVO – SANTANDER DATA

- **CREAMOS VARIABLES RUIDO**

- VAR_81_REP
- VAR_2_REP

- **CREAMOS VARIABLES RANDOM**

- RANDOMs

GAIN XGBOOST ASIGNA **BAJA IMPORTANCIA**



Importancia

