

Supplementary Methods

Discrimination of rosé wines using shotgun metabolomics with a genetic algorithm and MS ion intensity ratios

Mélodie Gil¹, Christelle Reynes², Guillaume Cazals³, Christine Enjalbal³, Robert Sabatier², Cédric Saucier^{1,*}

1- Univ Montpellier, SPO, INRAE, Montpellier Supagro, Montpellier, France

2- Univ Montpellier, IGF, CNRS INSERM, Montpellier, France

3- Univ Montpellier, IBMM, Montpellier, France

*Corresponding author at: Université de Montpellier, UMR 1083 INRAE Supagro Sciences pour l'Œnologie, 15 Avenue Charles Flahault, 34000 Montpellier, France. Tel.: +33411759567; fax: +33411759638. E-mail address: cedric.saucier@umontpellier.fr.

Appendix 1. In-house rosé wine compounds database

Molecule	Molecular formula	(MH) ⁺ (or M ⁺ for anthocyanins)
Organic acids		
Malic acid	C ₄ H ₆ O ₅	135.0924
Acetic acid	C ₂ H ₄ O ₂	61.029
Lactic acid	C ₃ H ₆ O ₃	91.0395
Tartric acid	C ₄ H ₆ O ₆	151.0243
Gluconic acid	C ₁₂ H ₁₂ O ₇	197.0661
Citric acid	C ₆ H ₈ O ₇	193.0348
Benzoic acids and derivatives		
Gallic acid	C ₇ H ₆ O ₅	171.0294
Protocatechuic acid	C ₇ H ₆ O ₄	155.0344
Syringic acid	C ₉ H ₁₀ O ₅	199.0607
Vanillic acid	C ₈ H ₈ O ₄	169.0501
Ethyl gallate	C ₉ H ₁₀ O ₅	199.0606
Ethyl protocatechuate	C ₉ H ₁₀ O ₄	183.0657
Gentisic acid	C ₇ H ₆ O ₄	155.0344
Hydroxybenzoic acid	C ₇ H ₆ O ₃	139.0395
Hydroxycinnamic acids and derivatives		
Caffeic acid	C ₉ H ₈ O ₄	181.0501
p-coumaric acid	C ₉ H ₈ O ₃	165.0552
Ferulic acid	C ₁₀ H ₁₀ O ₄	195.0657
Caftaric acid (cis/trans)	C ₁₃ H ₁₂ O ₉	313.056
Coutaric acid (cis/trans)	C ₁₃ H ₁₂ O ₈	297.0611
Fertaric acid	C ₁₄ H ₁₄ O ₉	327.0716
Ethyl caffeate	C ₁₁ H ₁₂ O ₄	209.0814
Ethyl coumarate	C ₁₁ H ₁₂ O ₃	193.0865
3 Caffeoylquinic acid	C ₁₆ H ₁₆ O ₉	355.1029
Sinapic acid	C ₁₁ H ₁₂ O ₅	225.0763
Ferulic acid ethyl ester	C ₁₂ H ₁₄ O ₄	223.0964
GRP	C ₂₃ H ₂₇ N ₃ O ₁₅ S	618.1242

Stilbenes		
Resveratrol (cis trans)	C ₁₄ H ₁₂ O ₃	229.0865
Viniferin	C ₂₈ H ₂₂ O ₆	455.1495
Piceid (cis trans)	C ₂₀ H ₂₂ O ₈	391.1393
Flavonols		
Quercetin Glc	C ₂₁ H ₂₀ O ₁₂	465.1033
Myricetin Glc	C ₂₁ H ₂₀ O ₁₃	481.0982
Myricetin glucuronide	C ₂₁ H ₁₈ O ₁₄	495.0775
Quercetin glucuronide	C ₂₁ H ₁₈ O ₁₃	479.0826
Quercetin	C ₁₅ H ₁₀ O ₇	303.0505
Kaempferol	C ₁₅ H ₁₀ O ₆	287.0556
Kaempferol Glc	C ₂₁ H ₂₀ O ₁₁	449.1084
Syringetin Glc	C ₂₃ H ₂₄ O ₁₃	509.1295
Isorhamnetin	C ₁₆ H ₁₂ O ₇	317.0661
Syringetin	C ₁₇ H ₁₄ O ₈	347.0766
Flavan-3-ols		
Phloroglucinol	C ₆ H ₆ O ₃	127.0395
Catechin	C ₁₅ H ₁₄ O ₆	291.0869
Epigallocatechin	C ₁₅ H ₁₄ O ₇	307.0818
Catechin gallate	C ₂₂ H ₁₈ O ₁₀	443.0978
Galloylated dimer	C ₃₇ H ₃₀ O ₁₆	731.1612
Trimer-1 (1 isomer)	C ₄₅ H ₃₈ O ₁₈	867.2137
Trimers-2 (several isomers)		
(epi)cat-ethyl-(epi)cat-1 (1 isomer)	C ₃₂ H ₃₀ O ₁₂	607.1816
(epi)cat-ethyl-(epi)cat-2 (2 co-eluted isomers)		
Epicatechin	C ₁₅ H ₁₄ O ₆	291.0869
Dimer B2	C ₃₀ H ₂₆ O ₁₂	579.1503
Dimer B1		
Dimer B3		
Dimer B4		
Dimer 1 epigallocatechin B	C ₃₀ H ₂₆ O ₁₃	595.1452
Dimer 1 catechin gallate B	C ₃₇ H ₃₀ O ₁₆	731.1612
Quadrimer cat epi	C ₆₀ H ₅₀ O ₂₄	1155.277

Dihydroflavonols		
Astilbin	C ₂₁ H ₂₂ O ₁₁	451.124
Taxifolin	C ₁₅ H ₁₂ O ₇	305.0661
Anthocyanins		
Pelargonin	C ₁₅ H ₁₁ O ₅	271.0607
Delphinidin	C ₁₅ H ₁₁ O ₇	303.0505
Cyanidin	C ₁₅ H ₁₁ O ₆	287.0556
Petunidin	C ₁₆ H ₁₃ O ₇	317.0661
Peonidin	C ₁₆ H ₁₃ O ₆	301.0712
Malvidin	C ₁₇ H ₁₅ O ₇	331.0818
Malvidin 3,5-diGlc	C ₂₉ H ₃₅ O ₁₇	655.1874
Delphinidin 3,5-diGlc	C ₂₇ H ₃₁ O ₁₇	627.1561
Cyanidin 3,5-diGlc	C ₂₇ H ₃₁ O ₁₆	611.1612
Petunidin 3,5-diGlc	C ₂₈ H ₃₃ O ₁₇	641.1718
Peonidin 3,5-diGlc	C ₂₈ H ₃₃ O ₁₆	625.1769
Pelargonin O Glc	C ₂₁ H ₂₁ O ₁₀	433.1135
Malvidin 3-O-Glc	C ₂₃ H ₂₅ O ₁₂	493.1346
Delphinidin 3-O-Glc	C ₂₁ H ₂₁ O ₁₂	465.1033
Cyanidin 3-O-Glc	C ₂₁ H ₂₁ O ₁₁	449.1084
Petunidin 3-O-Glc	C ₂₂ H ₂₃ O ₁₂	479.1190
Peonidin 3-O-Glc	C ₂₂ H ₂₃ O ₁₁	463.1240
Delphinidin 3-O-acetyl-Glc	C ₂₃ H ₂₃ O ₁₃	507.1139
Cyanidin 3-O-acetyl-Glc	C ₂₃ H ₂₃ O ₁₂	491.1190
Petunidin 3-O-acetyl-Glc	C ₂₄ H ₂₅ O ₁₃	521.1295
Peonidin 3-O-acetyl-Glc	C ₂₄ H ₂₅ O ₁₂	505.1346
Malvidin 3-O-acetyl-Glc	C ₂₅ H ₂₇ O ₁₃	535.1452
Delphinidin 3-O-coumaroyl-Glc (cis- and trans- isomers)	C ₃₀ H ₂₇ O ₁₄	611.1401
Cyanidin 3-O-coumaroyl-Glc (cis- and trans- isomers)	C ₃₀ H ₂₇ O ₁₃	595.1452
Petunidin 3-O-coumaroyl-Glc (cis- and trans- isomers)	C ₃₁ H ₂₉ O ₁₄	625.1557
Peonidin 3-O-coumaroyl-Glc (cis- and trans- isomers)	C ₃₁ H ₂₉ O ₁₃	609.1608
Malvidin 3-O-coumaroyl-Glc (cis- and trans- isomers)	C ₃₂ H ₃₁ O ₁₄	639.1714
Delphinidin 3-O-caffeoyl-Glc	C ₃₀ H ₂₇ O ₁₅	627.1350
Cyanidin 3-O-caffeoyl-Glc	C ₃₀ H ₂₇ O ₁₄	611.1401

Petunidin 3-O-caffeoyl-Glc	C ₃₁ H ₂₉ O ₁₅	641.1507
Peonidin 3-O-caffeoyl-Glc	C ₃₁ H ₂₉ O ₁₄	625.1557
Malvidin 3-O-caffeoyl-Glc	C ₃₂ H ₃₁ O ₁₅	655.1663
(epi)cat-ethyl-peonidin 3-O-Glc (4 isomers)	C ₃₉ H ₃₉ O ₁₇	779.2187
(epi)cat-ethyl-malvidin 3-O-Glc (4 isomers)	C ₄₀ H ₄₁ O ₁₈	809.2293
(epi)cat-ethyl-malvidin 3-O-coumaroyl-Glc (2 co-eluted isomers)	C ₄₉ H ₄₇ O ₂₀	955.2661
Delphinidin 3-O-Glc-(epi)cat	C ₃₆ H ₃₃ O ₁₈	753.1670
Cyanidin 3-O-Glc-(epi)cat	C ₃₆ H ₃₃ O ₁₇	737.1718
Petunidin 3-O-Glc-(epi)cat	C ₃₇ H ₃₅ O ₁₈	767.1823
Peonidin 3-O-Glc-(epi)cat	C ₃₇ H ₃₅ O ₁₇	751.1874
Malvidin 3-O-Glc-(epi)cat (2 isomers)	C ₃₈ H ₃₇ O ₁₈	781.1980
Malvidin 3-O-coumaroyl-Glc-(epi)cat (2 isomers)	C ₄₇ H ₄₃ O ₂₀	927.2348
(epi)cat-delphinidin 3-O-Glc A-F bicyclic	C ₃₆ H ₃₅ O ₁₈	755.1823
(epi)cat-cyanidin 3-O-Glc A-F bicyclic	C ₃₆ H ₃₅ O ₁₇	739.1874
(epi)cat-petunidin 3-O-Glc A-F bicyclic	C ₃₇ H ₃₇ O ₁₈	769.1980
(epi)cat-peonidin 3-O-Glc A-F bicyclic	C ₃₇ H ₃₇ O ₁₇	753.2031
(epi)cat-malvidin 3-O-Glc A-F bicyclic	C ₃₈ H ₃₉ O ₁₈	783.2137
Caftaric-peonidin 3-O-Glc (2 co-eluted isomers)	C ₃₅ H ₃₃ O ₂₀	773.1565
Caftaric-malvidin 3-O-Glc (2 co-eluted isomers)	C ₃₆ H ₃₅ O ₂₁	803.1671
Coutaric-malvidin 3-O-Glc (2 co-eluted isomers)	C ₃₆ H ₃₅ O ₂₀	789.1722
Pyranodelphinidin 3-O-Glc	C ₂₃ H ₂₁ O ₁₂	489.1033
Pyranocyanidin 3-O-Glc	C ₂₃ H ₂₁ O ₁₁	473.1084
Pyranopetunidin 3-O-Glc	C ₂₄ H ₂₃ O ₁₂	503.119
Pyranopeonidin 3-O-Glc	C ₂₄ H ₂₃ O ₁₁	487.124
Pyranomalvidin 3-O-Glc (vitisin B)	C ₂₅ H ₂₅ O ₁₂	517.1346
Carboxypyranodelphinidin 3-O-Glc	C ₂₄ H ₂₁ O ₁₄	533.0931
Carboxypyranocyanidin 3-O-Glc	C ₂₄ H ₂₁ O ₁₃	517.0982
Carboxypyranopetunidin 3-O-Glc	C ₂₅ H ₂₃ O ₁₄	547.1088
Carboxypyranopeonidin 3-O-Glc	C ₂₅ H ₂₃ O ₁₃	531.1139
Carboxypyranomalvidin 3-O-Glc (vitisin A)	C ₂₆ H ₂₅ O ₁₄	561.1244
Pyranodelphinidin 3-O-acetyl-Glc	C ₂₅ H ₂₃ O ₁₃	531.1139
Pyranocyanidin 3-O-acetyl-Glc	C ₂₅ H ₂₃ O ₁₂	515.119
Pyranopetunidin 3-O-acetyl-Glc	C ₂₆ H ₂₅ O ₁₃	545.1295

Pyranopeonidin 3-O-acetyl-Glc	C ₂₆ H ₂₅ O ₁₂	529.1346
Pyranomalvidin 3-O-acetyl-Glc	C ₂₇ H ₂₇ O ₁₃	559.1452
Carboxypyranodelphinidin 3-O-acetyl-Glc	C ₂₆ H ₂₃ O ₁₅	575.1037
Carboxypyranocyanidin 3-O-acetyl-Glc	C ₂₆ H ₂₃ O ₁₄	559.1088
Carboxypyranopetunidin 3-O-acetyl-Glc	C ₂₇ H ₂₅ O ₁₅	589.1194
Carboxypyranopeonidin 3-O-acetyl-Glc	C ₂₇ H ₂₅ O ₁₄	573.1244
Carboxypyranomalvidin 3-O-acetyl-Glc	C ₂₈ H ₂₇ O ₁₅	603.135
Pyranopeonidin 3-O-coumaroyl-Glc	C ₃₃ H ₂₉ O ₁₃	633.1608
Pyranomalvidin 3-O-coumaroyl-Glc	C ₃₄ H ₃₁ O ₁₄	663.1714
Carboxypyranopetunidin 3-O-coumaroyl-Glc	C ₃₄ H ₂₉ O ₁₆	693.1456
Carboxypyranopeonidin 3-O-coumaroyl-Glc	C ₃₄ H ₂₉ O ₁₅	677.1507
Carboxypyranomalvidin 3-O-coumaroyl-Glc	C ₃₅ H ₃₁ O ₁₆	707.1612
p-hydroxyphenylpyranopeonidin 3-O-Glc	C ₃₀ H ₂₇ O ₁₂	579.1503
p-hydroxyphenylpyranomalvidin 3-O-Glc	C ₃₁ H ₂₉ O ₁₃	609.1608
p-hydroxyphenylpyranopeonidin 3-O-acetyl-Glc	C ₃₂ H ₂₉ O ₁₃	621.1608
p-hydroxyphenylpyranomalvidin 3-O-acetyl-Glc	C ₃₃ H ₃₁ O ₁₄	651.1714
p-hydroxyphenylpyranopeonidin 3-O-coumaroyl-Glc	C ₃₉ H ₃₃ O ₁₄	725.187
p-hydroxyphenylpyranomalvidin 3-O-coumaroyl-Glc	C ₄₀ H ₃₅ O ₁₅	755.1976
Catechylpyranopeonidin 3-O-Glc	C ₃₀ H ₂₇ O ₁₃	595.1452
Catechylpyranomalvidin 3-O-Glc (pinotin A)	C ₃₁ H ₂₉ O ₁₄	625.1557
Catechylpyranopetunidin 3-O-acetyl-Glc	C ₃₂ H ₂₉ O ₁₅	653.1507
Catechylpyranopeonidin 3-O-acetyl-Glc	C ₃₂ H ₂₉ O ₁₄	637.1557
Catechylpyranomalvidin 3-O-acetyl-Glc	C ₃₃ H ₃₁ O ₁₅	667.1663
Catechylpyranopetunidin 3-O-coumaroyl-Glc	C ₃₉ H ₃₃ O ₁₆	757.1769
Catechylpyranopeonidin 3-O-coumaroyl-Glc	C ₃₉ H ₃₃ O ₁₅	741.182
Catechylpyranomalvidin 3-O-coumaroyl-Glc	C ₄₀ H ₃₅ O ₁₆	771.1925
Guaiacylpyranomalvidin 3-O-Glc	C ₃₂ H ₃₁ O ₁₄	639.1714
Guaiacylpyranomalvidin 3-O-acetyl-Glc	C ₃₄ H ₃₃ O ₁₅	681.182
Guaiacylpyranomalvidin 3-O-coumaroyl-Glc	C ₄₁ H ₃₇ O ₁₆	785.2082
Syringylpyranomalvidin 3-O-Glc	C ₃₃ H ₃₃ O ₁₅	669.182
Petunidin 3-O-acetyl-Glc-(epi)cat	C ₃₉ H ₃₇ O ₁₉	809.1929
Peonidin 3-O-acetyl-Glc-(epi)cat	C ₃₉ H ₃₇ O ₁₈	793.198
Malvidin 3-O-acetyl-Glc-(epi)cat	C ₄₀ H ₃₉ O ₁₉	823.2086

Pyranopeonidin 3-O-Glc-(epi)cat	C ₃₉ H ₃₅ O ₁₇	775.1874
Pyranomalvidin 3-O-Glc-(epi)cat	C ₄₀ H ₃₇ O ₁₈	805.198
Pyranomalvidin 3-O-coumaroyl-Glc-(epi)cat	C ₄₉ H ₄₃ O ₂₀	951.2348
Others (aminoacids and alcohols)		
Tyrosine	C ₉ H ₁₁ NO ₃	182.0817
Tyrosol	C ₈ H ₁₀ O ₂	121.0653
Hydroxytyrosol	C ₈ H ₁₀ O ₃	155.0708
Tryptophol	C ₁₀ H ₁₁ NO	162.0919
Tryptophan	C ₁₁ H ₁₂ N ₂ O ₂	203.0821
Oxidation product		
GSH	C ₁₀ H ₁₇ N ₃ O ₆ S	308.0916
GSSG	C ₂₀ H ₃₂ N ₆ O ₁₂ S ₂	613.1598
2-S-glutathionylcaftaric acid ethyl ester	C ₂₅ H ₃₂ N ₃ O ₁₅ S	646.1554
2-S-glutathionylcaffeic acid	C ₁₉ H ₂₃ O ₁₀ N ₃ S	486.1182

Appendix 2. Genetic algorithms

Genetic Algorithms (GAs) are inspired by nature and especially by natural selection (Goldberg, 1989), they are very useful in complex optimization issues. Here, the GA are used to find up optimal subsets of peak ratios. Thus, it can be considered as a feature selection issue, a field where GAs are widely used (Chandrashekar & Sahin, 2014), especially in metabolomics (Cavill *et al.*, 2009). The algorithm begins with a population constituted of several individuals which correspond to potential solutions in the optimization problem. Thus, in our context, the individuals will be subsets of peak ratios. Then, this population evolves according to three operators described in the next paragraphs: crossover, mutation and selection. Selection is a crucial step allowing to keep the best individuals with regard to the fitness function. Mutation and crossover are run independently from the optimization issue, they allow the exploration of solution space. Then, any GA can be described as follows:

Main steps of a GA:

1. Construction of the first generation
2. Selection

While stopping criteria not met do

3. Crossover
4. Mutation
5. Selection

End

Peak weighting

As the total number of potential solutions is huge (see main manuscript) we decided to associate a weight to each peak in order to favor peaks with two main properties: high values (that is to say high quantities) in at least one group and low within-group variability.

To favor peaks with high values in at least one group, we computed the mean intensities within each group:

$$\bar{x}_j^k = \frac{1}{n_k} \sum_{i \in k} x_{ij}$$

Where x_{ij} is the intensity value of the i -th sample for the j -th ($i \in \{1,2, \dots, 72\}$) peak and k ($k \in \{1,2,3\}$) is the class (here the region). Then we extracted the maximum mean across classes:

$$\bar{x}_j^{max} = \max_{k \in \{1,2,3\}} \bar{x}_j^k$$

Concerning variability, the coefficient of variation of each peak was computed within each group and the average value between groups was kept. In order to introduce a balance between those two criteria, each one was linearly normalized between 0 and 1 across the data. The average value of both normalized criteria was used as the final weight.

Fitness

Genetic algorithms are optimization methods. Hence, a criterion, called fitness, has to be defined in order to quantify the solution quality and to be optimized. As described in the main manuscript, a two-fold cross-validated (2-FCV) accuracy obtained with Linear Discriminant Analysis (LDA) is used to quantify the ability of a subset of peak ratios to discriminate between wine origins (denoted classes in the following). However, in order to avoid overfitting, we favor solutions including few peak ratios. Hence, a penalty term was added to the 2-FCV good classification rate:

$$fitness(S) = accuracy_{CV}(S) - 0.05 \times \frac{card(S) - C_{max}}{1 - C_{max}}$$

Where S denotes a subset of ratios of size $card(S)$ and C_{max} is the maximum number of ratios the user wants to put in a solution. The use of C_{max} only aims at hastening convergence and can be easily changed if the optimal solution reaches the boundary (in our applications $C_{max} = 10$). $accuracy_{CV}(S)$ is the average accuracy obtained by applying LDA to the subset S for 30 runs of 2-fold cross-validation

Initialization

The first population is randomly drawn in order to explore the set of potential solutions while favoring peaks with good inner properties (high intensities and low variability as described earlier). Then, a number of peak ratios is randomly drawn between 1 and C_{max} , and the corresponding number of peaks is chosen following the weights defined earlier. This process is repeated T_{pop} times (in our applications $T_{pop}=200$).

Selection

This step, based on the fitness values, is defined as in Reeves and Rowe (2003). The individuals are ranked according to their fitness value, the best one having the highest rank. Then, the probability to keep one solution in the next generation is proportional to its rank, so that the sum of probabilities over individuals sums to 1 and that the best individual is twice as likely to be selected as the individual with median rank. Hence, the selection probability of the k -th ranked solution is defined by:

$$p_{sel}(S_{(k)}) = 2 \times \frac{(T_{pop} - 2M) \times k + 1}{T_{pop}(3T_{pop} - 4M + 1)}$$

Where M is median of all ranks.

Using those probabilities, $T_{pop} - 1$ solutions are selected. In addition, elitism is applied, that is, the best solution of each generation is automatically introduced in the next one. We obtain a new population of size T_{pop} .

Crossover

The objective of this step is to produce new combinations of the previously retained ratios. It is performed through a usual single-point crossover (Reeves and Rowe, 2003) on ratios and applied to 50% of the individuals of a generation.

Mutation

This step brings the necessary hazard to efficiently explore the solution space. It assures that any point of this space can be reached within a finite number of generations. It is applied to 90% of the individuals of a generation. For each individual to be mutated, three kinds of mutations are possible:

- one of the peak in randomly replaced by another one (still using the weights previously defined) with a probability of 50% ;
- one of the ratios inside the solution is randomly chosen and removed with a probability of 25% ;
- one new ratios is added to the solution as performed in the initialization step with a probability of 25%.

Convergence and stopping criterion

As theoretically proved in Bhandari et al. (1996), two conditions are necessary and sufficient for GA to converge as the number of iterations goes to infinity:

- The best solution in the present population has a fitness value no less than the fitness values of the optimal strings from the previous populations (verified through our selection step).
- Each solution has a positive probability of going to an optimal string within a finite number of iterations (verified through our mutation step).

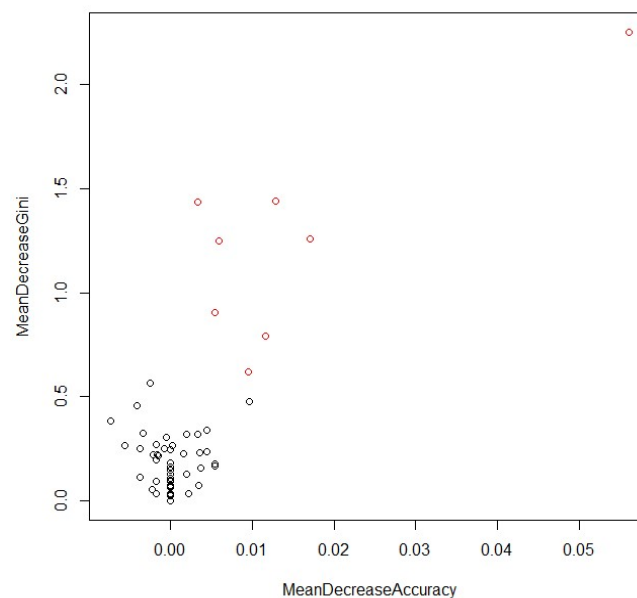
The designed operators meet those conditions but an infinite number of generations is not practicable so the maximum number of generations is set to 2000. It has been set by studying the evolution of fitness function for several runs and by choosing a number at least twice as large as the number of generations required to reach criterion stability.

Appendix 3. Random Forests

The Random Forest (RF) analysis was performed using the randomForest R package.

Use of importance parameters

In this package two measurements of variable importance are provided (Breiman, 2001): the mean decrease of accuracy (loss of accuracy when a variable is eliminated from trees) and the mean decrease of Gini Index (the same for the Gini Index which is used to choose the best variable and threshold value at each split). They were both used to choose the final set of selected ratios. For example, in the following figure (application to all the ratios), 8 ratios (red points) were selected as having high values of both indicators.



Bibliography

Bhandari, D., Murthy, C.A., Pal, S.K., 1996. Genetic algorithm with elitist model and its convergence. *International Journal of Pattern Recognition and Artificial Intelligence* 10 (6), 731–747.

Reeves, C.R., Rowe, J.E., 2003. *Genetic Algorithms—Principles and Perspectives, A Guide to GA Theory*. Kluwer Academic Publishers, London.

Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.