

Towards a Modular Decision Support System for Radiomics: A Case Study on Rectal Cancer

Roberto Gatta^a, Mauro Vallati^b, Nicola Dinapoli^c, Carlotta Masciocchi^a,
Jacopo Lenkowicz^a, Davide Cusumano^c, Calogero Casá^a, Alessandra
Farchione^d, Andrea Damiani^a, Johan van Soest^e, Andre Dekker^e, Vincenzo
Valentini^c

^a*Istituto di Radiologia, Università Cattolica del Sacro Cuore. Largo F.Vito 1, 00168, Roma, Italia*

^b*School of Computing and Engineering, University of Huddersfield, HD1 3DH Huddersfield, UK*

^c*Polo Scienze Oncologiche ed Ematologiche, Fondazione Policlinico Universitario Agostino Gemelli. Largo A.Gemelli, 8, 00168 Roma, Italia*

^d*Polo Scienze radiologiche e di laboratorio, Fondazione Policlinico Universitario Agostino Gemelli. Largo A.Gemelli 8, 00168, Roma, Italia*

^e*Department of Radiation Oncology (MAASTRO), GROW School for Oncology and Developmental Biology, Maastricht University Medical Centre, Netherlands*

Abstract

Following the personalized medicine paradigm, there is a growing interest in medical agents capable of predicting the effect of therapies on patients, by exploiting the amount of data that is now available for each patient. In disciplines like oncology, where images and scans are available, the exploitation of medical images can provide an additional source of potentially useful information. The study and analysis of features extracted by medical images, exploited for predictive purposes, is termed *radiomics*. A number of tools are available for supporting some of the steps of the radiomics process, but there is a lack of approaches which are able to deal with all the steps of the process.

In this paper, we introduce a medical agent-based decision support system capable of handling the whole radiomics process. The proposed system is tested on two independent data sets of patients treated for rectal cancer. Experimental results indicate that the system is able to generate highly performant centre-specific predictive model, and show the issues related to differences in data sets collected by different centres, and how such issues can

affect the performance of the generated predictive models.

Keywords: Decision Support Systems, Radiomics, Predictive Models, Image Feature Analysis

1. Introduction

Personalized medicine is a relatively new, but already well-established, paradigm based on the principle that each individual is born with unique biological and genetic characteristics [1, 2]. The foundation of this paradigm is formed by disciplines such as Genomics –the science of studying the genes in a genome and their interactions with each other–, and proteomics –which instead focuses on proteins. Furthermore, in disciplines like oncology, where images and scans are available, the exploitation of medical images can provide an additional source of potentially useful information. Thanks also to the recent advances in computer science, it is now possible to extract a huge number of “quantitative“ features from tomographic images (computed tomography [CT], magnetic resonance [MR], or positron emission tomography [PET] images), and such extracted features can then be automatically analysed in order to investigate their informativeness with regards to the evolution of the disease, or the response of the patient to a specific clinical treatment. This discipline is commonly termed *radiomics* [3], and is aimed at providing effective decision support to physicians and practitioners, and complementing the traditional “qualitative” analysis of images, commonly performed by human experts [4]. In this context, features represent a numerical synthesis of some properties of the considered image, which would not be possible to manually extract and analyse. Extracted features can then be combined with available clinical data into complex models to predict patient prognosis or benefit from a specific therapy.

Remarkably, evidence that radiomics can be helpful for predicting tumour control or clinical complications has been documented for most of the common modalities (CT, MRI, PET, etc.) and anatomical districts –such as lung, rectum, or brain– [5, 6, 7, 8].

The growing interest in radiomics lead to the development of several specifically-designed tools; examples include cGITA [9], TexRAD [10, 11], moddicom [12], Pyradiomics [13], and CERR [14]. In parallel with the grows of radiomics tools, initiative such as the Image Biomarker Standardisation

32 Initiative [15] and the Radiomics Ontology¹ become important to standardise
33 the different aspects of image processing and features extraction.

34 However, despite the growing number of radiomics tools, to the best of
35 our knowledge there is a lack of agents that can deal with all the steps of the
36 radiomics process, thus providing a complete and modular environment for
37 supporting the generation and analysis of the predictive models, and allow-
38 ing the exploitation of models in everyday medical routine. Existing tools
39 are mainly aimed at facilitating the extraction of features, and at extend-
40 ing the set of features that can be extracted from a medical image, only.
41 Pivotal steps, like feature selection and generation of the actual predictive
42 model –either via traditional statistical approaches or recent machine learn-
43 ing techniques– are ignored. More worryingly, existing tools do not natively
44 provide any support for the external validation of generated models. As a re-
45 sult, a crucial issue of the exploitation of radiomics predictive models, is that
46 their portability between centres or hospitals is unclear. This is also due to
47 the fact that different machines, particularly in MRI, provide medical images
48 with very different characteristics, particularly in terms of visual noise. Such
49 differences can strongly affect the predictive capabilities of the generated
50 models, and invalidate the results. A possible way for tackling this issue is
51 to extensively exploit external independent testing sets, and providing tools
52 that are supportive in this regards, in order to validate the generated models
53 [16, 17, 18]. Similarly, features can be analysed in order to identify those
54 which are more robust to common sources of image noise. However, despite
55 the fact that empirical investigations which rely on external validation are
56 deemed to be qualitative better than others by the TRIPOD guidelines [19],
57 this approach can not guarantee the reproducibility of the observed results
58 in every set [20].

59 The contribution of this paper is twofold. First, we introduce an approach
60 –under the form of a medical agent-based decision support system– for sup-
61 porting the whole radiomics process. In its current implementation, the agent
62 incorporates some of the ideas and functionalities of moddicom [12]. Given
63 a set of medical images, the proposed system is able to extract a wide range
64 of features, to analyse and select them with regards to the outcome to pre-
65 dict, and to generate an optimised predictive model. When data from a new
66 patient is provided as input, the proposed agent is able to collect features

¹<https://bioportal.bioontology.org/ontologies/RO>

67 from available medical images and patient’s data, and to return a predic-
68 tion about the clinical outcome of a proposed treatment. In other words,
69 the agent can be provided high level goals to achieve –such as, generate a
70 predictive model that shows some given properties–, and it able to reason
71 upon available knowledge in order to satisfy, whether possible, the goals.
72 This reduces the burden on human experts, and provides a valuable decision
73 support tool, that can also allow to investigate alternative approaches and
74 models. The agent is centre-specific, but has been designed in order to be
75 capable of exchanging models between agents in different centres and testing
76 generated models on different data, thus supporting external validation. As
77 a second contribution, we investigate the capabilities of the proposed agent
78 in a real-world scenario. We consider two sets of medical images acquired by
79 two different centres for treating patients affected by rectal cancer. By train-
80 ing the system on each set, we empirically demonstrate how the differences
81 in data sets collected by different centres can affect the performance of the
82 generated predictive models.

83 The remainder of this paper is organised as follows. Section 2 describes
84 the structure of the proposed system and gives details of the considered
85 features. In Section 3 the data sets are introduced, and empirical results are
86 then presented. An extensive discussion is provided in Section 4. Finally,
87 conclusions are given.

88 **2. The Proposed Agent**

89 The architecture of the proposed agent-based decision support system is
90 depicted in Figure 1, and the corresponding software is available at <https://github.com/robertogattabs/RadAgent>. The exploitation of an agent-
91 based approach has a number of advantages. An agent can cope with high
92 level goals, such as generate models that maximise given metrics, by taking
93 into account all the steps of the process. In facts, the agent can reason upon
94 overall and step-specific knowledge in order to modify the behaviour of the
95 corresponding modules, so that the overall goals are achieved. The modular
96 architecture supports the agent by (i) allowing the development and exploita-
97 tion of off-the-shelf modules that can be substituted without any modifica-
98 tion to the rest of the architecture; (ii) providing a standardised interface
99 between the modules; and (iii) allow to modify parameters and behaviour of
100 each component, and isolating the effects on the overall performance.
101

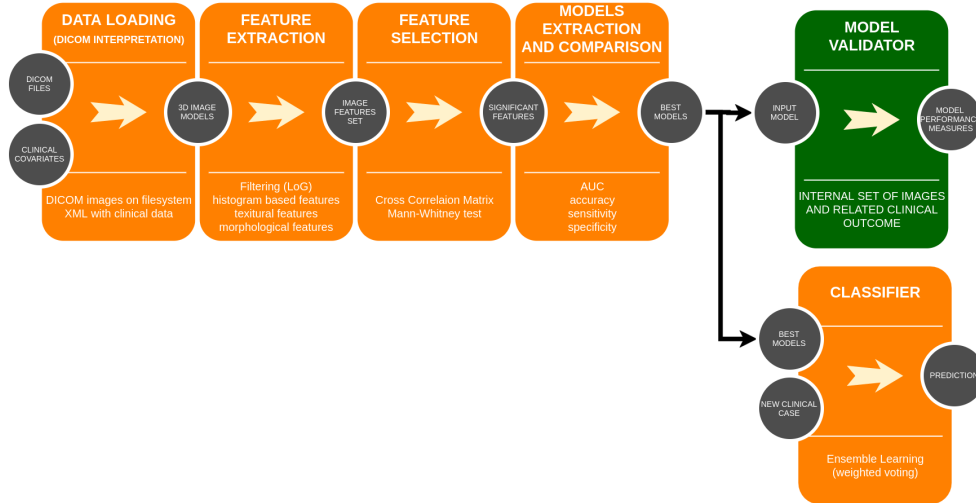


Figure 1: The overall architecture of the proposed decision support system for radiomics, in terms of modules and input/output. The modules included in the architecture correspond to the steps to be performed in the radiomics process. Generated models can be internally validated (green module) or exploited for predicting the outcome of a previously unseen clinical case, and support the physicians.

102 The system has been designed in order to being able to deal with all the
 103 steps of a radiomics analysis, and to provide useful information and support
 104 to the physicians. It is worth emphasising that clinicians are generally not
 105 very keen to exploit predictive models that can not be inspected and validated
 106 “clinically”. For this reason, in the rest of this paper we focus on machine
 107 learning techniques, such as logic regression or decision trees, that allows to
 108 generate predictive models that can be analysed by human experts –but that
 109 yet can provide reasonably high performance.

110 Main functionalities of the proposed system, that can be performed au-
 111 tomatically or required by the users, include:

- 112 • Features extraction from both original medical images, and images fil-
 113 tered using the well-known Laplacian over Gaussian convolution kernel
 114 (LoG) [16]. The LoG filter is commonly exploited in order to smooth
 115 the high frequency noise and enhance the variation of values among
 116 adjacent pixels in the images.
- 117 • The LoG filter can return images with very different appearance ac-
 118 cording the value of the σ parameter used. It is therefore important

119 to identify the σ values that lead to most significant and informative
120 features being extracted for the considered outcome to be predicted.
121 σ values are selected by a Mann-Whitney test with the clinical out-
122 come (to identify the most representative σ) and exploiting a cross-
123 correlation matrix, assessed via a p Pearson Test (to allow the use of
124 two different σ for the same feature in case of no-correlation between
125 the feature at the two σ values).

- 126 • Signatures of selected features, i.e. subsets of informative features, are
127 evaluated. Signatures are generated via greedy forward selection, and
128 are assessed according to metrics provided by the user. In our analysis
129 we considered the AUC with regards to a logistic regression.
- 130 • Signatures are then exploited for generating predictive models, and
131 are compared with regards to their predictive ability. The signature
132 that leads to the best predictive model is selected in order to be used
133 by the agent to support the decisions of the human expert, on new
134 testing cases. Optionally, performance and characteristics of the best
135 performing features can also be presented to the user, which can decide
136 to exploit a different set of features than those included in the signature
137 identified by the agent.

138 Noteworthy, the introduced agent has a high level of configurability, that
139 allows it to be optimised according to the characteristics of the images and
140 data sources of the centre. In order to maximise its compatibility with ex-
141 isting systems, it has been developed in R, which is one of the most used
142 environments for statistical analysis in medicine. Results of the analysis can
143 be exchanged between agents, in order to (externally) validate results or
144 evolve the generated predictive models.

145 The architecture is agnostic with regards to the element to be predicted
146 and to the available features. For the purposes of this work, we consider 92
147 types of features, that can be classified as follows:

- 148 • **BASIC**: first order image features [15] extracted by considering as-
149 pects such as Morphological (MRF), Statistical (STAT), and Intensity
150 Histogram (HI) of the image. Features in this set also include shape
151 properties, such as Volume, Surface, Surface to volume ratio, Com-
152 pactness, Sphericity, Centre of mass shift, Mean, Variance, Skewness,
153 Kurtosis, etc.;

- 154 • **GLCM**: Grey level co-occurrence based textural features [15]. Features
155 in this set include Mean, Variance, Skewness, Kurtosis, 10th and 90th
156 percentile, Robust mean absolute deviation, Energy, etc.;
- 157 • **GLRLM**: Grey level run length based textural features [15]. Examples
158 of features in this set are Short and Long runs emphasis, Short and Long
159 run low grey level emphasis, Grey level non-uniformity normalised, Run
160 entropy, etc.
- 161 • **GLDZM**: Grey level size zone based textural features [15]. Features
162 include Grey level non-uniformity, Zone size non-uniformity, Zone per-
163 centage, Zone size entropy, etc.

164 It should be noted that the value of a feature also depends on the con-
165 sidered σ used in the LoG filter. In this implementation of the system, for
166 the sake of efficiency, we consider 9 possible σ values: 0.35, 0.49, 0.54, 0.59,
167 0.64, 0.69, 0.74, 0.79, 0.84. Such values has been selected according to the
168 experimental results achieved in [21]. The use of 9 possible σ values leads
169 to a grand total of 734 features (Morphological and shape features, from the
170 Basic set, are not affected by changes in the LoG filter) considered by the
171 approach for generating the predictive model. The complete list of features
172 considered in this work is provided in appendix. For a detailed description
173 of the features, including the actual mathematical formulas, the interested
174 reader is referred to [15].

175 3. Experimental Analysis

176 The main purpose of this experimental analysis is to assess the useful-
177 ness of the proposed radiomics agent in supporting the different steps of a
178 radiomics investigation. It is therefore beyond the scope of this study to thor-
179 oughly compare the performance of differently generated predictive models.
180 For a clinical evaluation of mathematical predictive models, the interested
181 reader is referred to [16].

182 The experimental analysis considers two data sets of T2-weighted fast
183 spin-echo 2D oblique images MR scans, that are used for treating patients
184 affected by rectal cancer. The first data set includes scans of 173 patients
185 from the Gemelli polyclinic hospital in Rome, the second set is composed by
186 25 clinical cases treated at the Maastricht clinic of the Maastricht university
187 medical centre. Both the data sets of images include manual contouring

188 of the clinical target volume (CTV) [22]. CTV includes the gross tumour
189 volume, which is the region already affected by the tumour, as well as the
190 regions of direct, local subclinical spread of disease that must be treated in
191 order to stop the evolution of the tumour.

192 The different size of the sets provides an interesting test-bed for a ra-
193 diomics decision support system. Typical medical sets can show a significant
194 size variability, according to the typology of tumour considered and to the
195 characteristics of the medical centre.

196 The scanner used at the Gemelli polyclinic hospital is a MR 1.5 T unit
197 (Signa Excite GE Medical Systems), while the Maastricht clinic is equipped
198 with a Siemens Magnetom AVANTO machine. Acquisition parameters were
199 homogeneous for the two data sets, and are as follows:

- 200 • repetition time, 2500–5000 msec;
- 201 • inversion time, 100–110 msec;
- 202 • pixel spacing, ca. 0.7 mm;
- 203 • echo train length, 16–24;
- 204 • section thickness, 3 mm;
- 205 • no intersection gap;

206 Images have been acquired in a transverse plane orthogonal to the tu-
207 mour longitudinal axis. No intravenous contrast medium was administered.
208 The subsequent manual contouring was performed by an expert radiation
209 oncologist, using a radiotherapy delineation console (Eclipse, Varian Medical
210 System) for the definition of lesion outline as defined in ICRU n. 83.²

211 Figure 2 shows two MR slices from the Gemelli polyclinic data set (left),
212 and two MR slices acquired at the Maastricht clinic. Noteworthy, despite the
213 strict observance of acquisition procedures and acquisition parameters, it is
214 easy to notice that acquired images are significantly different. Qualitatively,
215 images acquired at the Gemelli polyclinic include more visual noise, particu-
216 larly at high frequencies, than those acquired by the other centre. Moreover,

²<https://icru.org/testing/reports/prescribing-recording-and-reporting-intensity-modulated-photon-beam-therapy-imrt-icru-report-83>

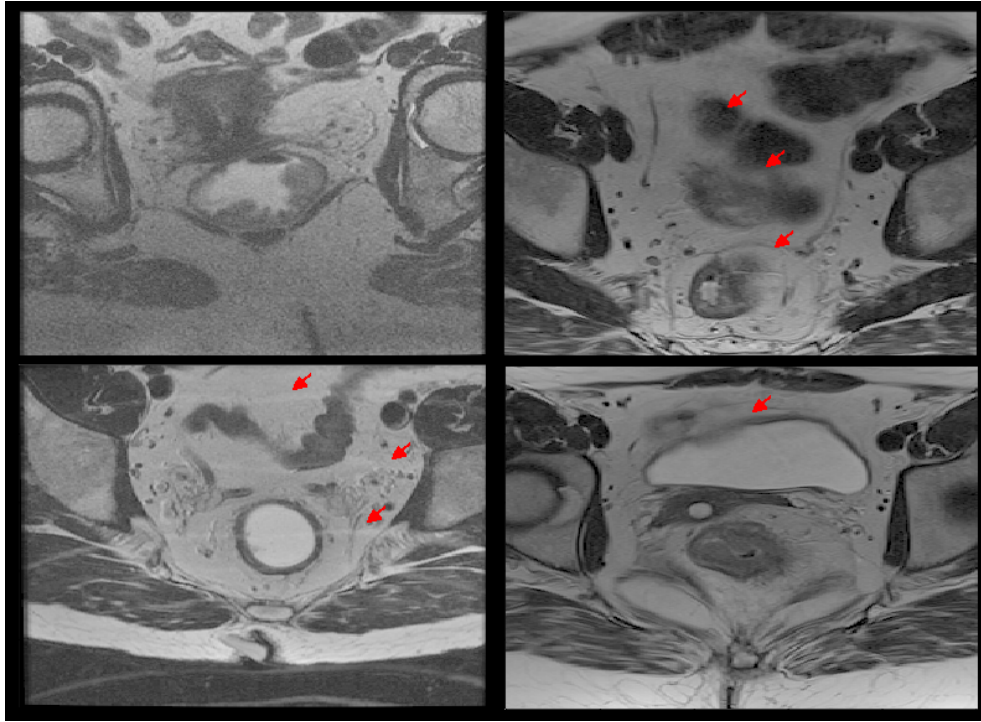


Figure 2: Two MR slices from the Gemelli polyclinic set (left) and from the Maastricht clinic (right). Images show significant differences in terms of high frequency noise (upper left), horizontal lines (bottom left), and spotted blurring artefact (upper and bottom right). Such artefacts have been highlighted using red arrows, for the sake of readability.

217 some horizontal interferences can be spotted (and are pointed in the fig-
 218 ure). On the other hand, images acquired at the Maastricht clinic may present
 219 blurring artefacts, as highlighted in the figure.

220 In order to provide the appropriate input for the proposed approach, MR
 221 scans have been processed using the moddicom R library [12]. Moddicom is
 222 an open source library that allows to: (i) deal with DICOM files in order to
 223 extract images and contouring information; (ii) process and store extracted
 224 data; and (iii) analyse stored data to extract morphological and structural
 225 features. DICOM (Digital Imaging and Communications in Medicine) is
 226 a standard for storing and transmitting medical images enabling the inte-
 227 gration of medical imaging devices such as scanners, servers, workstations,
 228 printers, network hardware, and picture archiving and communication sys-
 229 tems (PACS) from multiple manufacturers.

230 The clinical outcome to be predicted through the generation of radiomics-
231 based models is the pathological complete response (pCR) after surgery,
232 which indicates that there is no residual histological evidence of tumour after
233 surgery. pCR is increasingly found to be a reasonable surrogate for long-term
234 favourable outcomes [23]. In the considered datasets, 21–23% of the cases
235 show a positive pCR. The output of the proposed approach comes in the
236 form of probability of pCR; while the threshold can be provided as input by
237 the user, in this case we exploited a 50%-value threshold. Remarkably, the
238 probability value provides implicitly an estimation of the reliability of the
239 prediction: the closer the probability is to 50%, the lower the confidence.

240 With the aim of limiting the possibility of overfitting, predictive models
241 are evaluated using a 10-fold cross-validation strategy.

242 Performance are measured in terms of specificity and sensitivity. The
243 former measures the so-called true negative rate, i.e., the proportion of nega-
244 tive cases that are correctly identified as such. In our analysis, negative cases
245 correspond to the presence of residual histological evidence of tumour, and
246 the absence of a complete pathological response. Sensitivity (also called the
247 true positive rate) focuses on the proportion of correctly classified positive
248 cases.

249 3.1. Results

250 Hereinafter we will refer to the predictive model trained, using the pro-
251 posed framework, on the data set from the Gemelli polyclinic and the Maastro
252 clinic, as respectively, *Ag.G* and *Ag.M*. On the basis of the considered train-
253 ing data, the optimisation procedure included in the architecture lead to the
254 generation of differently-structured predictive models:

- 255 • The *Ag.G* model is based on a logistic regression built using cT (clinical
256 T stage), the zone size entropy [15] after the application of a LoG
257 with $\sigma = 0.35$ and the Skewness of the grey-level distribution after the
258 application of a LoG with $\sigma = 0.59$.
- 259 • The *Ag.M* predictive model is based on two covariates, the Grey level
260 co-occurrence correlation [15] obtained with a $\sigma = 0.84$ and the Grey
261 level co-occurrence joint entropy obtained with a $\sigma = 0.54$. The agent
262 decides automatically the number of covariates to consider according
263 to the size of the provided training set.

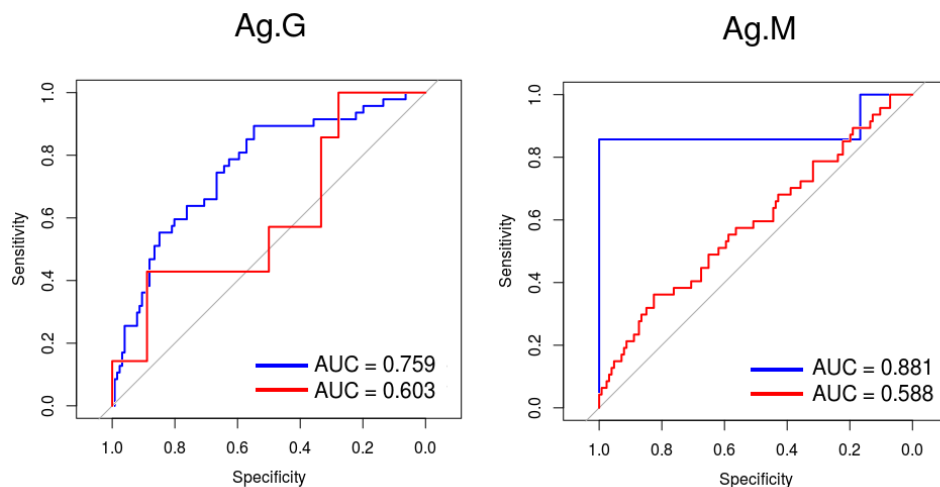


Figure 3: ROC curve of the predictive model trained on the Gemelli polyclinic’s data set (left) and on the data acquired by the Maastricht clinic (right). Blue is used to indicate the ROC observed on the training set, –in cross-validation. Red indicates the ROC obtained on the (external) testing set.

264 It should be noted that the automated optimisation is performed greedily,
 265 following the expected AUC value.

266 Figure 3 shows the receiver operating characteristic curve (ROC) of Ag.G
 267 and Ag.M on both the training set (blue) and the external testing set (red).
 268 Unsurprisingly, the performance of the models tend to be better on the train-
 269 ing set, rather than on the testing set. This is because the testing set images
 270 are affected by a different type of noise than the training set ones (examples
 271 have been discussed in Figure 2). The extremely good performance of Ag.M
 272 on the same data from the Maastricht clinic is possibly due to two main
 273 reasons: (i) the limited size of the set, which may result in some overfitting;
 274 and (ii) the fact that images acquired by the Maastricht clinic show a very
 275 limited noise, or a type of noise to which considered features are robust.

276 3.1.1. Features Analysis

277 To shed some light on the informativeness and the significance of con-
 278 sidered features in the two data sets, we performed an univariate analysis
 279 between each feature and the pCR outcome to be predicted. The analysis
 280 was performed using the Mann-Whitney test ($p < 0.05$). Table 1 presents
 281 the results of the investigation in terms of number of features that have a

Table 1: Number of features that, at least at one σ value and according to an univariate analysis performed using the Mann-Whitney test ($p < 0.05$), are correlated with the outcome to predict. Features are grouped according to the class they belong to. Results are provided for each considered data set, and also in terms of features which are relevant for both sets.

	BASIC	GLCM	GLRLM	GLSZM
Policlinico Gemelli	8	9	5	4
Maastric clinic	4	9	2	3
Common features	1	4	1	1

282 correlation with the outcome to predict. For each feature, only the most
 283 representative σ has been considered. Features are grouped according to the
 284 class they belong to. As a first remark, we observe that out of the total set of
 285 available features, a large subset (more than 20%) has a significant correlation
 286 with the pCR outcome to be predicted. Considering that the univariate
 287 analysis can not take into account combinations of features, this result seems
 288 to suggest that considered features can be very informative, as they carry
 289 useful information for predicting the required pCR outcome.

290 Results presented in Table 1 also highlight the limited overlap between
 291 the features deemed to be significant between the two data sets. In total, 7
 292 features are identified by the univariate analysis, for at least one σ value, in
 293 both the sets.

- 294 • Entropy, BASIC;
- 295 • Sum Entropy: Textural features, GLCM;
- 296 • Correlation: Textural features, GLCM;
- 297 • Sum variance: Textural features, GLCM;
- 298 • Cluster tendency: Textural features, GLCM;
- 299 • Run Entropy: Textural features, GLRLM;
- 300 • Large zone high grey level emphasis: Textural features, GLDZM.

301 Interestingly, most of the features (4) come from the GLCM class, which
 302 includes textural features about the grey level co-occurrence. This suggests

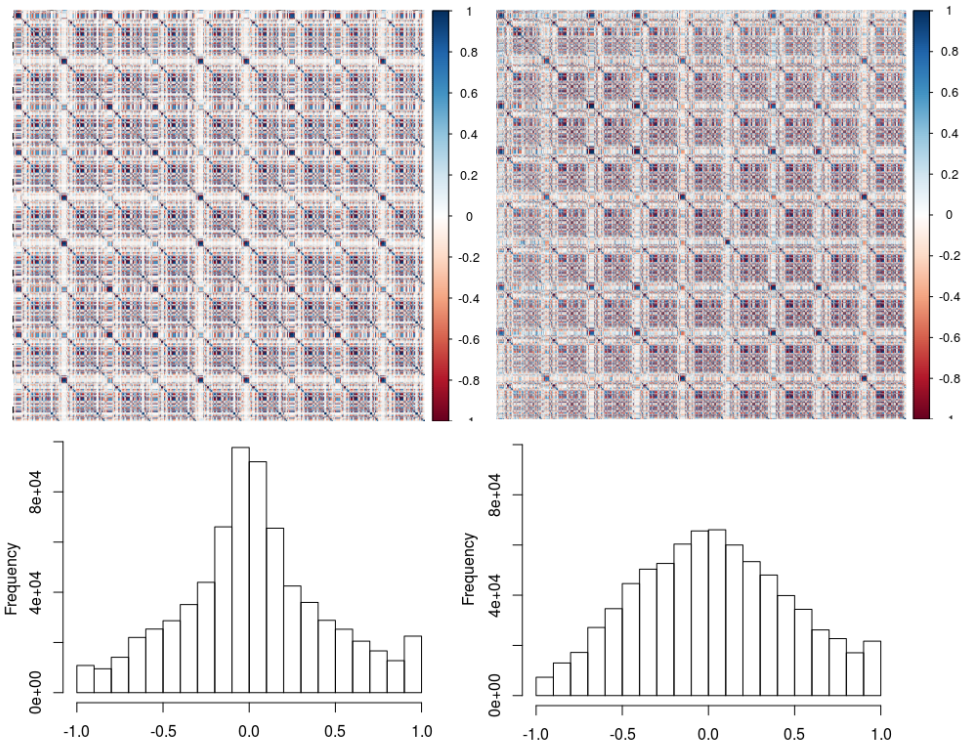


Figure 4: Cross-correlation matrices, using the Pearson correlation coefficient, obtained by analysing the data set of the Gemelli polyclinic (left) and Maastricht clinic (right) are presented in the top half. Bottom half shows the distribution of the coefficients under the form of histograms. 0.0 indicates that no correlation is found, while +1 (−1) identifies cases with strong direct (inverse) correlation.

303 that this class is, in general, more robust with regards to the kind of noise
 304 that affects the medical images acquired by the two considered centres.

305 Figure 4 shows the cross-correlation matrices of the extracted features,
 306 and the bivariate correlation –measured using the Pearson correlation coef-
 307 ficient. For the sake of readability, features in the histograms are ordered
 308 following the order used in the matrices. Evidence seems to indicate that in
 309 the Ag.G set, features have a lower correlation: the region around 0 is very
 310 populated. This is possibly due to the noisy of the images in the set, that
 311 may reduce the informativeness of extracted information. On the contrary,
 312 features in the Ag.M model show a higher level of correlation, as correlation
 313 values are evenly distributed among the scale.

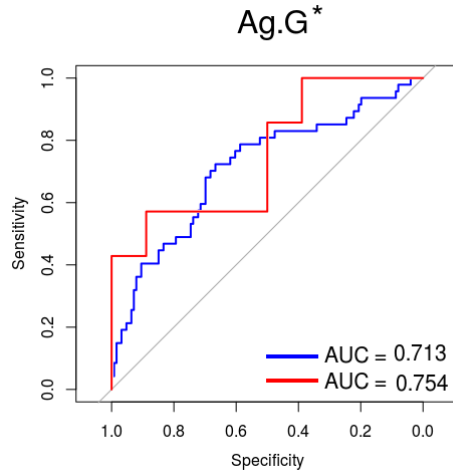


Figure 5: The ROC of a predictive model created by considering as training data images acquired by the Gemelli polyclinic. The model, called Ag.G*, shows to be more general and robust than the Ag.G and Ag.M models, but it delivers slightly worse performance. Blue is used to indicate the ROC observed on the training set, in cross-validation. Red indicates the ROC obtained on the Maastrro testing set.

314 *3.1.2. General Predictive Models*

315 It is worth reminding that the Ag.M and Ag.G models have been op-
 316 timised by the proposed system in order to maximise the performance on
 317 images from the corresponding medical centre. Results presented in Figure
 318 3 indicate that trained model perform poorly on a different data set. There-
 319 fore, the question naturally arises: *Is it possible to generate a more general*
 320 *and robust predictive model?* To answer this question, we configured the
 321 proposed system in order to generate a predictive model according to the
 322 approach proposed in [21]: their work was based on a very limited set of
 323 features, and showed to be portable and robust. We refer to the resulting
 324 model as Ag.G*, because it has been trained using data from the Gemelli
 325 clinic. The Ag.G* model is based on a logistic regression built using cT (clin-
 326 ical T stage), the entropy of the grey-level distribution after the application
 327 of a LoG with $\sigma=0.35$ and the Skewness of the grey-level distribution after
 328 the application of a LoG with $\sigma=0.49$.

329 Figure 5 shows the performance of the Ag.G* predictive model. Blue
 330 is used to indicate the ROC observed on the Gemelli training set, in cross-
 331 validation. Red indicates the ROC obtained on the Maastrro testing set.

332 The generated predictive model provides an interesting trade-off between
333 portability and performance: while the performance on the training set are
334 not as good as those delivered by the Ag.G or Ag.M models, the Ag.G*
335 approach is more robust when used on data from a different centre. This
336 seems to indicate that it is possible to generate a more general and robust
337 model, but at the cost of reduced performance on the specific set.

338 4. Discussion

339 According to the presented results, the proposed approach is able to deal
340 with all the steps of a radiomics analysis on data gathered by different cen-
341 tres. Specifically, the proposed framework showed to be capable of identifying
342 a suitable set of informative feature to maximise the performance –measured
343 in terms of AUC with regards to the outcome to be predicted– of a given
344 class of predictive models. In this work, we focused on logistic regression, but
345 the modularity of the framework allows to easily substitute logistic regression
346 with a different class of approaches, or even to consider more approaches at
347 once. We also highlighted how the framework can be exploited for comparing
348 predictive models generated for different data sets, and how the correspond-
349 ing features (and their characteristics) can be compared and analysed. Re-
350 markably, this analysis can potentially lead to identify issues in the machines
351 or in the environment, or even suggests the presence of procedural issues.

352 The empirical results presented in the previous section seem to confirm
353 the importance of centre-specific radiomics-based predictive models. Figure
354 3 suggests that the use of “general” predictive models can lead to very poor
355 predictive performance. However, results also clearly indicate the value of a
356 radiomics-based decision support system, that can provide useful information
357 to physicians and can lead to a more effective planning of the treatments for
358 patients. A trade-off between portability and performance is presented in
359 Figure 5: remarkably, the generated predictive model is less sensitive to the
360 difference in the data sets, for instance in terms of image noise. On the
361 other hand, general performance is worse than those that can be achieved by
362 exploiting centre-specific models.

363 Taking a different perspective, which is necessarily more speculative than
364 the analysis of the results presented in the previous section, we can identify a
365 number of ways in which the presented system can be exploited with regards
366 to radiomics:

- 367 • For the sake of the explainability of the predictive models, a number
368 of different models can be generated for predicting the same clinical
369 outcome. In particular, emphasis can be given to approaches that
370 generate models easy to investigate and analyse by humans, so that
371 an expert user can visualise the generated model, and can explore the
372 relevance of features with regards to the considered clinical outcome.
373 While the number of features can be extremely large, focusing on the
374 described classes of features can highlight the importance of a set of
375 feature, that can be used all together.
- 376 • The proposed framework can also allow users to provide as input a spe-
377 cific set of features to be analysed. Such features are then exploited for
378 generating predictive models, and can be compared in terms of relation
379 and correlation. This may allow to investigate features believed to be
380 informative in the relevant literature, and also to assess their usefulness
381 in the presence of images acquired by using different machines, settings,
382 or centres.
- 383 • Different data sets can also be compared, in terms of relevant features.
384 For instance, in the presence of large multi-centric studies, it may be
385 useful to identify centres which acquire images with similar properties;
386 that would reduce the noise of the analysis, and maximise the proba-
387 bility of generating an highly performant yet general –with regards to
388 the considered clinics– predictive model.

389 The physicians involved in the experimental analysis positively evaluated
390 the experience with the proposed agent. The agent allows the medical experts
391 to focus on the actual goals of their investigation and analysis: optimisation
392 and low-level details are optimised by the agent architecture without the
393 need of human guidance. The agent, given a range of alternative modules to
394 choose from, and the parametrisation of each module, can transparently test
395 different alternatives in order to achieve the specified goal. In the presented
396 experimental analysis, the goal was to generate a LR-based predictive model
397 of the pCR of patients treated for rectal cancer. A very important aspect
398 that the agent-based structure can support, but has not been integrated in
399 the proposed system yet, is the ability to *explain* results, and to *motivate*
400 the decisions. We are extremely interested in develop these aspects as part
401 of our future work.

402 An important aspect to consider, particularly in the case of agent-based
403 decision support system, is the ability to generalise on different data sets.
404 This has been partly covered in the experimental analysis by considering
405 images from two different centres. However, also due to the very limited
406 amount of contoured images available in the radiomics field, it is hard to
407 empirically demonstrate that the proposed agent will easily generalise on data
408 sets where different type of cancer are treated. On this matter, a preliminary
409 study performed by exploiting the proposed agent on a data set considering 15
410 patients affected by glioblastoma (a form of brain cancer) seems to indicate
411 that the agent, also due to its modularity, can generalise on significantly
412 different sets of MRI images [24].

413 The agent introduced in this work can play a central role in a distributed
414 learning scenario [25], where different agents cooperate to converge to a ro-
415 bust and shared predictive model while preserving the privacy of patients.
416 This can be achieved by exploiting an iterative approach, shown in Figure 6,
417 composed by four main steps: (a) Each centre trains a local model, (b) the
418 models are sent to a Master, (c) the Master calculates a model, considering
419 weighting the contribute of each centre with the cardinality of the locally
420 available sets, then calculates some new coefficients for each node, (d) the
421 coefficients are sent to each node and the process can be repeated until in a
422 (c) step a convergence criteria is reached.

423 However, we also envisage the use of the introduced agent in distributed
424 learning scenarios where federated learning approaches are exploited [26],
425 where there is no need for a master to coordinate learning and merge a
426 general model.

427 5. Conclusion

428 Radiomics is a topic that is gaining a significant interest in the scientific
429 community, as testified by the growing number of publications that can now
430 be found on the online library of medicine-related articles Pubmed. While
431 still in its infancy, a number of tools are now available for supporting ra-
432 diomics, as well as standardisation initiatives. These initiatives, such as
433 IBSI [15] are mainly aimed at maximising the reproducibility of results.

434 Despite the growing interest and the number of already available tools,
435 there is a lack of agents that can deal with all the steps of the radiomics
436 process. Existing tools are mainly aimed at facilitating the extraction of
437 features, and at extending the set of features that can be extracted from a

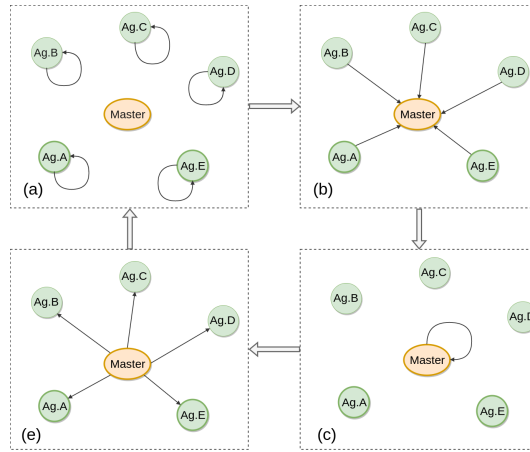


Figure 6: An example of a possible architecture of cooperative agents to converge to a robust and shared model by an iterative approach. Initially, each centre trains a local model using its local agent (a), that is then sent to a Master (b). The master agent merges the models into a general one (c), and send it back to each centre (d).

438 medical image, only. Crucial aspects, such as features selection, correlation
 439 between features and the outcome to predict, and the generation of the actual
 440 predictive model, are normally ignored.

441 In the light of the peculiarities of radiomics, such as the very different
 442 characteristics of acquired images according to the exploited machine or the
 443 location of the machine, two lines of evolution of radiomics can be envisaged:

444 • **General and Portable models.** By identifying features that are ro-
 445 bust with regards to different sources of image noise that can be found
 446 in images acquired in different centres, it could be possible to generate
 447 general and portable predictive models, which would allow to exploit
 448 the availability of numerous –even though sparse– sets of images. On
 449 the other hand, the focus on portability would lead to under performing
 450 (when compared to centre-specific) models, with clear negative reper-
 451 cussion on the quality of the treatment delivered to patients.

452 • **Centre-specific models.** By dropping any requirement related to the
 453 portability of models, a significant performance boost can be obtained
 454 by highly optimised centre-specific predictive models. This would allow
 455 every centre to train a model that is specific for the characteristics of
 456 the machines, and for the typology of noise which is included in the

457 acquired images. A significant drawback would then be that any change
458 in the environment, e.g. a new machine is bought to substitute and
459 obsolete one, may dramatically reduce the reliability of the generated
460 model. Furthermore, this approach does not allow to investigate, in a
461 general sense, the importance and robustness of features.

462 In this paper, we introduced a medical agent-based decision support sys-
463 tem which is capable of supporting the whole radiomics process³. The agent
464 can be given a high level goal, and is then able to reason in order to achieve
465 it. Given a set of medical images, the proposed system is able to extract a
466 wide range of features, to analyse and select them with regards to the out-
467 come to predict, and to generate an optimised predictive models. When data
468 from a new patient is provided as input, the proposed agent is able to collect
469 features from available medical images and patient’s data, and to return a
470 prediction about the clinical outcome of a proposed treatment. Our experi-
471 mental analysis demonstrated the ability of the system, and highlighted that
472 the proposed architecture is capable of supporting both the lines of research
473 mentioned above: predictive models can be optimised for a specific centre,
474 and then exchanged in order to analyse the differences. Furthermore, data
475 sets can be merged in order to generated general predictive models, or more
476 general approaches can be used for the creation of predictive models.

477 We see several avenues for future work. We are actively working on four
478 aspects.

- 479 1. The exploitation of additional data sets for testing the capability of the
480 proposed medical agent-based decision support system to generalise on
481 different types of images and contouring.
- 482 2. A graphical user interface, that would create a more comfortable envi-
483 ronment for researchers.
- 484 3. The development of additional modules for performing different kind of
485 features selection algorithms, and extend the set of techniques that can
486 be used for generating the actual predictive model. Specifically, we are
487 looking into neural networks [27], SVM [28], and decision trees [29].
488 Neural networks will need only a subset of the currently developed
489 modules of the proposed decision support agent, but this aspect is
490 already supported by the modularity of the system.

³The agent-based decision support system can be downloaded from <https://github.com/robertogattabs/RadAgent>

- 491 4. An approach for extracting information about the spectral components
492 (and other measurable aspects) of image noise of images included in the
493 considered data set. Such analysis will allow to assess the impact of
494 different sort of noise on the predictive capabilities of (some set of)
495 considered features, and to better counter-balance it. As a result, it
496 would be possible to generate more robust predictive models.
- 497 5. Improving the capabilities of the agent, so that it can explain the ob-
498 tained results and motivate the decisions taken.
- 499 6. An architecture to support multi-centric investigation based on the
500 distributed learning principles.

501 **Acknowledgement**

502 We want to thank Silvia Chiesa for providing us insights and data of MRI
503 images of brain cancer patients.

- 504 [1] F. Collins, *The language of life: DNA and the revolution in personalised*
505 *medicine*, Profile Books, 2010.
- 506 [2] L. Wen-Ling, T. Fuu-Jen, *Personalized medicine: A paradigm shift in*
507 *healthcare*, *BioMedicine* 3 (2013) 66 – 72.
- 508 [3] P. Lambin, E. Rios-Velazquez, R. Leijenaar, S. Carvalho, R. G. van
509 Stiphout, P. Granton, C. M. Zegers, R. Gillies, R. Boellard, A. Dekker,
510 et al., *Radiomics: extracting more information from medical images*
511 *using advanced feature analysis*, *European Journal of Cancer* 48 (2012)
512 441–446.
- 513 [4] R. J. Gillies, P. E. Kinahan, H. Hricak, *Radiomics: Images are more*
514 *than pictures, they are data*, *Radiology* 278 (2016) 563–577.
- 515 [5] E. Sala, E. Mema, Y. Himoto, H. Veeraraghavan, B. JD., et al., *Unrav-*
516 *elling tumour heterogeneity using next-generation imaging: radiomics,*
517 *radiogenomics, and habitat imaging*, *Clinical Radiology* 72 (2017) 3–10.
- 518 [6] M. Hatt, F. Tixier, L. Pierce, P. Kinahan, C. Le Rest, D. Visvikis,
519 *Characterization of pet/ct images using texture analysis: the past, the*
520 *present any future?*, *European Journal of Nuclear Medicine and Molec-*
521 *ular Imaging* 44 (2017) 151–165.

- 522 [7] G. Lee, H. Lee, H. Park, M. Schiebler, E. van Beek, Y. Ohno, J. Seo,
523 A. Leung, Radiomics and its emerging role in lung cancer research,
524 imaging biomarkers and clinical management: State of the art, Euro-
525 pean Journal of Radiology (2016).
- 526 [8] S. Alobaidli, S. McQuaid, C. South, V. Prakash, P. Evans, N. A., The
527 role of texture analysis in imaging as an outcome predictor and potential
528 tool in radiotherapy treatment planning, British Journal of Radiology
529 (2014).
- 530 [9] Y.-H. D. Fang, C.-Y. Lin, M.-J. Shih, et al., Development and eval-
531 uation of an open-source software package cgita for quantifying tumor
532 heterogeneity with molecular images, BioMed Research International
533 (2014).
- 534 [10] C. Chatwin, R. Young, B. Ganeshan, Texrad-feedback plc - cancer
535 management imaging software. project report (2015).
- 536 [11] M. Strzelecki, P. Szczypinski, A. Materka, A. Klepaczko, A software tool
537 for automatic classification and segmentation of 2d/3d medical images,
538 Nuclear Instruments and Methods In Physics Research 702 (2013) 137–
539 140.
- 540 [12] N. Dinapoli, A. Alitto, M. Vallati, R. Gatta, R. Autorino, L. Boldrini,
541 A. Damiani, V. Valentini, Moddicom: a complete and easily accessible
542 library for prognostic evaluations relying on image features, Conference
543 Proceeding IEEE Engineering in Medicine and Biology Society (2015)
544 771–774.
- 545 [13] J. J. van Griethuysen, A. Fedorov, C. Parmar, A. Hosny, N. Aucoin,
546 V. Narayan, R. G. Beets-Tan, J.-C. Fillion-Robin, S. Pieper, H. J. Aerts,
547 Computational radiomics system to decode the radiographic phenotype,
548 Cancer research 77 (2017) e104–e107.
- 549 [14] L. Zhang, D. Fried, X. Fave, L. Hunter, J. Yang, C. L., ibex: An
550 open infrastructure software platform to facilitate collaborative work in
551 radiomics, Medical Physics 42 (2015) 13411353.
- 552 [15] A. Zwanenburg, S. Leger, M. Vallières, S. Löck, Image biomarker stan-
553 dardisation initiative - feature definitions, CoRR abs/1612.07003 (2016).

- 554 [16] N. Dinapoli, C. Cas, B. Barbaro, G. Chiloiro, A. Damiani, M. Di Matteo,
555 A. Farchione, M. Gambacorta, R. Gatta, V. Lanzotti, C. Masciocchi,
556 V. Valentini, Radiomics for rectal cancer, *Translational Cancer Research*
557 5 (2016).
- 558 [17] B. Altazi, G. Zhang, D. Fernandez, M. Montejo, D. Hunt, J. Werner,
559 M. Biagioli, E. Moros, Reproducibility of f18-fdg pet radiomic features
560 for different cervical tumor segmentation methods, gray-level discretiza-
561 tion, and reconstruction algorithms, *Journal of Applied Clinical and*
562 *Medical Physics* (2017).
- 563 [18] A. J. Wong, A. Kanwar, A. S. Mohamed, C. D. Fuller, Radiomics in
564 head and neck cancer: from exploration to application, *Translational*
565 *Cancer Research* 5 (2016).
- 566 [19] G. Collins, J. Reitsma, D. Altman, K. Moons, Transparent reporting
567 of a multivariable prediction model for individual prognosis or diagnosis
568 (tripod): the tripod statement, *Annals of Internal Medicine* 6 (2015)
569 55–63.
- 570 [20] S. Powers, V. McGuire, L. Bernstein, A. Canchola, A. Whittemore,
571 Evaluating disease prediction models using a cohort whose covariate
572 distribution differs from that of the target population, *Stat Methods*
573 *Med Res* (2017).
- 574 [21] N. Dinapoli, J. van Soest, C. Masciocchi, C. Cas, V. Lanni, A. Dami-
575 ani, R. Gatta, B. Barbaro, M. Di Matteo, F. Cellini, M. Gambacorta,
576 A. Dekker, P. Lambin, V. V, Radiomics in magnetic resonance imaging
577 for prognosis in patients with rectal cancer: An independent external
578 validation, *Radiation Oncology* 96 (2016) E180–E181.
- 579 [22] N. G. Burnet, S. J. Thomas, K. E. Burton, S. J. Jefferies, Defining the
580 tumour and target volumes for radiotherapy.cancer imaging, *Cancer*
581 *Imaging* (2004).
- 582 [23] L. Ferrari, A. Fichera, Neoadjuvant chemoradiation therapy and patho-
583 logical complete response in rectal cancer, *Gastroenterology Report* 3
584 (2015) 277–288.
- 585 [24] S. Chiesa, M. Lupattelli, R. Gatta, I. Palumbo, M. Balducci, R. Tar-
586 ducci, R. Cusumano, C. Masciocchi, J. Lenkowicz, M. Martucci,

- 587 P. Floridi, N. Dinapoli, F. Beghella Bartoli, V. Valentini, C. Aristei,
588 C035 delta radiomica delle caratteristiche delle immagini per predire
589 gli outcomes nei pazienti con glioblastoma multiforme: studio prospet-
590 tico multicentrico- gli.f.a. project (english), in: Associazione Italiana
591 Radioterapia Oncologica (AIRO).
- 592 [25] A. Damiani, M. Vallati, R. Gatta, N. Dinapoli, A. Jochems, T. Deist,
593 J. van Soest, A. Dekker, V. Valentini, Distributed learning to protect
594 privacy in multi-centric clinical studies, in: Artificial Intelligence in
595 Medicine, AIME, pp. 66–75.
- 596 [26] T. S. Brisimi, R. Chen, T. Mela, A. Olshevsky, I. C. Paschalidis, W. Shi,
597 Federated learning of predictive models from federated electronic health
598 records, *International Journal of Medical Informatics* 112 (2018) 59 –
599 67.
- 600 [27] S. Dreiseitl, L. Ohno-Machado, Logistic regression and artificial neu-
601 ral network classification models: a methodology review, *Journal of*
602 *biomedical informatics* 35 (2002) 352–359.
- 603 [28] C. Cortes, V. Vapnik, Support vector machine, *Machine learning* 20
604 (1995) 273–297.
- 605 [29] W.-Y. Loh, *Classification and regression trees*, Wiley Interdisciplinary
606 *Reviews: Data Mining and Knowledge Discovery* 1 (2011) 14–23.

607 **Appendix A. Detailed List of Features**

608 Here we provide the list of 92 features exploited in this work. They
609 are described in the Image biomarker standardisation initiative Reference
610 manual [15]. In order to make it easier, for the interested reader, to identify
611 the features in the reference manual, the same id is used in the following list.

612 • **BASIC**

- 613 – 4.1.1 Volume
- 614 – 4.1.3 Surface area
- 615 – 4.1.4 Surface to volume ratio

- 616 – 4.1.5 Compactness 1
- 617 – 4.1.6 Compactness 2
- 618 – 4.1.7 Spherical disproportion
- 619 – 4.1.8 Sphericity
- 620 – 4.1.9 Asphericity
- 621 – 4.1.10 Centre of mass shift
- 622 – 4.1.11 Maximum 3D diameter
- 623 – 4.1.12 Major axis length
- 624 – 4.1.13 Minor axis length
- 625 – 4.1.14 Least axis length
- 626 – 4.1.15 Elongation
- 627 – 4.1.16 Flatness
- 628 – 4.3.1 Mean
- 629 – 4.3.2 Variance
- 630 – 4.3.3 Skewness
- 631 – 4.3.4 Kurtosis
- 632 – 4.3.5 Median
- 633 – 4.3.6 Minimum grey level
- 634 – 4.3.7 10th percentile
- 635 – 4.3.8 90th percentile
- 636 – 4.3.9 Maximum grey level
- 637 – 4.3.10 Interquartile range
- 638 – 4.3.11 Range
- 639 – 4.3.12 Mean absolute deviation
- 640 – 4.3.13 Robust mean absolute deviation
- 641 – 4.3.17 Energy
- 642 – 4.3.18 Root mean square
- 643 – 4.4.18 Entropy
- 644 – 4.4.19 Uniformity

- 645 • **Grey level co-occurrence based features – Texture features**
- 646 **(GLCM)**
- 647 – 4.6.1 Joint maximum
- 648 – 4.6.2 Joint average
- 649 – 4.6.3 Joint variance
- 650 – 4.6.4 Joint entropy
- 651 – 4.6.5 Difference average
- 652 – 4.6.6 Difference variance
- 653 – 4.6.7 Difference entropy
- 654 – 4.6.8 Sum average
- 655 – 4.6.9 Sum variance
- 656 – 4.6.10 Sum entropy
- 657 – 4.6.11 Angular second moment
- 658 – 4.6.12 Contrast
- 659 – 4.6.13 Dissimilarity
- 660 – 4.6.14 Inverse difference
- 661 – 4.6.15 Inverse difference normalised
- 662 – 4.6.16 Inverse difference moment
- 663 – 4.6.17 Inverse difference moment normalised
- 664 – 4.6.18 Inverse variance
- 665 – 4.6.19 Correlation
- 666 – 4.6.20 Autocorrelation
- 667 – 4.6.21 Cluster tendency
- 668 – 4.6.22 Cluster shade
- 669 – 4.6.23 Cluster prominence
- 670 – 4.6.24 First measure of information correlation
- 671 – 4.6.25 Second measure of information correlation
- 672 • **Grey level run length based features – Texture features (GLRLM)**

- 673 – 4.7.1 Short runs emphasis
- 674 – 4.7.2 Long runs emphasis
- 675 – 4.7.3 Low grey level run emphasis
- 676 – 4.7.4 High grey level run emphasis
- 677 – 4.7.5 Short run low grey level emphasis
- 678 – 4.7.6 Short run high grey level emphasis
- 679 – 4.7.7 Long run low grey level emphasis
- 680 – 4.7.8 Long run high grey level emphasis
- 681 – 4.7.9 Grey level non-uniformity
- 682 – 4.7.10 Grey level non-uniformity normalised
- 683 – 4.7.11 Run length non-uniformity
- 684 – 4.7.12 Run length non-uniformity normalised
- 685 – 4.7.13 Run percentage
- 686 – 4.7.14 Grey level variance
- 687 – 4.7.15 Run length variance
- 688 – 4.7.16 Run entropy
- 689 • **Grey level size zone based features – Texture features (GLDZM)**
- 690 – 4.8.1 Small zone emphasis
- 691 – 4.8.2 Large zone emphasis
- 692 – 4.8.3 Low grey level zone emphasis
- 693 – 4.8.4 High grey level zone emphasis
- 694 – 4.8.5 Small zone low grey level emphasis
- 695 – 4.8.6 Small zone high grey level emphasis
- 696 – 4.8.7 Large zone low grey level emphasis
- 697 – 4.8.8 Large zone high grey level emphasis
- 698 – 4.8.9 Grey level non-uniformity
- 699 – 4.8.10 Grey level non-uniformity normalised
- 700 – 4.8.11 Zone size non-uniformity

- 701 – 4.8.12 Zone size non-uniformity normalised
- 702 – 4.8.13 Zone percentage
- 703 – 4.8.14 Grey level variance
- 704 – 4.8.15 Zone size variance
- 705 – 4.8.16 Zone size entropy