

Alma Mater Studiorum · University of Bologna

Department of Physics and Astronomy

Master Degree in Physics

THESIS TITLE

Supervisor:

Prof. Enrico Giampieri

Co-Supervisor:

Dott.sa Lidia Strigari

Submitted by:

Lorenzo Spagnoli

Dedication...

Abstract

Contingency table per le segmentazioni sostituibili

		ICU Admission
		0 1
Death	0	
	1	8 1

Contingency table per le segmentazioni brutte o dubbie

		ICU Admission
		0 1
Death	0	
	1	32 4

		ICU Admission
		0 1
Death	1	
	0	11 2

Since the start of 2020 *Sars-COVID19* has given rise to a world-wide pandemic. In an attempt to slow down the fast and uncontrollable spreading of this disease various prevention and diagnostic methods have been developed. In this thesis, out of all these various methods, the attention is going to be put on Machine Learning methods used to predict prognosis that are based, for the most part, on data originating from medical images. The techniques belonging to the field of radiomes will be used to extract information from images segmented using a software available in the hospital that provided the clinical data as well as the images. The usefulness of different families of variables will be evaluated through their performance in the methods used, namely Lasso regularized regression and Random Forest. Dimensionality reduction techniques will be used to attain a better understanding of the dataset at hand. Following a first introductory chapter in the second a basic theoretical overview of the necessary core concepts that will be needed throughout this whole work will be provided and then the focus will be shifted on the various methods and instruments used in the development of this thesis. The third is going to be a report of the results and finally some conclusions will be derived from the previously presented results. It will be concluded that the segmentation and feature extraction step is of pivotal importance in driving the performance of the predictions. In fact, in this thesis, it seems that the information from the images adds no significant information to that derived from the clinical data. This can be taken as a symptom that the more complex *Sars-COVID19* cases are still too difficult to be segmented automatically, or semi-automatically by untrained personnel, which will lead to counter-intuitive results further down the analysis pipeline.

Contents

1	Introduction	1
2	Materials and methodologies	4
2.1	Theoretical background	4
2.1.1	Medical Images	4
2.1.2	Artificial Intelligence (AI) and Machine Learning(ML)	19
2.1.3	Combining radiological images with AI: Image segmentation and Radiomics	30
2.2	Data and objective	38
2.3	Preprocessing and data analysis	42
3	Results	48
3.1	LASSO	48
3.1.1	Death	49
3.1.2	ICU Admission	56
3.2	Ramdom forest	60
3.2.1	Death	61
3.2.2	ICU Admission	63
3.3	Dimensionality reduction	66

3.3.1 PCA	67
---------------------	----

¹ Chapter 1

² Introduction

³ Nowadays everybody knows of *Sars-COVID19* which, since the start of 2020, has
⁴ made necessary a few world-wide quarantines forcing everybody in self-isolation. It
⁵ is also well known that, among the main complications and features of this virus,
⁶ symptoms gravity as well as the rate of deterioration of the conditions are some
⁷ of the most relevant and problematic. In some cases asymptomatic or near to
⁸ asymptomatic people may, in the span of a week, get to conditions that require
⁹ hospital admission. This peculiarity is also what heavily complicates the triage
¹⁰ process, since trying to predict with some degree of accuracy the prognosis of the
¹¹ patient at admission is a thoroughly complex task. In this thesis the aim will be to
¹² use data, specifically including data that cannot be easily interpreted by humans,
¹³ to try various methods to predict a couple of clinical outcomes, namely the death of
¹⁴ the patient or the admission in the Intensive Care Unit (ICU), while assessing their
¹⁵ performance. These analyses will be carried out on a dataset of 434 patients with
¹⁶ different variables associated to every person. A part of the variables, which will
¹⁷ be called clinical and radiological, are defined by humans and are generally discrete
¹⁸ in nature but mostly boolean. The most part of the available variables, however,
¹⁹ will be image-derived following the approaches used in the field of radiomics. While
²⁰ the utility of clinical variables, such as age, obesity and history of smoking, is very
²¹ straightforward it's interesting and helpful to understand the basis behind the utility
²² of radiomic and radiological features. Generally speaking it's clear that images have
²³ the ability to convey a slew of useful images, this is especially true in the medical
²⁴ field where digital images are used to inspect also the internal state of the patient
²⁵ giving far more detailed information than that obtainable by visual inspection at
²⁶ the hand of medical professionals. Among the ways in which *Sars-COVID19* can
²⁷ manifest himself the one that is most relevant to the scopes of this thesis pneumonia
²⁸ and the complications that stem from it. Some of these complications, which are
²⁹ not specific of *Sars-COVID19* but can happen in any pneumonia case, display very
³⁰ peculiar patterns when visualizing the lungs through CT exams. These patterns
³¹ are due to the pulmonary response to inflammation which may lead to thickening
³² of the bronchial and alveolar structures up to pleural effusions and collapsed lungs.
³³ Without going too much in clinical detail what is of interest is how these condition

³⁴ manifest themselves in the CT exams:

³⁵ 1. **Ground Glass Opacity(GGO):**

³⁶ Small diffused changes in density of the lung structure cause a hazy look in the affected region. This complicates the individuation of pulmonary vessels.

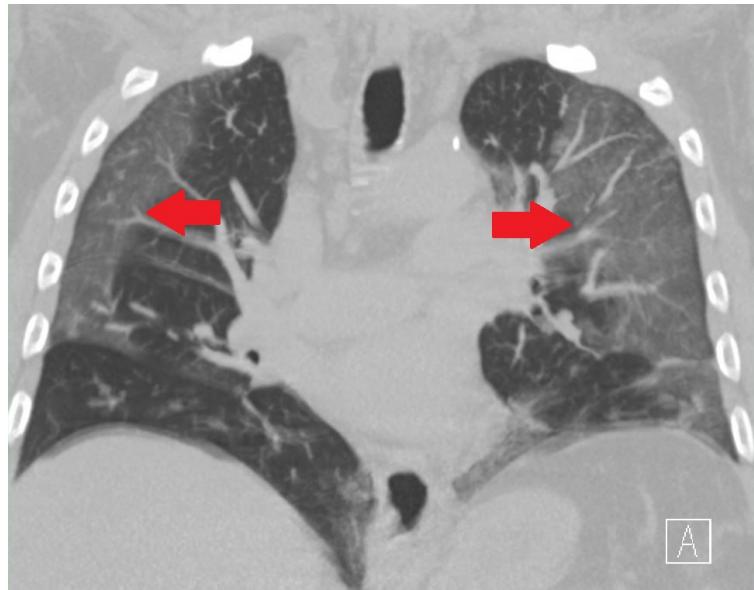


Figure 1.1: Example of GGO

³⁷

³⁸ 2. **Lung Consolidations:**

³⁹ ⁴⁰ Heavier damage reflects in whiter spots in the lung as the surface more closely ⁴¹ resembles outside tissue instead of normal air. The consolidation refer to presence of fluid, cells or tissue in the alveolar spaces

⁴² 3. **Crazy paving:**

⁴³ When GGOs are superimposed with inter-lobular and intra-lobular septal thickening.

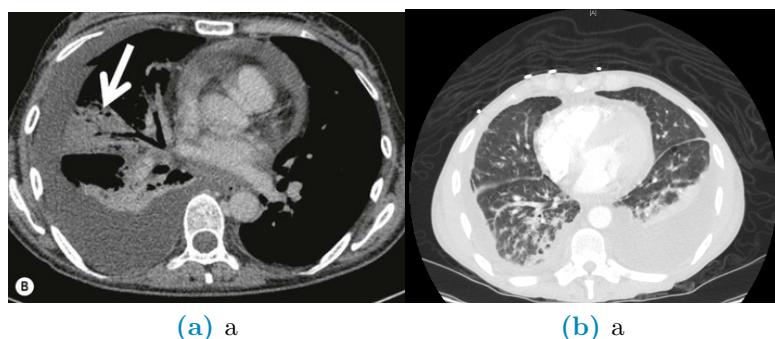


Figure 1.2: Differences between a collapsed lung (a) and pleural effusion(b)

⁴⁴

45 **4. Collapsed Lungs and Pleural Effusion:**

46 Both of these manifest themself as regions of the lungs that take the same
47 coloring as that of tissue outside the lung. The main difference between the
48 two is that collapsed lungs are somewhat rigid structures, they can occur
49 in singular lobes of the lung and stay where they occur. Pleural effusions,
50 however, are actually fluid being located in the lung instead of air. As such
51 these lesions usually are located 'at the bottom' of the lung in which they
52 happen and migrate to the lowest part of the lung according to the position
53 of the patient.

54 Having these manifestation it's clear that they are mainly textural and intensity-
55 like changes in the normal appearance of the lungs. However, whereas these proper-
56 ties can be easily described in a qualitative and subjective way, it's rather complex
57 to describe them in a quantitative and objective way. The field of radiomics, when
58 coupled with digital images and preprocessing steps which must include image seg-
59 mentation, is exactly what undertakes this daunting task. Radiomics comes from
60 the combination of radiology and the suffix *-omics*, which is characteristic of
61 high-throughput methods that aim to generate a large number of numbers, called
62 biomarkers or features, as such it uses very precise and strict mathematical defini-
63 tions to quantify in various ways either shape, textural or intensity based properties
64 of the radiological image under analysis. Given the large numerosity of the features
65 produced by radiomics it's necessary to analyze these kinds of data with methods
66 that rely on Machine Learning and their ability to address high-dimensional prob-
67 lems, be it in a supervised or unsupervised way, in a rather fast and accurate way.
68 Starting from these premises this thesis will be divided in a few chapters and sec-
69 tions. The first step will be taken by providing the general theoretical background
70 regarding the aforementioned topics and techniques, this will be followed by a dess-
71 cription of the data in use as well as a presentation of the analysis methods and
72 resources used. Finally the results of the methods described will be presented and
73 from them a set of concluding remarks will be set forth.

⁷⁴

Chapter 2

⁷⁵

Materials and methodologies

⁷⁶ In this section there's going to be an explanation of the dataset as well as instruments
⁷⁷ and methodologies used to analyze it's properties, as such the first step is going to
⁷⁸ be an in depth discussion of the data available and a general overview of the final
⁷⁹ use. The following step is going to be a description of the preliminary work done to
⁸⁰ the data itself and to the results of this preliminary analysis in order to select the
⁸¹ important features. The final step of this chapter is going to be an explanation of
⁸² the methods used to derive the final results and to evaluate them.

⁸³

2.1 Theoretical background

⁸⁴

2.1.1 Medical Images

⁸⁵ In this section the objective is to simply provide a set of basic definitions pertaining
⁸⁶ to images as well as a general introduction to the methods used to create said
⁸⁷ images. Firstly images are a means of representing in a visual way a physical object
⁸⁸ or set thereof, when talking about images it's common to refer specifically to digital
⁸⁹ images.

⁹⁰ **Definition 2.1.1** (Digital Image). A numerical representation of an object; more
⁹¹ specifically an ordered array of values representing the amount of radiation emitted
⁹² (or reflected) by the object itself. The values of the array are associated to the
⁹³ intensity of the radiation coming from the physical object; to represent the image
⁹⁴ these values need to be associated to a scale and then placed on a discrete 2D grid.
⁹⁵ To store these intensities the physical image is divided into regular rectangular
⁹⁶ spacings, each of which is called pixel¹, to form a 2D grid; inside every spacing is
⁹⁷ then stored a number (or set thereof) which measures the intensity of light, or color,

¹The term pixel seems to originate from a shortening of the expression Picture's (pics=pix) Element(el). The same hold for voxel which stands for Volume Element

98 coming from the physical space corresponding to that grid-spacing. The term digital
 99 refers to the discretization process that inherently happens in storage of the values,
 100 called pixel values, as well as in arranging them within the grid. It's possible to
 101 generalize from 2D images to 3D volumes, simply by stacking images of the same
 102 object obtained at different depths. In this context, the term pixel is substituted by
 103 voxel, however since they are used interchangeably in literature they will, from now
 104 on, be considered equivalent.

105 Generally pixel values stored as integers $p \in [0, 2^n - 1]$ with $p, n \in \mathbb{N}$ or as $p \in [0, 1]$
 106 with $p \in \mathbb{R}$, the type of value stored within each pixel changes the nature of the
 107 image itself. A single value is to be intended as the overall intensity of light coming
 108 from the part of the object contained corresponding to the gridspace and is used for
 109 a gray-scale representation, a set of three² or four³ values can be intended as a color
 110 image.

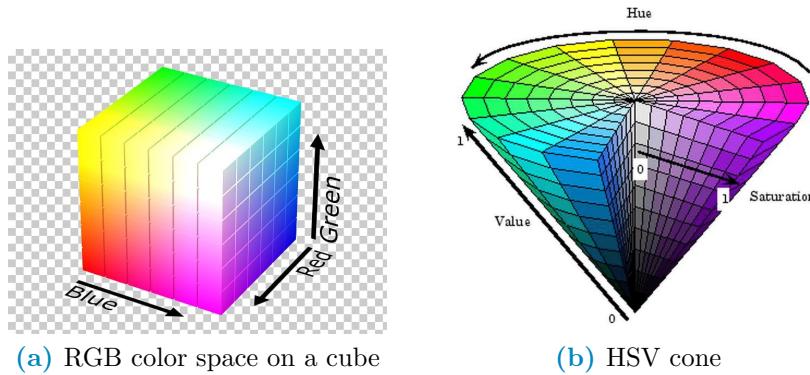


Figure 2.1: Examples of color spaces

111 There are a lot of possible scales for representation⁴, which are sometimes called
 112 color-spaces, however the most noteworthy in the scope of this work is the Hounsfield
 113 unit (HU) scale.

114 **Definition 2.1.2** (Hounsfield unit (HU)). A scale used specifically to describe ra-
 115 diodensity, frequently used in the context of CT (Computed Tomography) exams.
 116 The values are obtained as a transformation of the linear attenuation coefficient ??
 117 of the material being imaged and, since the scale is supposed to be used on humans,
 118 it's defined such that water has value zero and air has the most negative value -1000.
 119 For a more in depth discussion refer to [11]

²The three values correspond each to the intensity of a single color, the most commonly used set of colors is the RGB-scale (Red, Green, Blue). Further information can be found by looking into Tristimulus theory[23]

³Same as RGB but with four colors, the most common scale is CMYK (Cyan, Magenta, Yellow, black)

⁴Besides RGB and CMYK 2.1a the most common color spaces are CIE (Commision Internationale d'Eclairage) and HSV fig:2.1b (Hue,Saturation and Value). Refer to [9] for further details

$$HU = 1000 * \frac{\mu - \mu_{H_2O}}{\mu_{H_2O} - \mu_{Air}} \quad (2.1)$$

120 The utility of this scale is in its definition, since the pixel value depends on
 121 the attenuation coefficient it's possible to individuate a set of ranges that identify,
 122 within good reason, the various tissues in the human body: for example lungs are
 123 [-700, -600] while bone can be in the [500, 1900] range. A more in depth discussion
 124 of the topics relative to Hounsfield units is going to be carried out at a later point
 125 throughout this chapter, in the meantime it's necessary to clarify what are the most
 126 important characteristics of an image:

- 127 • Spatial Resolution: A measure of how many pixels are in the image or, equivalently,
 128 how small each pixel is; a larger resolution implies that smaller details
 129 can be seen better fig:2.2. Can be measured as the number of pixels measured
 130 over a distance of an inch ppi(Pixel Per Inch) or as number of line pairs that
 131 can be distinguished in a mm of image lp/mm (line pair per millimeter).
 132 • Gray-level Resolution: The range of the pixel values, a classic example is an
 133 8-bit resolution which yields 256 levels of gray. A better resolution allows a
 134 better distinction of colors within the image fig:2.2.

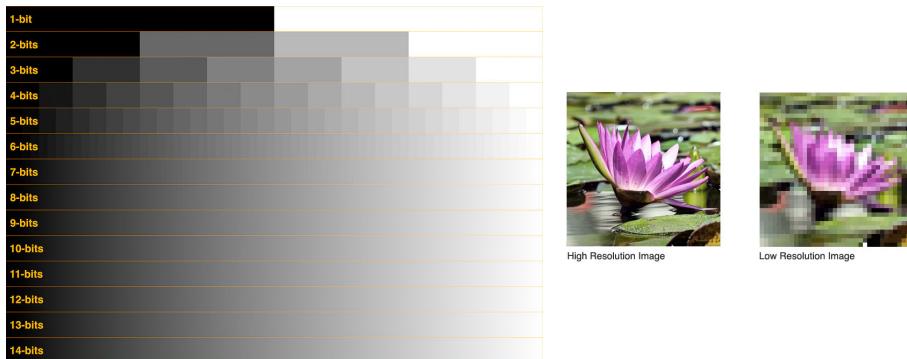


Figure 2.2: Example of visual differences in Gray-level (left) and spatial (right) resolution

- 135 • Size: Refers to the number of pixels per side of the image, for example in CT-
 136 derived images the coronal slices are usually 512x512. These numbers depend
 137 on the acquisition process and instrument but in all cases these refer to the
 138 number of rows and columns in the sampling grid as well as in the matrix
 139 representing the image.
 140 • Data-Format: How the pixel values are stored in the file of the image. The
 141 most commonly used formats are .PNG and .JPG however there are a lot of
 142 other formats. In the context of this work, which is going to be centered on
 143 medical images, the most interesting formats are going to be the nii.gz (Nifti)
 144 and the .dcm (DICOM). The first contains only the pixel value information

145 hence it's a lighter format, it originates in the field of Neuroimaging⁵, it is used
146 mainly in Magnetic resonance images of the brain but also for CT scans and,
147 since it contains only numeric information, it's the less memory consuming
148 option out of the two. The second contains not only the image data but also
149 some data on the patient, such as name and age, and details on how the exam
150 was carried out, such as machine used and specifics of the acquisition routine.
151 This format is heavier than the previous one and, for privacy purposes, is much
152 more delicate to handle which is why anonymization of the data needs to be
153 taken in consideration. For a thorough description of the DICOM standard
154 refer to [1].

155 The format in which the image is saved depends on the compression algorithm
156 used to store the information within the file. These algorithms can be lossy, in
157 which case some of the information is lost to reduce the memory needed for storage,
158 or lossless which means that all the information is kept at the expense of memory
159 space. The first set of methods is preferred for storage of natural images, these are
160 cases in which details have no importance, whereas the second set of methods is
161 used where minute details can make a considerable difference such as in the medical
162 field⁶. Given this set of characteristics it should now be clear that images can be
163 thought of as array of numbers, for this reason they are often treated as matrices
164 and, as such, there is a well defined set of valid operations and transformations
165 that can be performed on them. All these operations and transformations, in a
166 digital context⁷, are performed via computer algorithms which allow almost perfect
167 repeatability and massive range of possible operations. Given the list-like nature
168 of images one of the most natural things to do with the pixel values is to build an
169 histogram to evaluate some of the characteristic values of their distribution, such
170 as average, min/max, skewness, entropy... . The histogram of the image, albeit not
171 being an unambiguous way to describe images, is very informative. When looking
172 at an histogram it's immediately evident whether the image is well exposed and if
173 the whole range of values available is being used optimally.

⁵In fact Nifti stands for Neuroimaging Informatics Technology Initiative (NIfTI)

⁶A detailed description of compression algorithms is beyond the scopes of this thesis, for this reason please refer to [25] for more information

⁷As opposed to analog context, which would mean the chemical processes used at the start of photography to develop and modify the film on which the image was stored

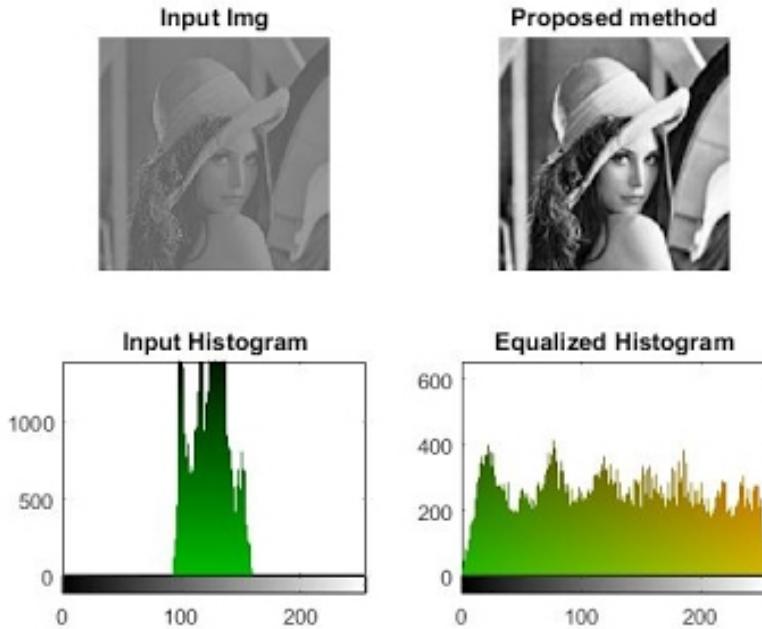


Figure 2.3: Example of differences in contrast due to histogram equalization

174 This leads us to the concept of *Contrast* which is a quantification of how well
 175 different intensities can be distinguished. If all the pixel values are bundled in a small
 176 range leaving most of the histogram empty then it's difficult to pick up the differences
 177 because they are small however, if the histogram has no preferentially populated
 178 ranges then the differences in values are being showed in the best possible way fig:2.3.
 179 Note also that if looking at the histogram there are two(or more) well separated
 180 distributions it's possible that these also identify different objects in the image, which
 181 will for example allow for some basic background-foreground distinction. Assuming
 182 they are being meaningfully used ⁸ all mathematical operations doable on matrices
 183 can be performed on images for this reason it would be useless to list them all.
 184 However, it's useful to provide a list of categories in which transformations can be
 185 subdivided:

186 1. Geometric Transformations involve the following steps:

- 187 (a) Affine transformations: Transformations that can be performed via
 188 matrix multiplication such as rotations, scaling, reflections and transla-
 189 tions. This step basically involves computing where each original pixel
 190 will fall in the transformed image
- 191 (b) Interpolation: Since the coordinates of the transformed pixel might not
 192 fall exactly on the grid it might become necessary to compute a kind
 193 of average contribution of the pixel around the destination coordinate

⁸For example adding/subtracting one image to/from another can be reasonably understood, multiplying/dividing are less obvious but still used e.g. in scaling/mask imposition and change detection respectively

194 to find a most believable value. Examples of such methods are linear,
195 nearest neighbour and bicubic.

196 2. Gray-level (GL) Transformations: Involve operating on the value stored within
197 the pixel, these can be further subdivided as:

198 (a) Point-wise: The output value at specific coordinates depends only on
199 the output value at those same specific coordinates. Some examples are
200 window-level operations, thresholding, negatives and non-linear opera-
201 tions such as gamma correction which is used in display correction. Taken
202 p as input pixel value and q as output and given a number $\gamma \in \mathbb{R}$, gamma
203 corrections are defined as:

$$q = p^\gamma \quad (2.2)$$

204 (b) Local: The output value at specific coordinates depends on a combination
205 of the original values in a neighbourhood around that same coordinates.
206 Some examples are all filtering operation such as edge enhancement, di-
207 lation and erosion. These filtering methods are based on performing con-
208 volutions in which the output value at each pixel is given by the sum of
209 pixel-wise multiplication between the starting matrix and a smaller (usu-
210 ally 3x3 or 5x5) matrix called kernel. The output image is obtained by
211 moving the kernel along the starting matrix following a predefined stride,
212 when moving near the borders the behaviour is defined by the padding of
213 the image. Stride, kernel shape and padding determine the shape of the
214 output matrix VALE LA PENA METTERE LA FORMULA E
215 SPIEGARLO MEGLIO?.

216 (c) Global: The output value at specific coordinates depends on
217 all the values of the original images. Most notable operation
218 in this category is the Discrete Fourier Transform and it's
219 inverse which allow switching between spatial and frequency
220 domains. It's worth noting that high frequency encode pat-
221 terns that change on small scales whereas low frequencies
222 encode regions of the image that are constant or slowly vary-
223 ing.

224 The aforementioned is surely not a comprehensive list of all that can be said on
225 images however it should be enough for the scopes of this work.

226 Having seen what constitutes an image and what can be done with one it
227 becomes interesting to explore how images are obtained. The following discussion
228 is going to introduce briefly some of the methods used to obtain medical images,

229 getting more in depth only on the modality used to obtain all the images used in
230 this thesis which is Computed Tomography.

- 231 1. Magnetic Resonance Imaging (MRI): This technique is based on the phe-
232 nomenon of Nuclear Magnetic Resonance(NMR) which is what happens when
233 diamagnetic atoms are placed inside a very strong uniform magnetic field are
234 subject to Radio Frequency (RF) stimulus. These atoms absorb and re-emit
235 the RF and supposing this behaviour can somehow be encoded with a posi-
236 tional dependence then it's possible to locate the resonant atoms given the
237 response frequency measured. Suffices to say that this encoding is possible
238 however the setup is very complex and the possible images obtainable with
239 this method are very different and can emphasize very different tissue/material
240 properties. Nothing more will be said on the topic since no data obtained with
241 this methodology will be used. More details can be found in [5]
- 242 2. Ultra-Sound (US): The images are obtained by sending waves of frequency
243 higher to those audible by humans and recording how they reflect back. This
244 technique is used mainly in imaging soft peripheral tissues and the contrast
245 between tissues is given by their different responses to sound and how they
246 generate echo. The main advantages such as low cost, portability and harm-
247 lessness come at the expense of explorable depth, viewable tissues, need for a
248 skilled professional and dependence on patient bodily composition as well as
249 cooperation.
- 250 3. Positron Emission Tomography (PET): In this case the images are obtained
251 thanks to the phenomenon of annihilation of particle-antiparticle, specifically
252 of electron-positron pairs. The positrons come from the β^+ decay of a radio-
253 nucleide bound to a macromolecule, which is preferentially absorbed by the
254 site of interest ⁹. Once the annihilation happens a pair of (almost) co-linear
255 photons having (almost) the same energy of 511 keV is emitted, the detection
256 of this pair is what allows the reconstruction of the image representing the
257 pharmaceutical distribution within the body. Once again there are a lot of
258 subtleties that are beyond the scopes of this thesis, suffices to say that: firstly
259 the exam is primarily used in oncology given the greater energy consumption,
260 hence nutrients absorption, of cancerous tissue and secondly this technique
261 can be combined with CT scans to obtain a more detailed representation of
262 the internal environment of the patient

263 The last technique that is going to be mentioned is Computed Tomography
264 however, given it's relevance inside this thesis work, it seems appropriate to describe
265 it in a dedicated section.

⁹Most commonly Fluoro-DeoxyGlucose FDG which is a glucose molecule labelled with a ^{18}F atom responsible of the β^+ decay. In general these radio-pharmaceuticals are obtained with particle accelerators near, or inside, the hospital that uses them. They are characterized by the activity measured as decay/s \doteq Bq (read Becquerel) and half-life $\doteq T_{\frac{1}{2}}$ which is how long it takes for half of the active atoms to decay

266 X-ray imaging and Computed Tomography (CT)

267 It's well known that the term x-rays is used to characterize a family electromagnetic
268 radiation defined by their high energy and penetrative properties. Radiation of this
269 kind is created in various processes such as characteristic emission of atoms, also
270 referred to as x-ray fluorescence, and Bremsstrahlung, braking radiation¹⁰. The
271 discovery that "A new kind of ray"[26] with such properties existed was carried out
272 by W.C.Roentgen in 1895, which allowed him to win the first Nobel prize in physics
273 in the same year. Clearly the first imaging techniques that involved this radiation
274 were much simpler than their modern counterpart, first of all they were planar and
275 analog in nature, as well as not as refined in image quality. The first CT image
276 was obtained in 1968 in Atkinson Morley's Hospital in Wimbledon. Tomography
277 indicates a set of techniques¹¹ that originate as an advancement of planar x-ray
278 imaging; these techniques share most of the physical principles with planar imaging
279 while overcoming some of it's major limitations, main of which being the lack of
280 depth information. X-ray imaging, both planar and tomographic, involves seeing
281 how a beam of photons changes after traversing a target, the process amounts to
282 a kind of average of all the effects occurred over the whole depth travelled. The
283 way in which slices are obtained is called focal plane tomography and, as the name
284 suggests, the basic idea is to focus in the image only the desired depth leaving
285 the unwanted regions out of focus. This selective focusing can be obtained either
286 by taking geometrical precautions while using analog detectors, such as screen-film
287 cassettes, or by feeding the digital images to reconstruction algorithms to perform
288 digitally the required operations¹². In both planar and tomographic setting the
289 rough description of the data acquisition process can be summarized as follows:
290 First x-rays are somehow generated by the machine, the quality of these x-rays
291 is optimized with the use of filters then focused and positioned such that they
292 mostly hit the region that needs imaging. The beam then exits the machine and
293 starts interacting with the imaged object¹³, this process causes an attenuation in
294 the beam which depends on the materials composing the object itself. Having then
295 travelled across the whole object it interacts with a sensor, be it film, semiconductor
296 or other, which stores the data that will then constitute the final image. In a digital
297 setting this final step has to be performed following a (tomographic) reconstruction
298 algorithm which given a set of 2D projections returns a single 3D image. In this
299 light the interesting processes are how the radiation is created and shaped before
300 hitting the patient and how said radiation then interacts with the matter of both
301 the patient's body and the sensor beyond it. To explore these topics it's necessary
302 to see:

- 303 • How these x-ray imaging machines are structured

¹⁰From the German terms *Bremsen* "to brake" and *Strahlung* "radiation"

¹¹from the greek *Tomo* which means "to cut" and suffix -graphy to denote that it's a technique to produce images

¹²In the first case the process is referred to as *Geometric Tomography* while in the second case as *Digital Tomosynthesis*

¹³In this work it's always going to be a patient, however this process is general and is also used in industry to investigate object construction

- 304 • How x-ray and matter interact as the first traverses the second

305 It would also be interesting to talk about reconstruction algorithms however
 306 since it's beyond the scopes of this work MAGARI DIRE CHE USER-
 307 EMO SOLO RICOSTRUZIONI DI UN TIPO QUINDI NON
 308 CI INTERESSA VEDERE LE DIFFERENZE?, refer to [13] and
 309 [30].

310 **Generation and management of radiation: digital CT scan-
 311 ners**

312 As of the writing of this thesis, seven generations of CT scanners with different
 313 technologies used. The conceptual structure of the machines is mostly the same,
 314 and the differences between generations also make evident those between machines.
 315 Exploiting this fact the structural description is going to be only one followed by a
 316 brief list of notable differences between generations.

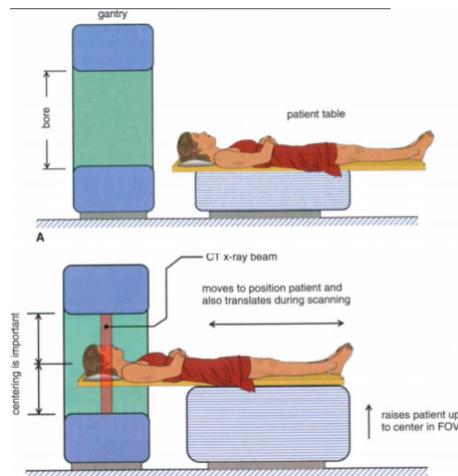


Figure 2.4: General set-up of a CT machine

317 The beam is generally created by the interaction of high energy particles with
 318 some kind of material, so that the particle's kinetic energy can be converted into
 319 radiation. In practice this means that an x-ray tube is encapsulated in the ma-
 320 chine. Inside this vacuum tube charged particles¹⁴ are emitted from the cathode,
 321 accelerated by a voltage differential and shot onto a solid anode¹⁵. This creation
 322 process implies that the spectrum of the produced x-rays is composed of the almost
 323 discrete peaks of characteristic emission, due to the atoms composing the target,
 324 superimposed with the continuum Bremsstrahlung radiation.

¹⁴Most commonly electrons

¹⁵Typical materials can be Tungsten, Molybdenum

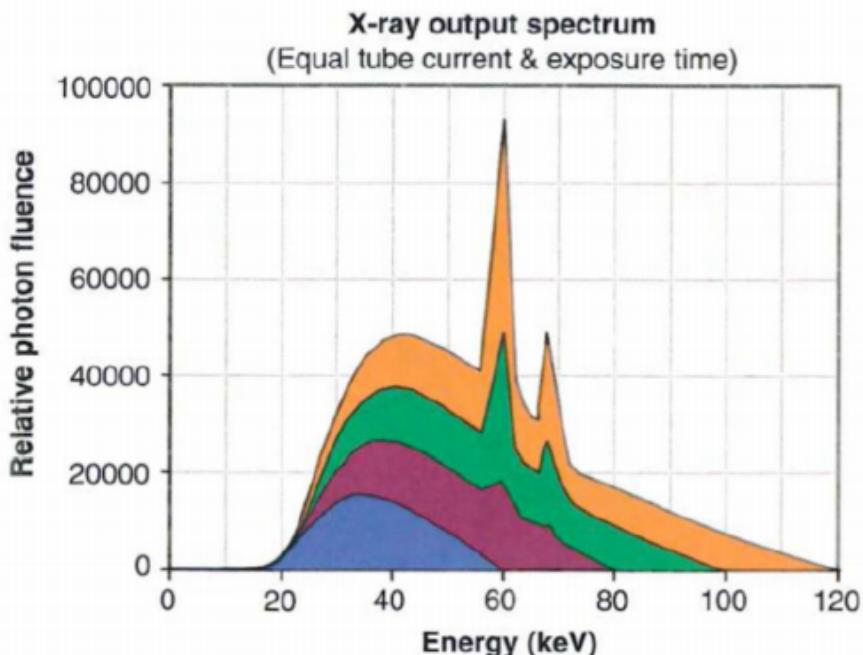


Figure 2.5: X-ray spectrum, composed of characteristic peaks and Bremsstrahlung continuum, computed at various tube voltages

Some of the main characteristics of the x-ray beam are related to this stage in the generation, the Energy of the beam is due to the accelerating voltage in the tube whereas the photon flux is determined by the electron current in the tube. Worth noting, en passant, that these two quantities can be found in the DICOM image of the exam as *kilo Volt Peak (kVP)* and *Tube current mA* and can be used to compute the dose delivered to the patient. Other relevant characteristics in the tube are the anode material, which changes the peaks in the x-ray spectrum and time duration of the emission, which is called exposure time and influences dose as well as exposure¹⁶. The electron energy is largely wasted ($\sim 99\%$) as heat in the anode, which then clearly needs to be refrigerated. The remaining energy, as said before, is converted into an x-ray beam which is directed onto the patient. To reduce damage delivered to the tissues it's important that most of the unnecessary photons are removed from the beam. Exploiting the phenomenon of beam hardening a filter, usually of the same material as the anode, is interposed between the beam and the patient to block lower energy photons from passing through thereby reducing the dose conveyed to the patient. At this point there may also be some form of collimation system which allows further shaping of the dose delivered. Having been collimated the beam traverses the patient and gets to the sensor of the machine, which nowadays are usually solid-state detectors. Naturally the description of these technologies can be done on a much finer level, however for the scopes of this work this description is deemed

¹⁶Exposure is a term used to identify how much light has gotten in the imaging sensor. Too high an exposure usually means the image is burnt, i.e. too bright and white, while lower exposures are usually associated to darker images. Exposure is proportional to the product of tube current and exposure time, measured in mA*s. Generally the machine handles the planning of exposure time according to treatment plan

345 sufficient FORSEANCHE TROPPO? O TROPPO POCO?? At
 346 this point is where the differences between generations arise which, loosely speaking,
 347 can be found in the emission-detection configuration and technology.

- 348 ● 1st generation-Pencil Beam: A single beam is shot onto a single sensor, both
 349 sensor and beam are translated across the body of the patient and then rotated
 350 of some angle. The process is repeated for various angles. Main advantages
 351 are scattering rejection and no need for relative calibration, main disadvantage
 352 is time of the exam
- 353 ● 2nd generation-fan Beam: Following the same process as the previous genera-
 354 tion the main advantage is the reduction of the time of acquisition by intro-
 355 ducing N beam and N sensors which don't wholly cover the patient's body so
 356 still need to translate.

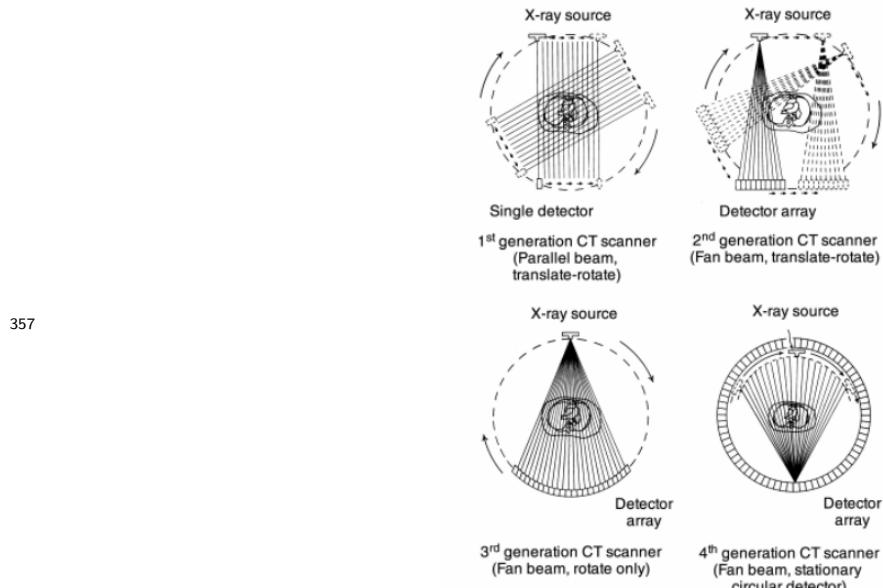


Figure 2.6: First four generation of CT scanners

- 357 ● 3rd generation-Rotate Rotate Geometry: Enlarging the span of the fan of
 358 beams and using a curved array of sensor a single emission of the N beams
 359 engulfs the whole body so the only motion necessary is rotation of the couple
 360 beam-sensor array around the patient.
- 361 ● 4th generation-Rotate Stationary Geometry: The sensors are now built to com-
 362 pletely be around the patient so that only the beam generator has to rotate
 363 around the body
- 364 ● 5th generation-Stationary Stationary Geometry: The x-ray tube is now a large
 365 circle that is completely around the patient. This is only used in cardiac
 366 tomography and as such will not be described further [12]
- 367 ● 6th generation-Spiral CT: Supposing the patient is laying parallel to the axis
 368 of rotation, all previous generations acquired, along the height of the patient,

371 a single slice at a time. In this generation as the tube rotates around the
 372 patients the bed on which they're laying moves along the rotation axis so that
 373 the acquisition is continuous and not start-and-stop. This further reduces the
 374 acquisition time while significantly complicating the mathematical aspect of
 375 the reconstruction. It's necessary to add another important parameter which
 376 is the pitch of the detector¹⁷. This quantifies how much the bed moves along
 377 the axis at each turn the tube makes around the patient. Pitches smaller
 378 than one indicate oversampling at the cost of longer acquisition times, pitches
 379 greater than one indicate shorter acquisition times at the expense of a sparser
 380 depth resolution.

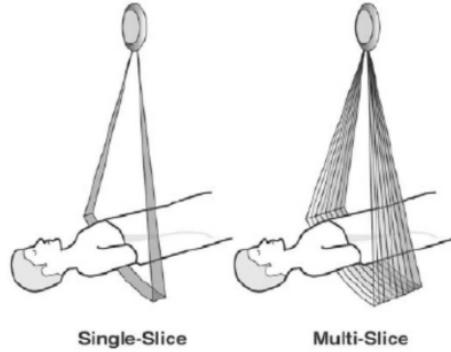


Figure 2.7: 7th generation setup

- 382 • 7th generation-MultiSlice: Up to this seventh generation height-wise slice ac-
 383 quisition was of a singular plane, be it continuous or in a start and stop
 384 motion. In this final generation mutliple slices are acquired. Considering
 385 cylindrical coordinates with z along the axis of the machine the multiple slice
 386 acquisition is obtained by pairing a fanning out along θ and one along z of
 387 both sensor arrays and beam. This technique returns to a start and stop
 388 technology in which only $\sim 50\%$ of the total scan time is used for acquisition

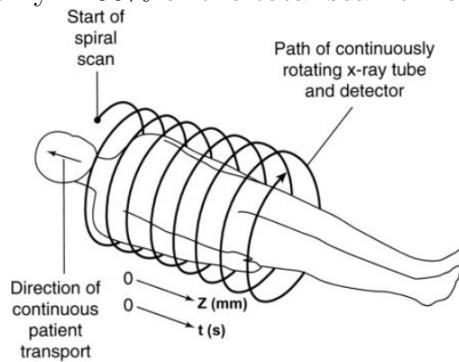


Figure 2.8: 7th generation setup

390 The machines used to obtain the images used in this thesis, all belonging to the
 391 Ospedale S. Orsola, were distributed as follows shown in 2.9:

¹⁷Once again the important parameters, such as this, can be accessed in the DICOM file resulting from the exam,

		counts	freqs
KVP	100.000	23	5,28%
	120.000	398	91,28%
	140.000	15	3,44%
	A	15	3,44%
Convolutional kernel	B	2	0,46%
	BONE	13	2,98%
	BONEPLUS	130	29,82%
	LUNG	27	6,19%
	SOFT	1	0,23%
	STANDARD	8	1,83%
	YB	29	6,65%
	YC	210	48,17%
	YD	1	0,23%
	Ingenuity CT	245	56,19%
Machine	LightSpeed VCT	179	41,06%
	iCT SP	12	2,75%
	1	257	58,94%
Slice Thickness	1,25	179	41,06%

Figure 2.9: Acquisition parameter and machine distribution

- 392 1. Ingenuity CT (Philips Medical Systems Cleveland): $\sim 56\%$ of the exams were
 393 obtained with this machine
- 394 2. Lightspeed VCT (General Electric Healthcare, Chicago-Illinois): $\sim 41\%$ of the
 395 exams in study come from this machine
- 396 3. ICT SP (Philips Medical Systems Cleveland): $\sim 3\%$ of the exams were per-
 397 formed with this machine

398 **Radiation-matter interaction: Attenuation in body and mea-
 399 surement**

400 Having seen the apparatus for data collection the remaining task is to see how the
 401 information regarding the body composition can be actually conveyed by photons.
 402 Let's first consider how a monochromatic beam of x-rays would interact with an
 403 object while passing through it. All materials can be characterized by a quantity
 404 called attenuation coefficient μ which quantifies how waves are attenuated traversing
 405 them, this energy dependent quantity is used in the Beer-Lambert law which allows
 406 computation of the surviving number of photons, given their starting number N_0
 407 and μ :

$$N(x) = N_0 e^{-\mu(E)*x} \quad (2.3)$$

408 At a microscopic level the absorption coefficient will depend on the probability
 409 that a photon of a given energy E interacts with a single atom of material. This can
 410 be expressed using atomic cross section σ as:

$$\mu(E) = \frac{\rho * N_A}{A} * (\sigma_{Photoelectric}(E) + \sigma_{Compton}(E) + \sigma_{PairProduction}(E)) \quad (2.4)$$

411 Where ρ is material density, N_A is Avogadro's number, A is the atomic weight
 412 in grams and the distinction among the various possible interaction processes for a
 413 generic photon of high energy E is made explicit. Overall the behaviour of the cross
 414 section is the following

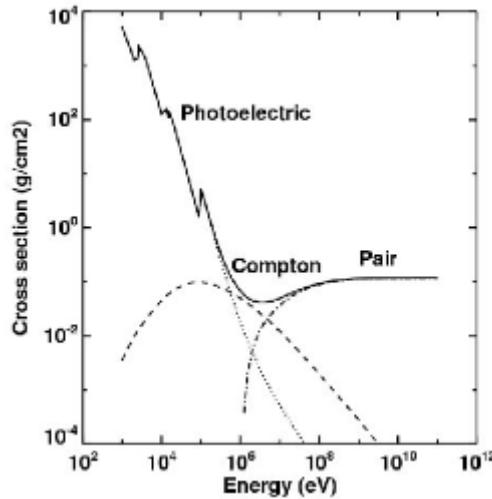


Figure 2.10: Photon cross-section in Pb

415 Overall, given eq:2.3, it's clear to see that the attenuation behaviour of a monochro-
 416 matic beam would be linear in semi-logarithmic scale hence the name for μ "Linear
 417 Attenuation Coefficient". The first complication comes from the fact that, given
 418 their generation method, the x-rays are not monochromatic but rather polychro-
 419 matic. This introduces a further complication which is the phenomenon of beam
 420 hardening: lower energy x-rays interact much more likely than those at higher en-
 421 ergies which implies that as it crosses some material the mean energy of the whole
 422 beam increases. This behaviour is exploited still within the machine, filters are in-
 423 terposed between anode and patient to reduce the useless part of the spectrum as
 424 shown in fig: 2.11a:

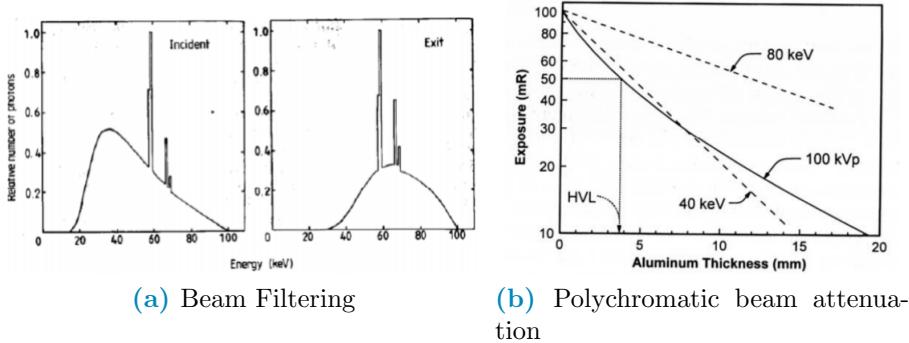


Figure 2.11: Polychromatic beam behaviour

Another effect of the beam being polychromatic is that the graphical behaviour of the attenuation instead of being linear gets bent as shown in fig: 2.11b. Since the image brightness is related to the number of photons that get on the sensor it's still possible to define the contrast between two pixel p_1, p_2 as:

$$C(p_1, p_2) = \frac{N_{\gamma,p_2} - N_{\gamma,p_1}}{N_{\gamma,p_2}} \quad (2.5)$$

This formula, that connects the beam to the image, together with eq: 2.3, which connects the beam property to the patient's composition, make clear the processes by which the beam carries patient information. Another complication arises in the context of this last equation due to the phenomenon of scattering which reduces the contrast by changing the direction of the beam and introducing an element of noise. Anti-scattering grids are positioned right before the sensors to reduce this effect by allowing to reach the sensor to only the photons with the correct direction.

The biological effects of radiation won't be treated in this thesis. Suffices to say that damage can be classified as primary, due to ionization events within the nucleus of the cell, or secondary, due to chemical changes in the cell environment. The energy deposited per unit mass is called dose and is measured in Gy(Gray) and, as said before, depends on exposure time, current and kVP of the tube. Most of contemporary machines for CT self-regulate exposure time during the acquisition automatically using Automatic Exposure Control(AEC). Having the dose it's possible to estimate the fraction of surviving cells and, to do so, various models are used. In the clinical practice it's common to find, still within the DICOM image metadata, the information regarding Dose delivered such as CTDI (Computed Tomography Dose Index) from which it's possible to obtain the DLP (Dose Length Product) taking into consideration the total length of irradiated body. For an introduction to one of these models, the Linear Quadratic (LQ) refer to [20].

449

450 451 2.1.2 Artificial Intelligence (AI) and Machine Learning(ML)

452 Having clarified the type of data that will be used in this work, and having seen
453 the general procedure used to gather it, it becomes interesting to discuss what
454 kind of techniques will be used to analyze it. Starting from the definition given by
455 John McCarthy in [18] "[AI] is the science and engineering of making intelligent
456 machines, especially intelligent computer programs. It is related to the similar
457 task of using computers to understand human intelligence, but AI does not have
458 to confine itself to methods that are biologically observable.". Machine Learning
459 (ML) is a sub-branch of AI and contains all techniques that make the computer
460 improve performances via experience in the form of exposure to data, practically
461 speaking this finds its application in classification problems, image/speech/pattern
462 recognition, clustering, autoencoding and others. The general workflow of Machine
463 Learning is the following: given a dataset, the objective is to define a model or
464 function which depends on some parameters which is able to manipulate the data
465 in order to obtain as output something that can be evaluated via a predefined
466 performance metric. The parameters of the model are then automatically adjusted
467 in steps to minimize or maximize this performance metric until a stable point at
468 which the model with the current parameters is considered finalized; one of the main
469 problems in this procedure is being sure that the stable point found is global and not
470 local. MAGARI FORSE POTREBBE ESSERE CARINO SPIEGARE UN PO' DI
471 DISCESA DEL GRADIENTE. The whole procedure is carried out keeping in mind
472 that the resulting model needs to be able to generalize its performance on data
473 that it has never seen before, for this reason usually ML is divided in a training
474 phase and a testing phase. The training phase involves looking at the data and
475 improving the performance of the model on a specific dataset¹⁸, the testing phase
476 involves using brand new data to evaluate the performance of the model obtained
477 in the preceding phase. Machine Learning techniques can be further grouped into
478 the following categories:

- 479 1. Supervised Learning: In this type of ML the model is provided with the input
480 data as well as the correct expected output, which is hence called *label*. The
481 objective of the model is to obtain an output as similar to the labels as possible,
482 while also retaining the best possible generalization ability in predicting never
483 seen before data. Some problems that benefit from the use of these techniques
484 are regression and classification problems.
- 485 2. Unsupervised Learning: As the name suggests this category of models trains on
486 the data alone, without having the labels available by minimizing some metric
487 defined from the data. For example clustering techniques try to find a set of
488 groups in the data such that the difference within each group is minimal while

¹⁸Usually in studies there is a single dataset which is split into a train-set and a test-set, in some cases if the model is good it can be validated prospectively, which means that its performance is evaluated on data that did not yet exist at the time of birth of the model

489 the difference among groups is maximal, ideally producing dense groups, called
490 clusters, that are each well separated from all the others. Other techniques in
491 this family are Principal Component Analysis (PCA) and autoencoding but
492 the general objective is to infer some kind of structure within the data and
493 the relation between data points.

494 3. Reinforcement Learning: This kind of ML is well suited for data which has
495 a clear sequential structure in which the required task is to develop good
496 long term planning. Broadly speaking the general set-up is that given a set of
497 $(\text{state}_t, \text{action}, \text{reward}, \text{state}_{t+1})$ these techniques try to maximize the cumulative
498 reward¹⁹. The main applications of these techniques are in Autonomous
499 driving and learning how to play games

500 The following methods are those directly involved in this work.

501 Regression, Classification and Penalization

502 Regression and Classification are methods used to make predictions via supervised
503 learning by understanding the input-output relation in continuous and discrete cases
504 respectively. The most basic example of regression is linear regression which consists
505 in finding the slope and intercept of a line passing through a set of points. A
506 branch of classification that's similar to linear regression is logistic regression, this
507 translates in finding out whether or not each point belongs to a certain category
508 given its properties and can be practically thought of, for a binary classification, as
509 a fitting procedure such that the output can be either 0 for one category and 1 for
510 the other, this is usually done using the logistic function, also called sigmoid, which
511 compresses \mathbb{R} in $[0, L]$ as seen in 2.12. L represents the maximum value desired,
512 x_0 is the midpoint and k is the steepness. Note that for multiple categories it would
513 conceptually suffice to add sigmoids with different centers to obtain a step function
514 in which each step corresponds to a category, in this case the procedure is usually
515 called multinomial regression.

$$\sigma(x) = \frac{L}{1 + e^{-k(x-x_0)}} \quad (2.6)$$

¹⁹i.e. the sum of the rewards obtained at all previous time steps

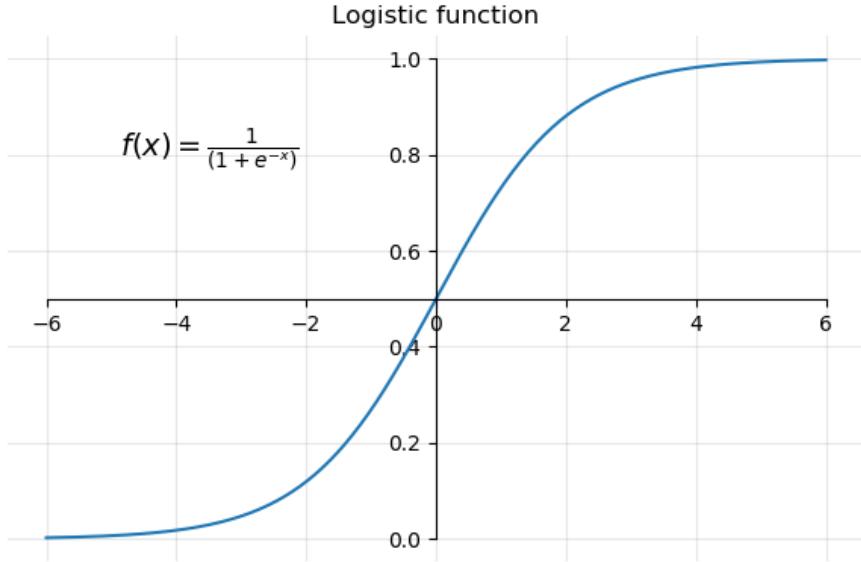


Figure 2.12: sigmoid function with L=1, $x_0=0$, k=1

516 To give a general intuition of the procedure, without going in unnecessary details,
 517 let's only consider linear regression; the theory on which it is founded is based on
 518 the assumption that the residuals, i.e. the distance model-label, are normally dis-
 519 tributed. Under this assumption the parameters can be found by changing them and
 520 trying to minimize the sum of squared residuals. When each datapoint is character-
 521 ized by a lot of different features the procedure is called multiple linear regression,
 522 the task becomes finding out how much each feature contributes in predicting the
 523 output within a weighted linear combinaion of features. Practically speaking this
 524 can be done in matricial form, supposing that each of the m datapoints has n fea-
 525 tures associated let \mathbf{X} be a matrix with $n+1^{20}$ columns and m rows and let \mathbf{Y} be
 526 a vector with m entries. Let also θ be a vector of $n+1$ entries, one for each feature
 527 plus the intercept θ_0 , then we are supposing that:

$$y_i = \sum_{j=0}^{n+1} x_{i,j} * \theta_j + \epsilon_i \Rightarrow \mathbf{Y} = \mathbf{X} * \theta + \epsilon \quad (2.7)$$

528 Where ϵ is the array of residuals of the model which, as said before, is supposed
 529 to contain values that are normally distributed. Minimizing the squared residuals
 530 corresponds to minimizing the following cost function:

$$J_\theta = (\mathbf{Y} - \mathbf{X} * \theta)^T * (\mathbf{Y} - \mathbf{X} * \theta) \quad (2.8)$$

531 By setting $\frac{\delta J}{\delta \theta} = 0$ it can be shown that the best parameters θ^* are :

²⁰ $n+1$ because we have n features but we also want to estimate the intercept of the line, so in practice the first column will be of all ones to have the correct model shape in the following matrix multiplication

$$\theta^* = (\mathbf{X}^T * \mathbf{X})^{-1} * \mathbf{X}^T \mathbf{Y} \quad (2.9)$$

532 It's evident that to obtain a result from the previous operation it's necessary that
 533 $(\mathbf{X}^T * \mathbf{X})$ be invertible, which in turn requires that there be no correlated features
 534 and that the features be less than the datapoints. To solve this problem the first
 535 step is being careful in choosing the data that goes through the regression, which
 536 may even involve some preprocessing. The second step is called Regularization, it
 537 involves adding a penalty to the cost function by adding a small quantity along the
 538 diagonal of the matrix. The nature of this small quantity changes the properties
 539 of the regularization procedure, the most famous penalties are Lasso, Ridge and
 540 ElasticNet. In practical terms the shape of the penalty determines how much and
 541 how fast the slopes relative to the features can be shrunk.

- 542 1. Ridge: Adds $\delta^2 * \sum_{j=1}^{n+1} \theta_j^2$, is called also L^2 regularization since it adds the
 543 L^2 norm of the parameter vector. This penalty can only shrink parameters
 544 asymptotically to zero but never exactly, which means that all features will
 545 always be used, even with very small contributions
- 546 2. Lasso: Adds $\frac{1}{b} * \sum_{j=1}^{n+1} |\theta_j|$, is called also L^1 regularization since it adds the L^1
 547 norm of the parameter vector. This penalty can shrink parameters to exactly
 548 zero, getting rid of the useless variables within the model.
- 549 3. ElasticNet: Adds $\lambda * [\frac{1-\alpha}{2} * \sum_{j=1}^{n+1} \theta_j^2 + \alpha * \sum_{j=1}^{n+1} |\theta_j|]$, evidently this is a midway
 550 between the Lasso and Ridge methods, where the balance is dictated by the
 551 value of α .

552 Following regularization procedures, specifically Lasso regularization, as a byprod-
 553 uct of avoiding divergences a selection and ranking of the important features is ob-
 554 tained however it's still important to preprocess the data to simplify the job of the
 555 regularization.

556 forse questo è meglio in materiali e metodi?
 557 Since the whole foundation is the normality of the residuals it's impor-
 558 tant that the data behaves somewhat nicely in this regard. Changing
 559 the data to modify the residuals is not straightforward, a proxy for
 560 this procedure is to preprocess the data by manipulating the distribu-
 561 tion of the features to make them as close to normal as possible. To
 562 this end some of the operations that can be done are:

- 563 1. Standard Scaling: This amounts to subtracting the mean of the
 564 distribution to the feature and dividing by the standard deviation
 565 so that the resulting distribution is somewhat centered around
 566 zero and has close to unitary standard deviation

567 2. Boxcox transform: When distributions are heavy-tailed, like a
568 gamma distribution would be, this finds the best exponent to
569 transform via a power-law the data. Since the exponent $\lambda \in \mathbb{R}$
570 the input data needs to be strictly positive and not constant.

571 Practical application of these procedures can be found in cap:2.
572 These considerations conclude the theoretical background on regres-
573 sion, classification and penalization thereof.

574 Decision Trees and Random Forest

575 Apart from logistic and multinomial regression there are various other
576 supervised classifying methods such as Support Vector Machines (SVM),
577 Neural Networks and Decision Trees. In this thesis, among the afore-
578 mentioned algorithms, the chosen one was a particular evolution of
579 DecisionTrees called RandomForest (RF). To provide some insight in
580 the method it's necessary to first explain how decision trees work,
581 specifying what problems they face and what are their strong points.
582 Since it's a classification method let's consider the simple case of bi-
583 nary classification with categorical²¹ features. The task of the decision
584 tree is to approximate to the best of it's abilities the labels contained
585 in the training set using all the features associated with each dat-
586 apoint, this is done by building a graph-like structure in which the
587 nodes represent the features and the links departing from it are the
588 possible values the feature takes. This graph is built in a top-down
589 approach by choosing at every step the feature that best separates the
590 data in the label categories, this process is done along each branch
591 until a node in which the separation of the preceding feature is bet-
592 ter then that provided by all remaining features or all features have
593 been considered. The first node is called root node, while the nodes
594 that have no branches going out of them are called leaves, the graph
595 represents a tree hence the name of the method. Note that at every
596 node only the subset of data corresponding to all previous feature cat-
597 egories is used. At each node the separation between the two label
598 categories c_1, c_2 due to the feature is commonly measured with Gini

²¹Categorical is to be intended as features with discrete value, opposed to continuous variables which can potentially take any value in \mathbb{R}

599 impurity coefficient, which is computed as:

$$\begin{aligned} G &= 1 - p_1^2 - p_2^2 \\ &= 1 - \left[\frac{N_{c_1 \in node}}{N_{samples \in node}} \right]^2 - \left[\frac{N_{c_2 \in node}}{N_{samples \in node}} \right]^2 \end{aligned} \quad (2.10)$$

600 The Gini impurity for a feature is then computed as an average of
601 the gini coefficients of all the deriving nodes weighted by the num-
602 ber of samples in each of the nodes. This method can be obviously
603 generalized to cases in which features are continuous by threshold-
604 ing the features choosing the value that best improves the separation
605 of the deriving node, a way to choose possible thresholds is to take
606 all the means computed with all adjacent measurements. It's impor-
607 tant to note firstly that there is no restriction on using, along differ-
608 ent branches, different thresholds for the same feature and, secondly,
609 that the same feature can end up at different depths along different
610 branches. The strength of this method is its performance on data it
611 has seen however it has very poor generalization abilities, Volendo si
612 può citare elements of statistical learning

613

614 Random Forests algorithms are born to overcome this problem. As
615 the name suggests the idea is to build an ensemble of Decision Trees
616 in which the features used in the nodes are chosen among random
617 subsamples of all the available features, the final result is obtained
618 as a majority vote over all trained trees. Each tree is also trained
619 on a bootstrapped dataset created from the original, this procedure
620 might exacerbate some problems of the starting dataset by changing
621 the relative frequency of classes seen by each tree and can be corrected
622 by balancing the bootstrap procedure. This method vastly improves
623 the performance and robustness of the final prediction while retaining
624 the simplicity and ease of interpretation of the decision trees, naturally
625 there are methods, such as AdaBoost, to deploy and precautions, such
626 as having balanced dataset, to take to further improve the performance
627 of RF classifiers. In the context of this thesis it will become necessary
628 to take care of balancing the input dataset, to do so one could ran-
629 domly oversample, by duplicating instances in the minority class, or

undersample, by removing instances within the majority class. The choice that was made was to use Synthetic Minority Oversampling TEchnique (SMOTE)[6] to rebalance the dataset.

633 Synthetic Minority Oversampling TEchnique (SMOTE)

634 This technique considers a user-defined number of nearest neighbours
 635 of randomly chosen points in the minority class and populates the
 636 feature space by generating samples on the lines that connect the
 637 chosen sample with a random neighbour²².

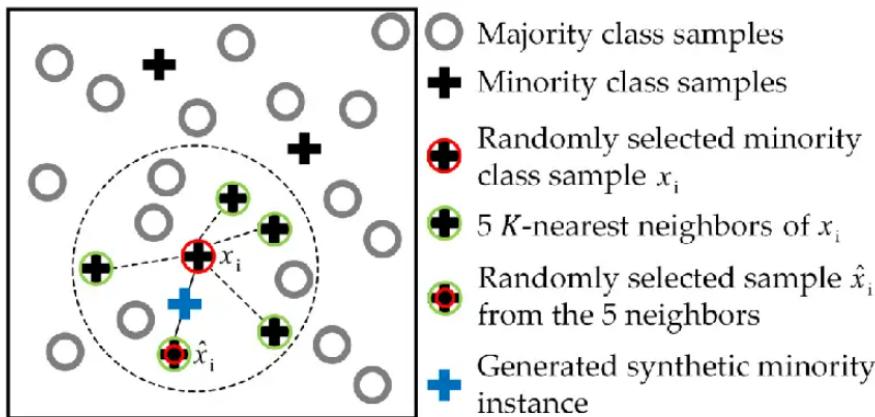


Figure 2.13: Example of SMOTE with 5 nearest neighbours

638 Worth noting that this has been done using the library imblearn in
 639 python [17]. To preview some of the data, looking at the performance
 640 of a vanilla random forest implementation with all preset parameters,
 641 the effect of the position of oversampllling is the following.

²²This procedure is formally called convex combination, which is a peculiar linear combination of vector in which the coefficients sum to one. Particularly all convex combination of two points lay on the line that connects them, and for three point lay within the triangle that has them as vertices

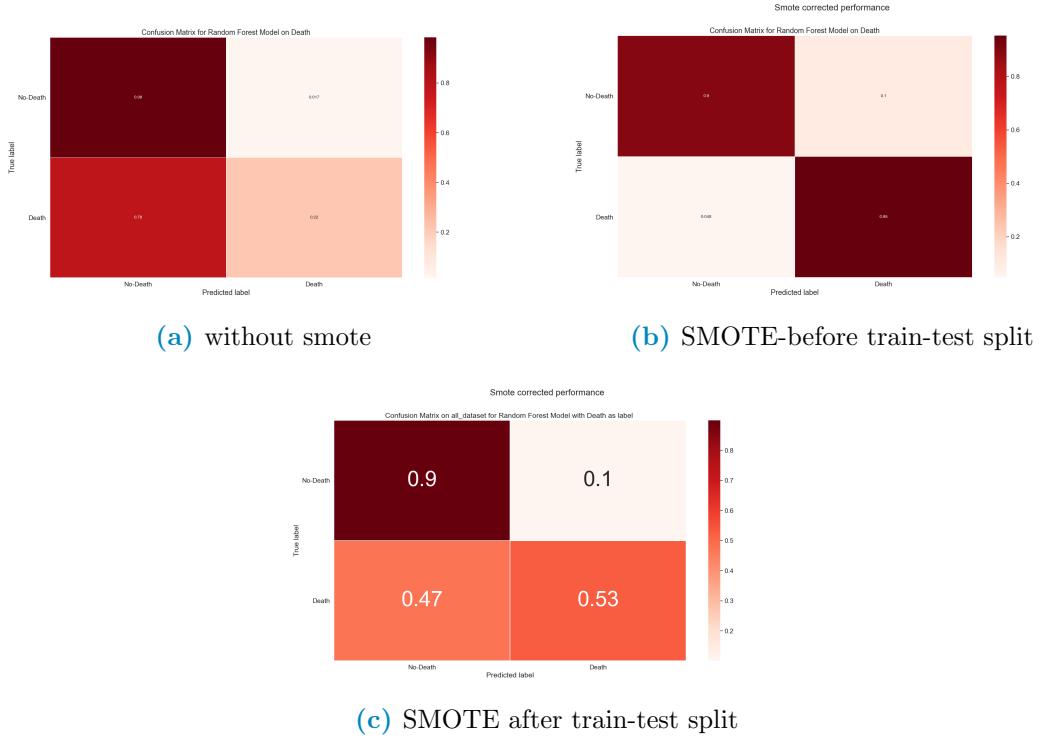


Figure 2.14: Confusion matrix used to evaluate performance done using datasets with various combinations of SMOTE position

It's clear to see that in unbalanced case, like the one analysed in this thesis in which the label classes are 15%-85%, oversampling the data always determines an improvement in performance. However performing it in the wrong place it clearly makes a far too optimistic evaluation of the performance. This highlights the point that all preprocessing should be done with care when train-test splitting the data, specifically it's very important that the oversampling, as well as all other data handling, be performed after the train-test split of the data on the train data alone. It's clear to see that what's being observed is a leakage phenomenon, in which the testing data contains information regarding the training data and viceversa. In practice, especially because the points are created as convex combination of existing data²³, when the synthetically generated points end up in the testing they depend, at least partially, on the data used in the training procedure.

²³Note that this would happen with any other over/under-sampling methods, such as random oversampling which randomly duplicates data points

656 **Dimensionality reduction and clustering**

657 When dataset are composed of many features, namely more than
658 three, it becomes difficult if not impossible to visualize the distribution
659 of data. There are various techniques that can mitigate or solve this
660 problem, they are grouped under the umbrella term of "Dimensional-
661 ity reduction techniques" and they are generally based on ML learning
662 methods. As such , following the same reasoning and definitions given
663 in regard to ML, it's possible to introduce a sub-categorization of the
664 whole family of techniques in supervised and unsupervised techniques.
665 Dimensionality reduction, beyond providing useful insight in the gen-
666 eral structure of the data, can also be used as a preprocessing step in
667 some analysis pipelines. A first example could be to find more mean-
668 ingful features before analyzing with methods such as Random Forests
669 or Regression techniques, a second example could be using the reduced
670 representation that keeps the most information possible to ease in clus-
671 tering analysis by highlighting the differences in subgroups within the
672 dataset. The most common dimensionality reduction techniques are:

673 **1. Unsupervised methods**

- 674 • Principal Component Analysis (PCA): Linearly combines the
675 pre-existing features to obtain new ones and orders them by
676 decreasing ability to explain total variance of data. The first
677 principal component is the combination of features that most
678 explains the variance in the original dataset. Given it's na-
679 ture it focuses on the global characteristics of the dataset.

- 680 • t-distributed Stochastic Neighbor Embedding (t-SNE): Keeps
681 a mixture of local and global information by using a distance
682 metric in the full high dimensional space and trying to repro-
683 duce the distances measured in the lower dimensional space
684 which can be 2- or 3-dimensional. NON SO SE VA SPIE-
685 GATO MEGLIO O PIU' IN DETTAGLIO

686 Let's define similarity between two points as the value taken
687 by a normal distribution in a point away from the centre, cor-
688 responding to the first point of interest, the same distance as
689 the two points taken in consideration. Fixing the first point
690 the similarities with all other points can be computed and

691 normalized. Doing this procedure for all points it's possible
692 to build a normalized similarity matrix. Then the whole
693 dataset is projected in the desired space, usually \mathbb{R}^i with
694 $i=1,2$ or 3 , in which a second similarity matrix is built us-
695 ing a t-distribution instead of a normal distribution ²⁴. At
696 this point, in iterative manner, the points are moved in small
697 steps in directions such that the second matrix becomes more
698 similar to the one computed in the full space.

699 2. Supervised methods

- 700 • Partial Least Squared Discriminant Analysis (PLS-DA): This
701 technique is the classification version of Partial Least Squares
702 (PLS) regression. Much like PCA the idea is to find a set
703 of orthonormal vectors as linear combination of the original
704 features in the dataset, however in PLS and PLS-DA is nec-
705 essary to add the constraint that the new component, besides
706 being perpendicular to all previous ones, explains the most
707 variability in a given target variable, or set thereof. For an in
708 depth description refer to [4] while for a more modern review
709 refer to [16] Hanno senso queste due in bibliografia?

710 3. Mixed techniques

- 711 • Uniform Manifold Approximation and Projection (UMAP):
712 Builds a network using a variable distance definition on the
713 manifold on which the data is distributed then uses cross-
714 entropy as a metric to reproduce a network with the same
715 structure in the space with lower dimension. This technique
716 maintains very local information on the data and allows a
717 complete freedom in the choice of the final embedding space
718 as well as the definition of distance metrics in the feature
719 space²⁵, it's also implemented to work an a generic pandas
720 dataframe in python so it can take in input a vast range of
721 datatypes. The math behind this method is much beyond the

²⁴The use of this t-distribution gives the name to the technique, the need for this choice is to avoid all the points bunching up in the middle of the projection space since t-distributions have lower peaks and are more spread out than normal distributions. For more details refer to [?]

²⁵Actually to keep the speed in performance the distance function needs to be Numba-jet compilable

722 scopes of this thesis, as such refer to [19] for more in depth
723 information.

724 It should be noted that dimensionality reduction is not to be taken
725 as a necessary step, however it can reduce the noise in the data by
726 extrapolating the most informative features while easing in visualiza-
727 tion and reducing computational costs of subsequent data analysis.
728 These upsides become particularly relevant in the field of clustering,
729 where the objective is to group data in sets with similar features by
730 minimizing the differences within each group while maximizing the
731 differences between different groups. Clustering techniques can be
732 roughly divided in:

- 733 1. Centroid based techniques: A user defined number of points is
734 randomly located in the data space, datapoints are then assigned
735 to groups according to their distance from the closest center. The
736 main technique in this category is k-means clustering, and some,
737 if not most, other techniques include it as step in the processing
738 pipeline²⁶. The main problems are firstly that these techniques
739 require prior knowledge, or at the very least a good intuition, on
740 the number of clusters in the dataset while also assuming that the
741 clusters are distributed in spherical gaussian distribution, which
742 is not always the case.
- 743 2. Hierarchical clustering: The idea is to find the hierarchical struc-
744 ture in the data using a bottom-up or a top-down approach, as
745 such this category further subdivides in agglomerative and di-
746 visive methods. In the first each point starts by itself and then
747 points are agglomerated using a similarity or distance metric, this
748 build a dendrogram in which the k^{th} level roughly corresponds to
749 a k -centroid clustering. In the latter the idea is to start with a
750 unique category and then divide it in subgroups.
- 751 3. Density based techniques: These methods use data density to
752 define the groups by looking for regions with larger and lower
753 density as clusters and separations.

²⁶An example of such techniques is: affinity propagation

In order to evaluate the performance of these methods it's necessary to define metrics that evaluate uniformity within clusters and separation among them, the choice in the definition of metric should be taken in careful consideration since the different task may imply very different optimal metrics or, put differently, the optimal technique for the task at hand may very well depend on the metric used. It's worth mentioning that when working in two, or at most three, dimensions humans are generally good at performing clustering yet it's nearly impossible to do in more dimensions. Computers, on the other hand, require more careful planning even in low dimension but are much more performing even in higher dimensions because the generally good human intuition on the definition of cluster is not so easily translated in instruction to a machine. So to obtain good results with clustering techniques it's necessary to work with care, especially with a good understanding of the dataset in use and its overall structure. In the context of this thesis, supposing the data suggested a clear cut distinction of two populations in the dataset, as could be male-females or under- vs normal- vs over-weight individuals, then it might become necessary to analyse these groups as different cohorts in order to more accurately predict their clinical outcome.

2.1.3 Combining radiological images with AI: Image segmentation and Radiomics

Having seen the kind of data that will be of interest throughout this thesis, and having a set of techniques used to describe and make prediction on the data at hand, the final step in this theoretical background chapter will be to combine these two notions in describing first how images can be treated in general terms and then, more specifically, how medical images can be analysed to exploit as much as possible the vast range of information they contain. Image analysis seems, at first glance, very intuitive since for humans it's very easy to infer qualitative information from images. However upon closer inspection this matter becomes clearly non-trivial due to the subjectivity involved in the process as well as in the intuition behind it. More specifically, in the context of this thesis, the same image of damaged lungs contains very different informations to the eyes of trained professionals versus

789 those of an ordinary person as well as to the eyes of different profes-
790 sionals. The first big obstacle in this task is the definition of region of
791 interest: not all people will see the same boundary in a damaged or-
792 gan, sometimes the process of defining a boundary between organ and
793 tissue may need to account for the final objective it has to achieve. If
794 the objective is to evaluate texture of a damaged lung then the lesion
795 needs to be included whereas in other cases these regions may only be
796 unwanted noise. Generally speaking finding regions of interest in an
797 image is a process called image segmentation. The next step would
798 be to quantify the characteristics of the region identified, as such it
799 should be clear that finding ways to derive objective information from
800 images it's of paramount importance, especially when this information
801 can aid in describing the health of a patient. It should also be clear
802 that medical images are a kind of high dimensional data as such, as it's
803 fashion with fields that occupy themselves with big biological data, the
804 field that studies driving quantitative information out of radiological
805 images is called radiomics.

806 **Image Segmentation**

807 Generally speaking image segmentation is a procedure in which an
808 image is divided in smaller sets of pixels, such that all pixel inside a
809 certain set have some common property and such that there are no
810 overlaps between sets. These sets can then be used for further analysis
811 which could mean foreground-backgroud distinction, edge detection
812 as well as object detection, computer vision and pattern recognition.
813 Image segmentation can be classified as:

- 814 • Manual segmentation: The regions of interest are manually de-
815 fined usually by a trained individual. The main advantage of this
816 it's the versatility, on the other hand this process can be very
817 time consuming
- 818 • Semi-automatic segmentation: A machine defines as best as it
819 can the shape of the region of interest, however the process is
820 then thought to receive intervention of an expert to correct the
821 eventual mistakes or refine the necessary details. This provides

822 the best compromise between time needed and accuracy obtained
823 and becomes of interest in fields in which finer details are impor-
824 tant.

- 825 • Automatic segmentation: A machine performs the whole segmen-
826 tation without requiring human intervention

827 On a practical level these techniques can be used in various fields, as
828 illustrated by the variety of aforementioned tasks, but the one that in-
829 terests this work the most is the medical field. Nowadays in medicine,
830 where most of imaging exams are stored in digital form, the abil-
831 ity to automatically discern specific structures within the images can
832 provide a way to aid clinical professionals in their everyday decision
833 making their workload lighter and helping in otherwise difficult cases.
834 A staggering example connected to this thesis is the process of organ
835 segmentation in CT scans: usually these scans are n^{27} stacked images
836 in 512x512 resolution. To have a contour of the lungs in a chest CT a
837 radiologist would need to draw by hand the contour of the lung in each
838 slice of the scan. Even if some shorcuts exist to reduce the number
839 of slices to draw on this very boring, time consuming and repetitive
840 task can occupy hours if not days of work to a human while a machine
841 can take minutes to complete a whole scan. Then considering lesion
842 detection in medical exams having a machine that consistently finds
843 lesions that would otherwise be difficult to discern by a human eye
844 can be of paramount importance in diagnosis as well as treatment.
845 In the medical field the difficulty comes from the fact that different
846 exams have different types of image formats which means that an al-
847 gorithm that works well on CT may not work as intended on MRI or
848 other procedures. Automatic image segmentation can be performed
849 in various ways:

- 850 1. Artificial Neural Network (ANN): These techniques belong to a
851 sub-field in ML called Deep Learning, they involve building net-
852 work structures with more layers in which each node is a process-
853 ing unit that takes a combination of the input data and gives an
854 output according to a certain activation function. These struc-
855 tures are called Neural Networks because they resemble and are

²⁷ n clearly depends on the exam required and slice thickness, some common values for thoracic CTs are around 200-300 but can range up to 900 slices

modeled after the workings of neuron-dendrite structures while the deep in deep learning refers to the fact that various layers composed of various neurons are stacked one after the other to complete the structure. Learning is obtained by changing how each neuron combines the inputs it receives. In the case of images these structures are called Convolutional Neural Networks because the first layers, which are intended to extract the latent features or structures in the images, perform convolution operation in which the parameters are the values pixel values of convolutional kernels

2. Thresholding: These techniques involve using the histogram of the image to identify two or more groups of pixel values that correspond to specific parts/objects within the image. An obvious case would be a bi-modal distribution in which the two sets can be clearly identified but there are no requirements on the histogram shape. In the case of CT this has a very simple and clear interpretation since HU depend on tissue type it's reasonable to expect that some tissues can be differentiated with good approximation by pixel value alone
3. Deformable models and Region Growing: Both these techniques involve setting a starting seed within the image, in the first case the seed is a closed surface which is deformed by forces bound to the region of interest, such as the desired edge. In the second case the seed is a single point within the region of interest, step by step more points are added to a set which started as the seed alone according to a similarity rule or a predefined criteria.
4. Atlas-guided: By collecting and summarizing the properties of the object that needs to be segmented it's possible to compare the image at hand with these properties to identify the object within the image itself.
5. Classifiers: These are supervised methods that focus on classifying via by focusing on a feature space of the image. A feature space can be obtained by applying a function to the image, an example of feature space could be the histogram. The main distinction from other methods is the supervised approach

891 6. Clustering: Having a starting set of random clusters the procedure
892 computes the centroids of these clusters, assigns each point to the
893 closest cluster and recomputes centroids. This is done iteratively
894 until either the point distribution or the centroid position doesn't
895 change significantly between iterations.

896 For a more in depth review of the main methods used in medical
897 image segmentation refer to [24]. Worth noting, at this point, that the
898 semi-automatic segmentation software used in this work uses probably
899 a mixture of region growing and thresholding methods, maybe guided
900 by an atlas. The segmentation process generally produces a boolean
901 mask which can be used to select, using pixel-wise multiplication, the
902 region of interest in the image. The next step in image analysis would
903 be to derive information from the region defined during segmentation,
904 in general this step is called feature extraction²⁸ and it's objective is to
905 find non-redundant quantities that meaningfully summarize as much
906 properties of the original data as possible.

907 Radiomics

908 When the images are medical in nature and when referring to high-
909 throughput quantitative analysis the task of finding these features fall
910 in the realm of radiomics, which uses mathematical tools to describe
911 properties of the images that would otherwise be unquantifiable to
912 the human eye. The features that can be computed from images are
913 of various types, some of them can be understood somewhat easily
914 through intuition since they are close to what humans generally use
915 to describe images, others are much more complex in definition and
916 quantify more difficulty perceived properties of the image. Features
917 can be then roughly classified in different families:

918 1. Morphological features: These features describe only the shape
919 of the region of interest, as such they are independent of the pixel
920 values inside the region and hence, to be computed, require only
921 a boolean mask of the segmented region. These features can be

²⁸Even if the term is used in pattern recognition as well as machine learning in general

further subcategorized as two or three dimensional features based on whether they focus on single slices or whole volumes. Most of these features compute volumes, lengths, surfaces and shape properties such as sphericity, compactness, flatness and so on.

2. First order features: These features depend strictly on the gray levels within the region of interest since they evaluate the distribution of these values, as such they need that the boolean mask of the segmentation be multiplied pixel-wise with the origian image to obtain a new image with only the interesting part in it. Most of these features are commonly used quantities, such as Energy, Entropy, Minimum and Maximum value which have been adapted to the imaging context using the histogram of the original image or by considering intensities within an enclosed region.

3. Higher order features: All the other fatures fall in this macro-category which can be clearly subdivided in other smaller categories in which the features are obtained following the same guiding principle or starting point. Generally these describe more texture-like properties of the image and, to do so, use particular matrices derived from the original image which contain specific information regarding order and relationships in pixel value positioning within the image. These matrices have very precise definition, as such only the general idea behind them will be reported here redirecting to [32] for a more strict and in detail description. The matrices from which the features are computed also give name to the smaller categories in this family, these categories are:

- Gray Level Co-occurrence Matrix (GLCM) features: This matrix expresses how combination of pixel values are distributed in a 2D or 3D region by considering connected all neighbouring pixel in a certain direction with respects to the one in consideration. Using all possible directions for a set distance, usually $\delta=1$ or $\delta=2$, various matrices are obtained and from these a probability distribution can be built and evaluated. It should be noted that before computing these matrices the intensities in the image are discretized.
- Gray Level Run Length Matrix (GLRLM) features: Much like before the task of these features is to quantify the distri-

bution of relative values in gray levels trhoughout the image, as the name suggests what this matrix quantifies is how long a path can be built by connecting pixel of the same value along a single direction. This time information from the matrices computed by considering different directions are aggregated in different ways to improve rotational invariance of the final features

- Gray Level Size Zone Matrix (GLSZM) features: This matrix counts the number of zones in which voxel have the same discretized gray level. The zones are defined by a notion of connectedness most commonly first neighbouring voxel are considered as connectable if they have the same value, this leads to a 26 neighbouring voxel in 3 dimensions and to 8 connected pixels in 2 dimensions²⁹. The matrix contains in position (i,j) the number of zones of size j in which pixel have value i.
- Neighbouring Gray Tone Difference Matrix (NGTDM) features: Born as an alternative to GLCM these features rely on a matrix that contains the sum of differences between all pixel with a given pixel value and the average of the gray levels in a neighbourhood around them.
- Gray Level Dependence Matrix (GLDM) features: The aim of these features is to capture in a rotationally invariant way the texture and coarseness of the image. This matrix requires the already seen concept of connectedness with a given distance as well as dependence among pixel. Two voxel in a neighbourhood are dependent if the absolute value of the difference between their discretized value is less then a certain threshold. The number of dependent voxel is then counted with a particular approach to guarantee that the value be at least one

In talking about the previous feature groups the concept of discretization of data which is already digital, and hence a discretized, has emerged. This is often a required step to make the computations

²⁹To visualize, imagine a 2D grid: the 8 pixel are the four at the sides of each square and the four at the corners.

of the matrices tractable and consists in further binning together the pixel values, which is commonly done in two main ways: either the number of bins is fixed or the width of the bins is fixed preceding the discretization process. It should be evident that the results of the feature extraction procedure is heavily dependent on the choices made in all the steps that precede it, main of which being the segmentation, the eventual re-discretization of the image leading to the algorithm used to compute the feature themselves. For this reason recently the International Biomarker Standardization Initiative (IBSI [32]) wrote a "*reference manual*" which details in depth the definitions of the features, description of data-processing procedures as well as a set of guidelines for reporting results. This was done in an attempt to reduce as much as possible the variability and lack of reproducibility of radiomic studies. In the past the main attention in radiological research was focused on improving machine performances and evaluating acquisition sequence technologies, however the great developments in artificial intelligence and performance of computers have brought a lot of attention to the field. Various papers, such as [15] and [3] have been written with the objective of presenting the general workflow. By design the topics in this thesis have been presented to resemble the general order of the pipeline, which can be summarised as:

- Data acquisition
- Definition of the Region Of Interest (ROI)
- Pre-processing
- Feature extraction
- Feature selection
- Classification

Another interesting possibility offered by the biomarkers computed following radiomics is the ability to quantify the differences between successive exams of the same patient. This specific branch of radiomics is called Delta-radiomics, referencing to the time differential that become the main focus of the analysis.

NON SAPREI SE VA BENE COSÌ' O BISOGNA AGGIUNGERE

1026 ALTRI DETTAGLI

1027

1028 With this the theoretical background has been laid out and the the-
1029 sis can finally proceed to practically and more specifically describing
1030 the data at hand as well as the specific techniques used to elaborate
1031 it.

1032 **2.2 Data and objective**

1033 The objective of this thesis will be to compare how different methods
1034 perform in predicting clinical outcomes in covid patients, while also de-
1035 termining if different kind of input data imply different performances
1036 of the same methods. All images are non segmented, as such all of
1037 them are going to be semi-automatically segmented via a new soft-
1038 ware being tested in the medical physics department called *Sophia Ra-*
1039 *domics*, which will be briefly explained shortly. Statistical analysis are
1040 going to be performed mainly in python using libraries such as scikit-
1041 learn, pandas, numpy, scipy while the graphical part of the analysis
1042 is done with either seaborn or matplotlib. The starting dataset was a
1043 list of all the patients that, from 02/2020 to 05/2021, were hospital-
1044 ized as COVID-19 positive inside the facilities of *Azienda ospedaliero-*
1045 *universitaria di Bologna - Policlinico Sant'Orsola-Malpighi*. As far as
1046 exclusion criteria go the main exclusion criteria, except unavailability
1047 of the feature related to the patient, was visibly damaged or lower
1048 quality images, for example images with cropped lungs. The first set
1049 of selection criteria were:

- 1050 • All patients that had undergone a CT exam which was retriev-
1051 able via the PACS (Picture Archiving and Communication Sys-
1052 tem) of *Azienda ospedaliero-universitaria di Bologna - Policlinico*
1053 *Sant'Orsola-Malpighi*
- 1054 • All patients that had all of the clinical and laboratory features,
1055 listed in fig:2.15, suggested by Lidia Strigari ³⁰

³⁰She is, as of the writing of this thesis, the head of the Medical Physics department in the S. Orsola hospital. These suggestions were given to her by clinical professionals.

- 1056 • Since all patient had at least 2 CT exams only the closest date
 1057 to the hospital admission date was taken. When more exams
 1058 were performed on the same date all of them were initially taken.
 1059 At first only chest or abdomen CTs were taken regardless of the
 1060 acquisition protocol used.

Feature	counts	freqs	categories	Feature2	counts3	freqs4	categories5	Feature7	counts8	freqs9	categories10
Obesity	363	83%	0	HRCT performed	479	100%	1	MulBSTA score total	6	1%	0
Obesity	73	17%	1	High Flow Nasal Cannulae	382	88%	0	MulBSTA score total	7	2%	2
qSOFA	226	52%	0	High Flow Nasal Cannulae	54	12%	1	MulBSTA score total	7	2%	1
qSOFA	178	41%	1	Bilateral Involvement	33	8%	0	MulBSTA score total	50	11%	5
qSOFA	29	7%	2	Bilateral Involvement	403	92%	1	MulBSTA score total	5	1%	6
qSOFA	3	1%	3	Respiratory Failure	231	53%	0	MulBSTA score total	72	17%	7
SOFA score	28	6%	0	Respiratory Failure	205	47%	1	MulBSTA score total	9	2%	8
SOFA score	114	26%	1	DNR	413	95%	0	MulBSTA score total	111	25%	9
SOFA score	144	33%	2	DNR	23	5%	1	MulBSTA score total	1	0%	10
SOFA score	72	17%	3	ICU Admission	359	82%	0	MulBSTA score total	61	14%	11
SOFA score	42	10%	4	ICU Admission	77	18%	1	MulBSTA score total	7	2%	12
SOFA score	23	5%	5	Sub-Intensive care unit admission	336	77%	0	MulBSTA score total	70	16%	13
SOFA score	7	2%	6	Sub-intensive care unit admission	100	23%	1	MulBSTA score total	17	4%	15
SOFA score	3	1%	7	Death	358	82%	0	MulBSTA score total	2	0%	16
SOFA score	3	1%	8	Death	78	18%	1	MulBSTA score total	8	2%	17
CURB65	141	32%	0	Fever	186	43%	0	MulBSTA score total	1	0%	18
CURB65	141	32%	1	Febbre	250	57%	1	MulBSTA score total	2	0%	19
CURB65	120	28%	2	Sex	151	35% Female	Lung consolidation	211	48%	0	
CURB65	30	7%	3	Sex	284	65% Male	Lung consolidation	225	52%	1	
CURB65	4	1%	4	Sex	1	0%	Hypertension	194	45%	0	
MEWS score	27	6%	0	Ground-glass	54	12%	0	Hypertension	242	56%	1
MEWS score	143	33%	1	Ground-glass	382	88%	1	History of smoking	348	80%	0
MEWS score	127	29%	2	Crazy Paving	337	77%	0	History of smoking	88	20%	1
MEWS score	82	19%	3	Crazy Paving	99	23%	1	NIV	371	85%	0
MEWS score	33	8%	4	O2-therapy	66	15%	0	NIV	65	15%	1
MEWS score	13	3%	5	O2-therapy	370	85%	1				
MEWS score	10	2%	6	cPAP	368	84%	0				
MEWS score	1	0%	7	cPAP	68	16%	1				

Figure 2.15: Description of clinical label dataset, in sets of four columns there's the clinical feature, the total count of occurrences, the percentage over the final dataset and the possible values the feature could take

1061 Most clinical features are pretty self-explanatory, for those that
 1062 were obscure to me as an outsider and not otherwise explained in ??,
 1063 I'm going to provide a very simple explanation:

1. DNR : Acronym for "Do Not Resuscitate", used to indicate the wish of the patient or their relatives that cardiac massage not be performed in case of cardiac arrest.
2. NIV: Acronym for "Non Invasive Ventilation", it's a form of respiratory aid provided to patients.
3. cPAP: Acronym for "continuous Positive Airway Pressure", another form of respiratory aid.
4. ICU: Acronym for "Intensive Care Unit". When patients are in really severe conditions they are treated in these facilities.

1073 5. Clinical Scores: When available values from laboratory analyses
1074 and/or patient conditions are summarised in scores that represent
1075 the gravity of the state of the patient, as such these can be some-
1076 what correlated and will be treated as comprehensive values to
1077 substitute an otherwise large set of obscure clinical features. At
1078 admission, or closely thereafter, a set of clinical questions regard-
1079 ing the patient receives a yes or no answer, each answer has an
1080 additive contribution towards the final value of the score. These
1081 scores differ in how much they add for each condition and the set
1082 of symptoms the check for.

- 1083 (a) MulBSTA: This score accnts for **M**ultilobe lung involve-
1084 ment, absolute **L**ymphocyte count, **B**acterial coinfection, his-
1085 tory of **S**moking, history of hyper**T**ension and **A**ge over 60
1086 yrs. [10]
- 1087 (b) MEWS: Modified Early Warning Score for clinical deterio-
1088 ration. Computed considering systolic blood pressure, heart
1089 rate, respiratory rate, temperature and AVPU(Alert Voice
1090 Pain Unresponsive) score. [28]
- 1091 (c) CURB65: **C**onfusion, blood **U**rea Nitrogen or Urea level,
1092 **R**espiratory Rate, **B**lood pressure, age over **65** years. This
1093 score is specific for pneumonia severity [31]
- 1094 (d) SOFA: **S**equential **O**rgan **F**ailure **A**sessment score. Consid-
1095 ers various quantities from all systems to assess the overall
1096 state of the patient, $\text{PaO}_2/\text{FiO}_2^{31}$ for respiratory system,
1097 Glasgow Coma scale³² for nervous, mean pressure for cardio-
1098 vascular, Bilirubin levels for liver, platelets for coagulation
1099 and creatine for kidneys [2]
- 1100 (e) qSOFA: quick SOFA. Only considers pressure, high respi-
1101 ratory rate and the low values in the Glasgow scale.

1102 This procedure produced a starting cohort of ~ 700 patients which,
1103 having all various images available, created a huge set of ~ 2200 CT
1104 scans. Since this analysis is focused on radiomics there is an evident

³¹Very unrefined yet widely used indicator for lung dysfunction

³²GCS for short, proposed in 1974 by Graham Teasdale and Bryan Jennet. Evaluates what kind of stimulus is necessary to obtain motor and verbal reactions in the patient as well as what's necessary for the patient to open their eyes

need for as much consistency as possible in the images analysed. For this reason all CTs taken with medium of contrast were excluded, since they would have brightnesses not indicative of the disease, and for every patient only images with thin slice reconstruction were considered. More specifically only images with slice thickness of 1 o 1.25 mm³³ along the z-axis were taken into consideration, which meant excluding all the 1.5,2,2.5 and 5 mm slice thicknesses. Since all these images were segmented by me and two other students using a semiautomatic segmentation tool provided by the hospital it sometimes happened that manual corrections were necessary, these were performed only in very obvious and simple cases while, in all remaining cases, the patient was dropped out of the study. Overall this left the final study cohort to be composed of 435 patients, all descriptions and analyses are related to this cohort. The same software used for segmentation allowed the extraction of the radiomic features form the segmented volumes, even if it did not allow the extraction of the segmentation masks nor any changes in the segmentation parameters, which seems a region growth algorithm mixed with thresholding. For this reason all the image analysis in this thesis is reliant on said software which has been treated as a black-box. NON SO SE POSSO SCRIVERE CHE IL PROGRAMMA ERA DELUDENTE NE' QUANTO POSSO SBILANCIA RMI A DARE INFO CHE NON HO SULLA SEGMENTAZIONE, SE CON LA STRIGARI HANNO PUBBLICATO L'ARTICOLO SULLA IBSI-COMPLIANCE SI POTREBBE CITARE QUELLO So, having segmented all the images and extracted all the features supported in the software, the next step is the definition of the actual analysis pipeline. The whole dataset was comprised of ~200 features, which were divided in three subgroups as follows:

1. Clinical: All these features are derived from the admission pro-

³³This meant that only exams called 'Parenchima' or 'HRCT' were included. Throughout the internship 'parenchima' has always appeared in contrast with 'mediastino'. These two keywords are used in the phase of reconstruction of the raw data to identify reconstructions with specific properties. Parenchima is used for finer reconstruction of lung specifically, the requiring professional uses these images to look for small nodules with very high contrast and, to do so, the reconstruction allows some noise to achieve the best resolution possible. Mediastino is used in the lung, as well as other regions, to look for bigger lesions but with low contrast. As such the 'mediastino' reconstruction compromises a worse spatial resolution for a better display of contrast, visually speaking the first images are more coarse and noisy while the second are smoother. It should be noted that even with the same identifier, be it HRCT parenchima or others, the machines on which the exams were made were different and had different proprietary convolutional kernels used for reconstruction.

1134 cedure in the hospital.

- 1135 • The continuous are Age taken at the date of the CT exam and
1136 Respiratory Rate defined as number of breaths in a minute.
 - 1137 • The discrete one were the aforementioned scores and the
1138 boolean ones were sex of the patient, obesity status, if the
1139 patient had a fever³⁴ as of hospital admission and whether or
1140 not the patient suffered of Hypertension
 - 1141 • The remaining features, namely those in 2.2 as well as the
1142 death status of the patient, were either used as labels or
1143 not used at all because they refer to treatments used and
1144 not characteristics of the patient. As such these features,
1145 while plausibly correlated to the clinical outcome, are not
1146 really descriptive of the patient as of admission and are not
1147 information that can be used to aid professionals at admission
1148 to assess the situation
- 1149 2. Radiomic: These features were all the ones supported by the seg-
1150 mentation software and are pretty much most of those described
1151 in [32] with the addition of fat and muscle surface, computed as
1152 cm² by counting pixel identified via threshold as fat or muscle
1153 tissue in thoracic slices taken at height of vertebra T-12
- 1154 3. Radiological: These features are those that can be derived from
1155 CT exams by humas. Namely acquisition parameters, such as
1156 KVP and Current, were used to search for eventual correlations
1157 between image quality and predictive power of the feature derived
1158 from the image while boolean features, such as Bilaterality of
1159 lung damage, presence of Groun Glass Opacities (GGO), lung
1160 consolidations as well as crazy paving were used to see if they
1161 were sufficient in determining outcome.

1162 2.3 Preprocessing and data analysis

1163 Before any preprocessing a choice was made to exclude all of the clin-
1164 ical scores, this was done to avoid having them mask other, more

³⁴Defined as body temperature>38°

1165 straightforward variables. **QUESTO ANDRA MESSO IN UNA
1166 FOOTNOTEIn APPENDICE SE MAI RIESCO A SCRIVERLA**
1167 an alternative choice will be made, namely to keep only the most
1168 strongly predicting one out of all of them. The first step in the anal-
1169 ysis of this data is going to be a lasso regularized regression using
1170 either death or ICU admission as target. As mentioned before when
1171 operationg with regressions it's a necessity that the residuals be nor-
1172 mally distributed, a common way to get as close as possible to this
1173 hypothesis is to boxcox transform the data. Apart from the data
1174 which contained negative values, which cannot be fed into the boxcox
1175 transform, all variables have been transformed using this method.

1176 The quatile-quantile plots have been used for a visual check of nor-
1177 mality, normally distributed data will populate the bisector of the
1178 graph while deviations are symptoms of non-normality. Heavy and
1179 light tails are visible as deviations respectively above and belox the bi-
1180 sector. It should be noted that when the data is normally distributed
1181 then the transform doesn't change much the distribution while, in
1182 cases with more heavy tailed distributions, the improvement is clear
1183 to see as it can be seen in Figures 2.16 and 2.17 The next prepro-
1184 cessing step has been to apply a *StandardScaler* to all of the features,
1185 which corresponds to subtracting the mean and dividing by the stan-
1186 dard deviation, in order to center all the features around zero. The
1187 final step in preprocessing has been to reduce the features by using
1188 a correlation threshold which means that all variables that correlate
1189 with another more, in absolute value, than a certain threshold, which
1190 has been set to 0.6 in this work, are dropped a priori. A rather im-
1191 portant thing to notice is that correlation has been computed using
1192 Spearman correlation and not Pearson since the first is invariant un-
1193 der monotone transformations, such as boxcox, and the second is not.
1194 Given the large number of features it's very plausible that at least
1195 one of the eliminated features is correlated with all other dropped fea-
1196 tures but with none of the remaing ones, since a Lasso regularization
1197 will be used introducing a few redundant features is not too damag-
1198 ing and the possible benefits outweigh the risks. For this reason a
1199 redrawing method has been implemented to add one of the dropped
1200 features, this has been done by choosing the one that most correlates
1201 with the label being used. A pivotal point in all of this analysis is

Distribution and q-q plot of Angular second moment

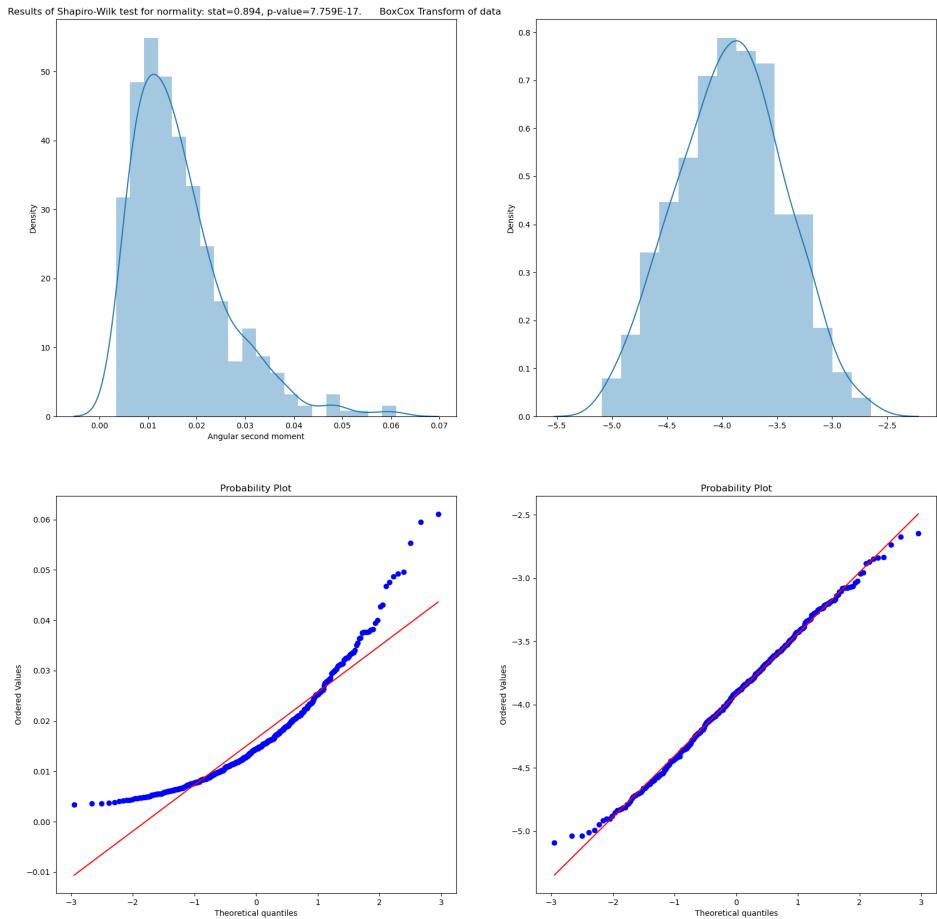


Figure 2.16: Example of boxcox applied to the radiomic feature *Angular second moment* which has a heavy tailed distribution. The graphs contain the original distribution (top-left) the transformed distribution (top-right) and the two respective quantile-quantile plots

Distribution and q-q plot of Difference average

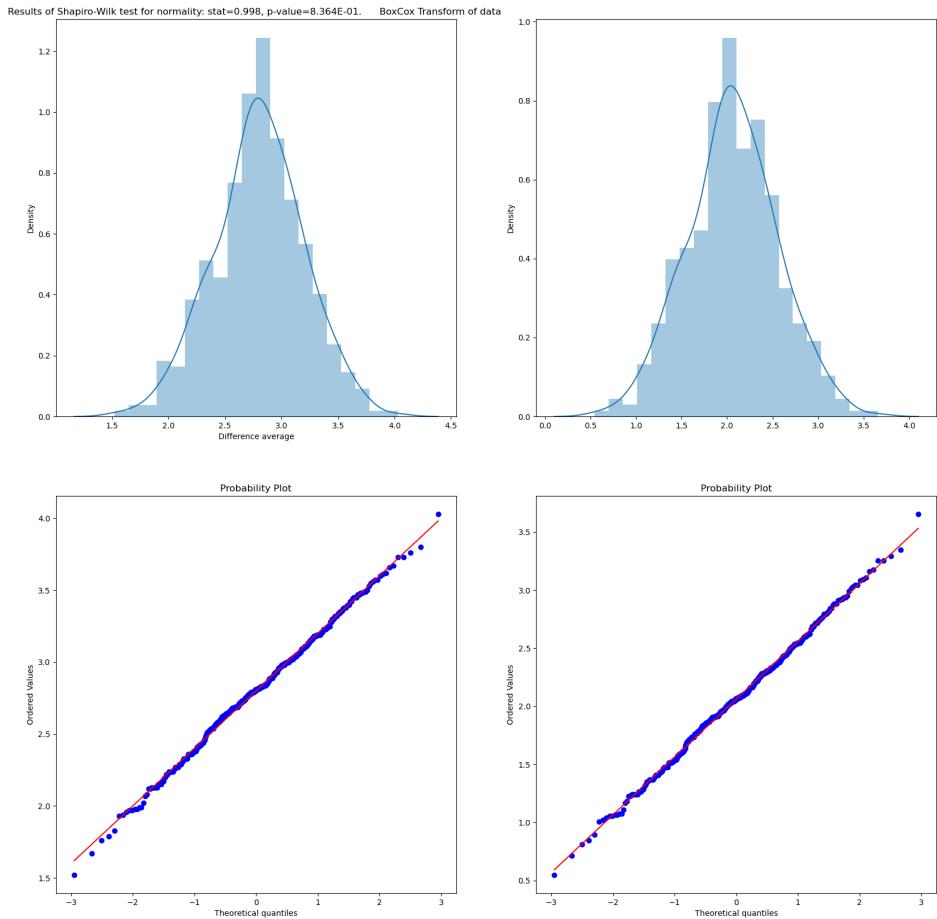


Figure 2.17: Example of boxcox applied to the radiomic feature *Difference Average* which has a close to normal distribution. The graphs contain the original distribution (top-left) the transformed distribution (top-right) and the two respective quantile-quantile plots

that, to obtain reasonable values in the cross-validation procedures and to avoid leakage³⁵ problems, all of the preprocessing steps have been done after train-test splitting the data on the train set and then applied as defined during training on the test dataset. When it comes to cross-validation procedure the choice was made to use a stratified k-fold approach with $k=10$, the data is split in 10 parts with the same percentages of labels³⁶ then a model is built by training on 9 of the folds and its performance is then tested on the remaining fold. To use the whole dataset for testing a prediction of it has been built by combining the predictions on the 10th fold for ten different models trained on the respective 9 remaining folds. The lasso model from training on the 9 folds is actually chosen as the model with the hyperparameters that give the best performance with another 10-fold crossvalidation. This has been obtained by using *cross_val_predict* on a model obtained by including a *LassoCV* step inside a *pipeline* from scikit-learn library in python. The performance of the cross-validated predictions that, when built this way, is much more representative of the real-world performance of the model, has been evaluated using ROC curves and AUC. Different models have been compared with a Delong test[7] for the significance of difference in the ROC curves.

The second analysis method used was RandomForest. As said before in this case no preprocessing was needed nor has been done, however particular care was taken in handling the imbalances in the dataset by using SMOTE [6] once again being careful to avoid leakage. The performance of this model was evaluated using confusion matrices which, at a glance, provide very much information on the situation of the data.

Finally a few dimensionality reduction techniques, namely the unsupervised PCA[8] and Umap [19] and the supervised PLS-DA[4], have been used to understand better the state of the data and further explain some of the obtained results.

This concludes the discussion on the methodologies used and leads

³⁵This term is used in the field of Machine Learning. It refers to models being created on information that comes from outside the training data.

³⁶The stratified in the name refers to this property. This method is useful when dealing with unbalanced datasets, such the one under analysis, in which the label has an uneven 15-85% frequency of occurrences of the two labels

¹²³⁴ perfectly in the discussion of the results.

¹²³⁵ Chapter 3

¹²³⁶ Results

¹²³⁷ In this chapter the result obtained with the methods explained in the
¹²³⁸ previous chapters will be briefly presented to ease in the discussion in
¹²³⁹ the final chapter.

¹²⁴⁰ 3.1 LASSO

¹²⁴¹ When it comes to lasso regression usually graphs are reported that
¹²⁴² show the convergence of the parameters to the final value. Since the
¹²⁴³ real information of this process is the value to which the coefficients
¹²⁴⁴ converge only these values will be presented in tables and the ROC
¹²⁴⁵ curves, with respective AUCs, will be provided. An example of the
¹²⁴⁶ graph i'm talking about is the one visible in Figure 3.1.

¹²⁴⁷ All of the graphs with respective captions will be, maybe, put in
¹²⁴⁸ the final appendix. Finally it seems useful to report the following
¹²⁴⁹ contingency table that gives an idea on how superimposed the labels
¹²⁵⁰ on death and ICU admission are.

		ICU Admission	
		0	1
Death	0	311	47
	1	48	30

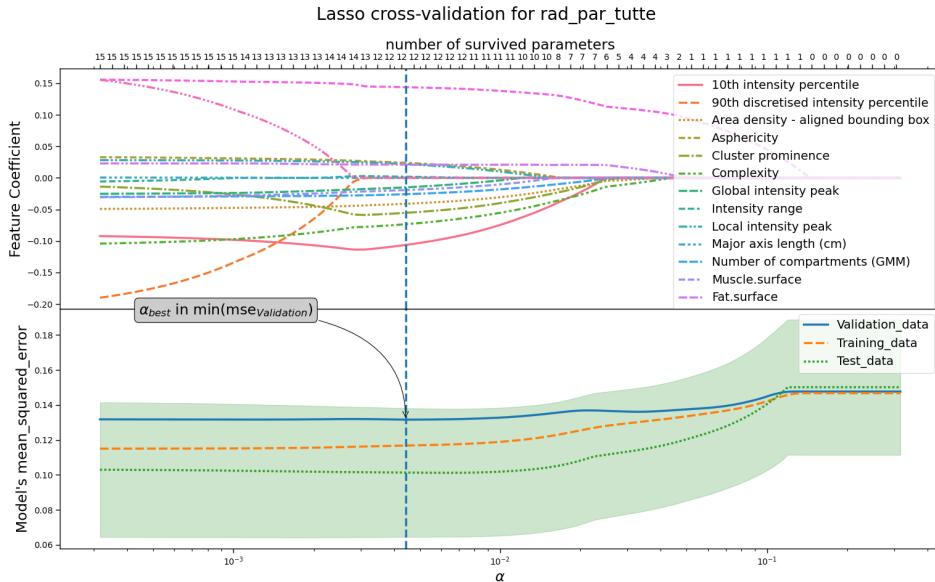


Figure 3.1: Example of graph representing the convergence of the coefficients in a lasso procedure. The final model is chosen by looking at the lowest value in mean squared error computed on the testing dataset. This graph is relative to only radiomic features with the model using death as label.

1252 3.1.1 Death

1253 Starting to predict the death outcome of the patient different groups
 1254 of features have been used, first of all only the radiomic features have
 1255 been used. The ROC curve obtained with the radiomic features is the
 1256 one in Figure 3.2. The curve in bold is an average curve obtained by
 1257 aggregating the ten curves relative to each of the folds used in testing
 1258 and the gray band represents a ± 1 standard deviation.

1259 The model was built using the coefficients reported in 3.1 reaches
 1260 a $AUC = 0.76 \pm 0.09$. The features inside the tabular, as well as those
 1261 in the tabulars that follow, will have the coefficients in descending
 1262 order by absolute value so that the top features are the most relevant
 1263 within it's relative model.

1264 Before commenting more in depth the performance of this model it
 1265 seems appropriate to see at least the other models built with singular
 1266 features groups, so, when it comes to the clinical features, the results
 1267 reported in Figure 3.3 and Table 3.7 have been obtained. All the results
 1268 will be put together for ease in Table 3.5.

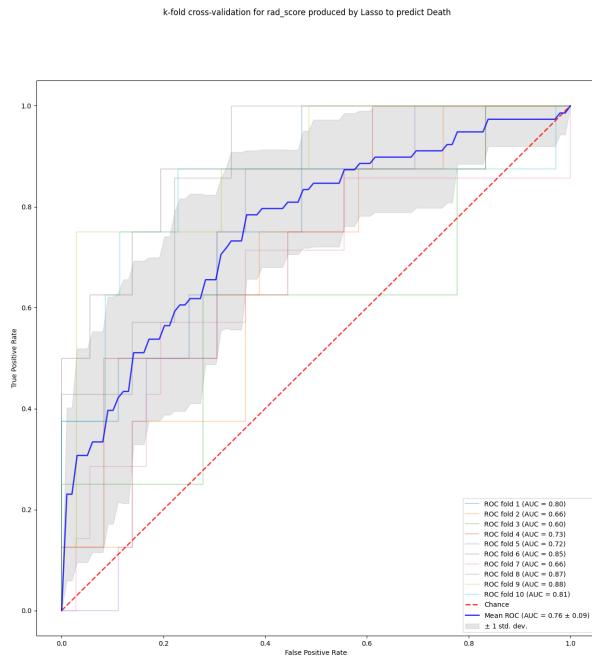


Figure 3.2: ROC curves obtained with crossvalidation procedure using the radiomic features alone.

Table 3.1: Coefficients used in the linear combination estimated by a Lasso regularization relative to the radiomic features

	lassocvL1 penalty	Importance
Intercept		0.178899
10th intensity percentile		-0.125094
Intensity-based interquartile range		0.103349
Complexity		-0.102924
Cluster prominence		-0.064690
Area density - aligned bounding box		-0.039374
Entropy		0.033002
Number of compartments (GMM)		-0.032441
Asphericity		0.028517
Local intensity peak		0.028478
Global intensity peak		-0.024832
Intensity range		0.012509
Fat.surface		0.007267
Major axis length (cm)		0.000000
Number of voxels of positive value		0.000000

Table 3.2: Coefficients used in the linear combination estimated by a Lasso regularization relative to the clinical features

	lassocvL1 penalty Importance
Intercept	0.178899
Age (years)	0.116771
Respiratory Rate	0.082292
Sex	-0.037591
Febbre	-0.022923
Hypertension	-0.000000
History of smoking	-0.000000
Obesity	0.000000

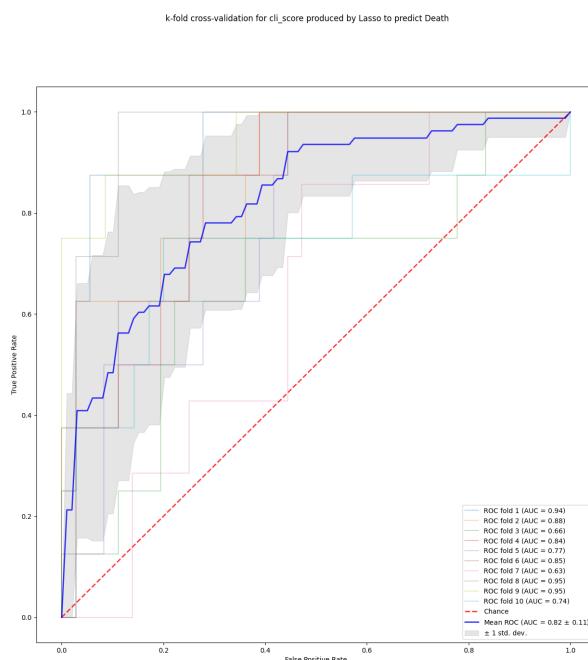


Figure 3.3: ROC curves obtained with crossvalidation procedure using the clinical features alone

Table 3.3: Coefficients used in the linear combination estimated by a Lasso regularization relative to the radiological features

	lassocvL1 penalty Importance
Intercept	0.178899
Ground-glass	-0.043875
Lung consolidation	0.038143
XRayTubeCurrent	-0.017264
KVP	0.004995
Crazy Paving	-0.000000
Bilateral Involvement	0.000000
SliceThickness	0.000000

1269 And finally, considering the radiological features Figure 3.4 and Ta-
 1270 ble 3.8 are obtained.

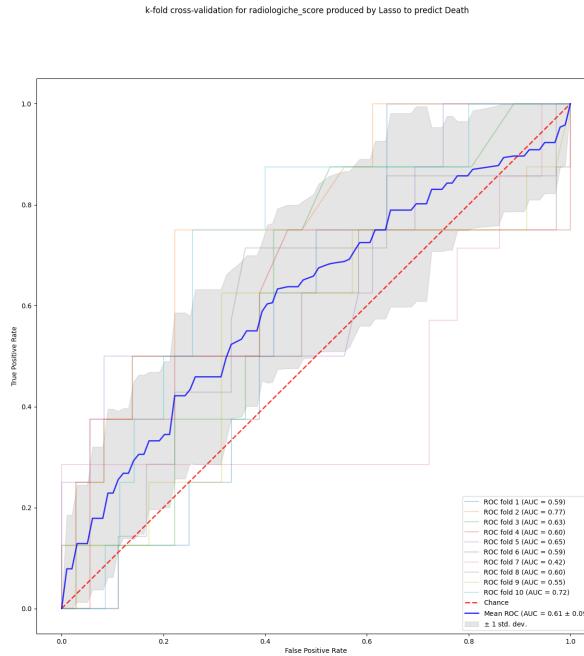


Figure 3.4: ROC curves obtained with crossvalidation procedure using the radiological features alone

1271 The first thing to notice is that the radiological features, when con-
 1272 sidered alone, have close to null predictive power. This is reasonable
 1273 for at least part of the features because there is no reason for ac-
 quisition parameter to actually influence the outcome of the patient.
 1274 When it comes to the radiologically determined quantities, such as
 1275 GGO, Crazy paving, lung consolidation and bilaterlaity even if one

would expect these to be relevant their distribution across the dataset is not condusive the good predictions. In fact 88% of patients had GGO, 50% of all patient had Lung consolidation, 77% of all patients did not have Crazy paving and 92% had bilateral involvment. When it comes to clinical features, categorys that performs better when considered singularly, nothing ground breaking has been obtained. Age, Respiratory rate and sex are the most relevant features and are all very much in concordance with what is expected. Finally radiomic features preform slightly worse than the clinical features. To see if the feature obtained are, at least, reasonable a quick explanation is needed. This will be done only for the top performing features:

- 10th intensity percentile and Intensity based interquartile range are both intensity based statistics. The first indic
- Complexity: A complex image is one that presents many rapid changes in intensity and is heavily non-uniform because it has a lot of primitive components
- Cluster Prominence: GLCM feature which measures the symmetry and skewness of the matrix from which it derives. When this is high the image is not symmetric
- Area density aligned bounding box: This is a ratio of volume to surface
- Entropy: Measures the average quantity of information needed to describe the image. In other words it quantifies randomness in the image, the more random the more info is needed to describe it.

To summarize, since all of the features are computed on the whole lung segmentation, it seems that the some information on the distribution of gray levels as well as some textural information inside the whole lung are important. It also seems that some information on the shape of the organ itself is also relevant. One would expect that when combining all of the available features, i.e. by building a model using the previous clinical, radiomic and radiological features, the performance should somewhat rise especially given the fact that clinical and

Table 3.4: Coefficients used in the linear combination estimated by a Lasso regularization relative to all available features features

	lassocvL1 penalty Importance
Intercept	0.178899
Age (years)	0.092963
Intensity-based interquartile range	0.057260
Respiratory Rate	0.049603
Ground-glass	-0.031423
Sex_bin	-0.028895
Complexity	-0.028606
Lung consolidation	0.017272
Febbre	-0.016933
XRayTubeCurrent	-0.016908
Area density - aligned bounding box	-0.009676
Cluster prominence	-0.006663
Fat.surface	0.004984
Number of compartments (GMM)	-0.001448
Local intensity peak	0.000195
Obesity	0.000000
Number of voxels of positive value	0.000000
Hypertension	0.000000
Intensity range	0.000000
Global intensity peak	-0.000000
Asphericity	0.000000
Crazy Paving	-0.000000
Bilateral Involvement	-0.000000
SliceThickness	0.000000
KVP	0.000000
10th intensity percentile	-0.000000
Entropy	0.000000
History of smoking	-0.000000

₁₃₁₀ radiomic features have almost the same performance. The combined
₁₃₁₁ results can be seen in Figure 3.5 and Table 3.9.

₁₃₁₂ In spite of what the expectations were the performance not only
₁₃₁₃ doesn't improve but it doesn't even change. To be sure of this claim a
₁₃₁₄ Delong test was used to compare pairwise the receiver operator curves
₁₃₁₅ and their respective AUCs. The null hypothesis of this test is that
₁₃₁₆ the two models are the same, hence a p-value smaller than 0.05 means
₁₃₁₇ that the curves and their aucs are statistically different. The results
₁₃₁₈ from this analysis can be seen in 3.6

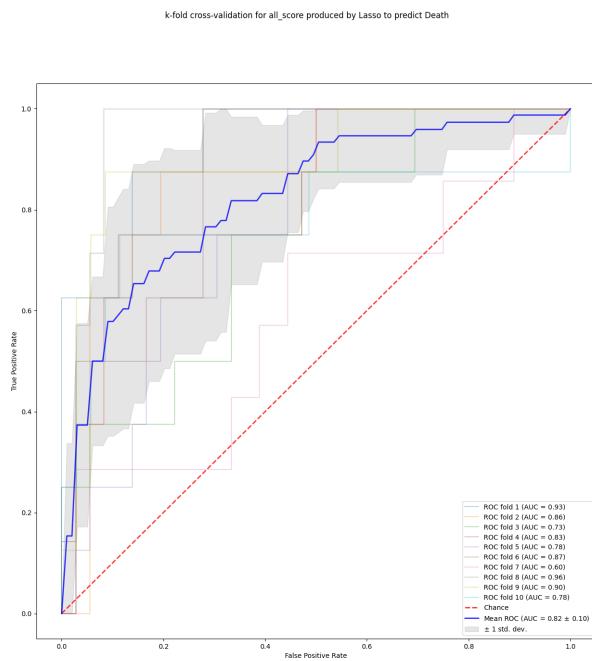
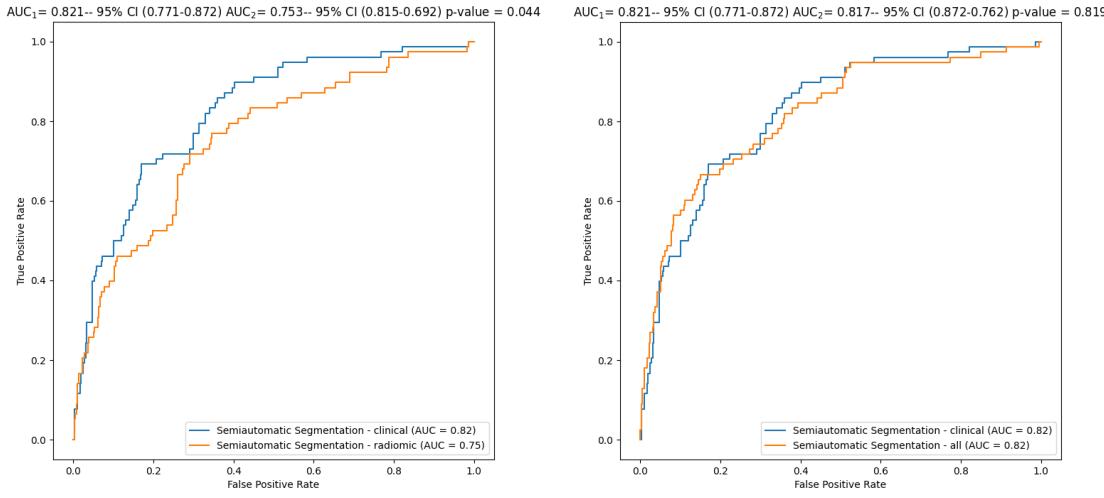


Figure 3.5: ROC curves obtained with crossvalidation procedure using all the available features

Table 3.5: Recap table with the performance of the various families of features

Features used	mean AUC ± std
Radiomic	0.76 ± 0.09
Clinical	0.82 ± 0.11
Radiological	0.61 ± 0.09
All	0.82 ± 0.10



(a) Comparison clinical-radiomic

(b) Comparison clinical-all

Figure 3.6: Comparison for the curves that are not evidently different

As one would have expected the models from radiomic and clinical features are different while the ones built using clinical and all features are not statistically different. Inspecting the coefficient of the parameters there are a few perplexing things to notice and a few reassuring ones. First of the reassuring facts is that the most relevant clinical features are still relevant in this combined model. Then, as one would expect, the radiological features retain some importance when combined with the others. However, when it comes to perplexing behaviours, the most concerning fact is that the radiomic features have mostly lost all relevance in the model which is surely unexpected. Some possible explanations will be given in the concluding remarks at the end of this subsection. as well as in the final chapter of the thesis.

3.1.2 ICU Admission

In trying to predict if the patient will be admitted in the Intensive Care Unit of the hospital the same procedure as before has been used. The ROC curves obtained with the various features are reported in Figure 3.7. Just like before the curve in bold is an average curve obtained by aggregating the ten curves relative to each of the folds used in testing and the gray band represents a ± 1 standard deviation. Following the

blueprint of the previous subsection, all of the results will be presented
 and then briefly discussed

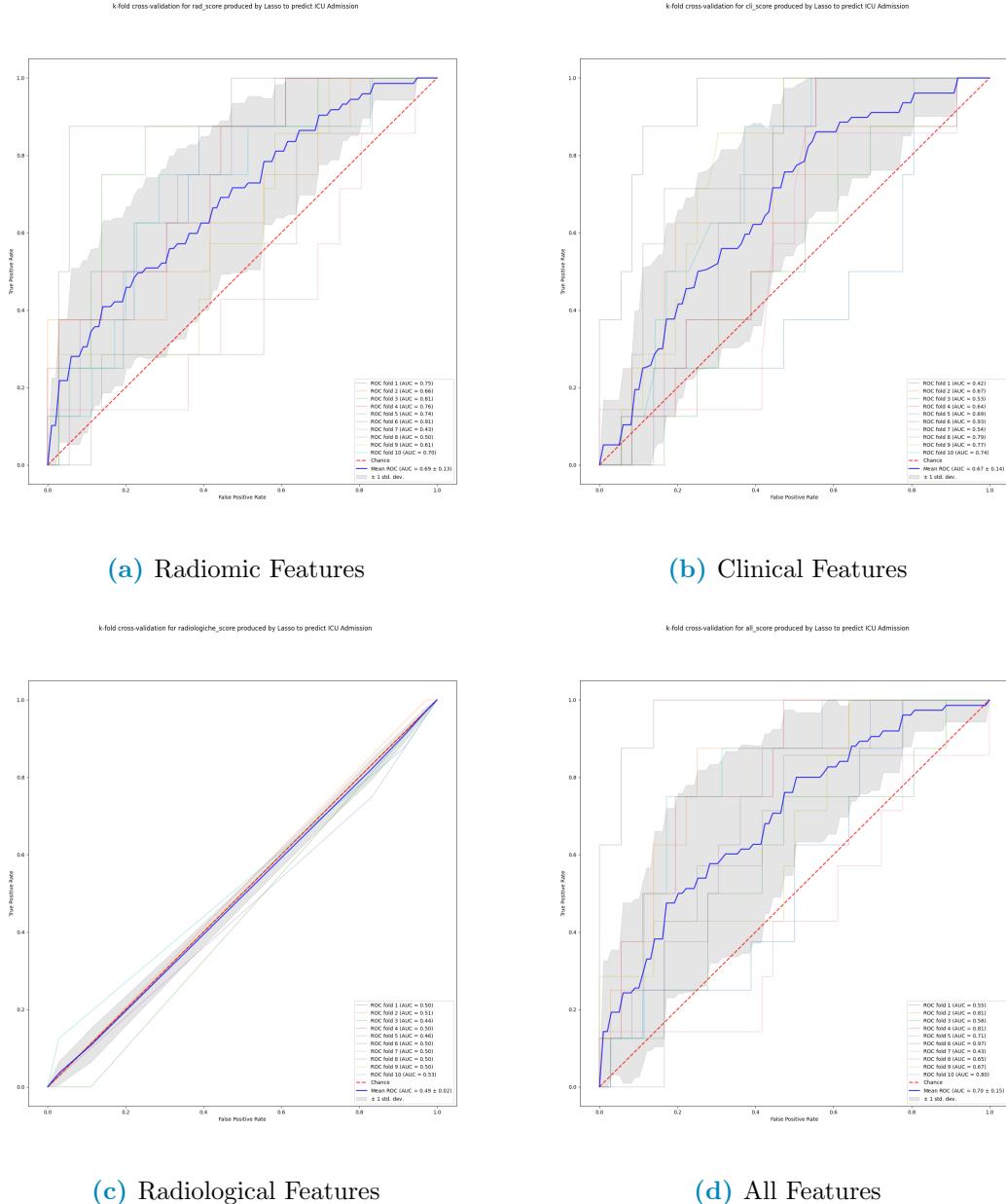


Figure 3.7: Performances of all the models

Compared to before the performance is definitely worse. The radiological features have the same performance of a random variable, which is not that concerning given their expected impact on gravity of the clinical picture of the patient. Even if superfluous a DeLong test was used to confirm that the hypothesis of the curves being equal could not be rejected. When it comes to the relevant features in each

Table 3.6: Coefficients used in the linear combination estimated by a Lasso regularization relative to the radiomic features

	lassocvL1 penalty Importance
Intercept	0.176605
Number of voxels of positive value	0.160751
Intensity range	-0.144834
Entropy	0.128999
Cluster prominence	-0.122290
Complexity	-0.093416
10th intensity percentile	-0.081133
Area density - aligned bounding box	-0.037373
Major axis length (cm)	-0.035723
Dependence count entropy	-0.029603
Fat.surface	0.027308
Asphericity	-0.023645
Local intensity peak	-0.019619
Global intensity peak	-0.016209
Number of compartments (GMM)	-0.000157

Table 3.7: Coefficients used in the linear combination estimated by a Lasso regularization relative to the clinical features

	lassocvL1 penalty Importance
Intercept	0.176606
Respiratory Rate	0.045510
Febbre	0.038332
History of smoking	0.036888
Hypertension	0.034547
Sex_bin	-0.031504
Obesity	0.030716
Age (years)	-0.014646

Table 3.8: Coefficients used in the linear combination estimated by a Lasso regularization relative to the radiological features

	lassocvL1 penalty Importance
Intercept	1.766055e-01
XRayTubeCurrent	6.111319e-18
Lung consolidation	0.000000e+00
Ground-glass	0.000000e+00
Crazy Paving	0.000000e+00
Bilateral Involvement	0.000000e+00
SliceThickness	-0.000000e+00
KVP	0.000000e+00

Table 3.9: Coefficients used in the linear combination estimated by a Lasso regularization relative to all available features features

	lassocvL1 penalty Importance
Intercept	0.176605
Number of voxels of positive value	0.109947
Dependence count entropy	0.070527
Cluster prominence	-0.069526
Intensity range	-0.057924
Febbre	0.044149
Hypertension	0.039127
SliceThickness	-0.037174
Complexity	-0.033393
History of smoking	0.032812
Age (years)	-0.032579
XRayTubeCurrent	-0.032182
Respiratory Rate	0.028504
Obesity	0.027580
Local intensity peak	-0.026135
Area density - aligned bounding box	-0.023027
Asphericity	-0.022999
Global intensity peak	-0.018280
Fat.surface	0.014179
Sex.bin	-0.004316
Crazy Paving	-0.003428
Ground-glass	0.003077
Lung consolidation	-0.001329
Bilateral Involvement	-0.001047
Number of compartments (GMM)	-0.000000
KVP	0.000000
10th intensity percentile	-0.000000

Table 3.10: Recap table with the performance of the various families of features

Features used	mean AUC ± std
Radiomic	0.69 ± 0.13
Clinical	0.67 ± 0.14
Radiological	0.49 ± 0.02
All	0.70 ± 0.15

₁₃₄₆ model the following can be deduced:

- ₁₃₄₇ • For the clinical features all of them have a role in the prediction.
- ₁₃₄₈ The only surprising fact, even if it keeps a certain degree of plausibility, is that age is the less relevant out of the available features
- ₁₃₄₉ when it comes to ICU admission.
- ₁₃₅₀
- ₁₃₅₁ • None of the radiological features have virtually any impact
- ₁₃₅₂ • The radiomic features still value intensity measurements and disorder in the image as primary origins of information. However it seems that shape of the lung now has more relevance in the whole
- ₁₃₅₃ model.
- ₁₃₅₄
- ₁₃₅₅

₁₃₅₆ 3.2 Random forest

₁₃₅₇ When it comes to random forests the approach followed was similar, ₁₃₅₈ in the sense that the various combinations of features were tried, yet ₁₃₅₉ different, because the preprocessing step consisted in simply applying ₁₃₆₀ a standard scaler to the data. The performance of the models has ₁₃₆₁ been looked at using both confusion matrices and ROC curves, the ₁₃₆₂ last of which has the only objective of comparing RF classifiers to ₁₃₆₃ the previously described Lasso regression. Once again, following the ₁₃₆₄ description scheme used in the section dedicated to Lasso, the results ₁₃₆₅ will be divided using the two labels ICU Admission and Death.

3.2.1 Death

For the sake of brevity all of the results will be reported and then discussed. Since RF classifiers use all of the available features it is very space consuming to report a table with all of the importances for the radiomic features as well as the model with all the features. These two will either be reported in the appendix or a link will be provided to my github page where they will be available.

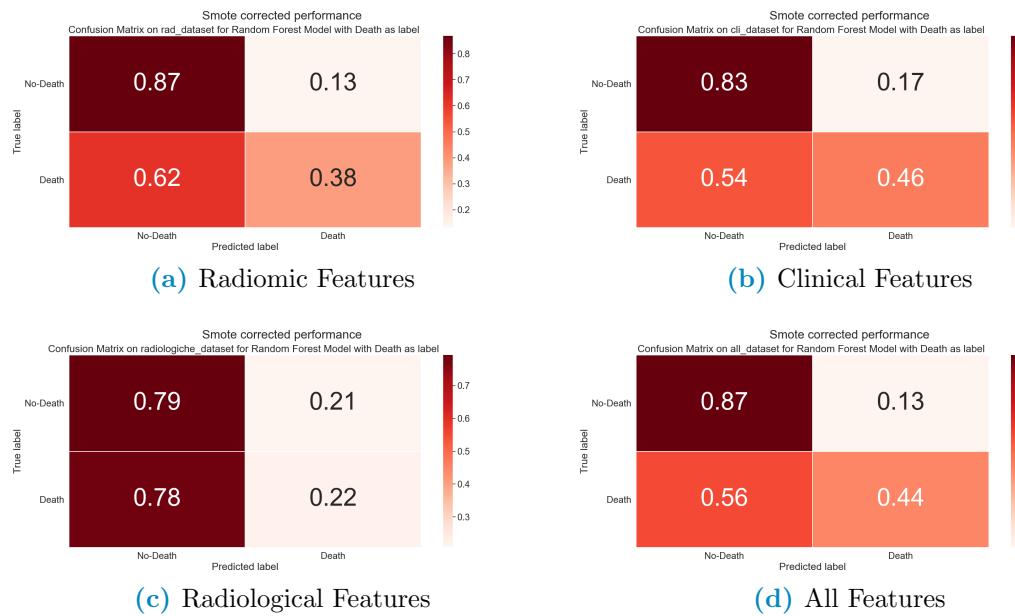


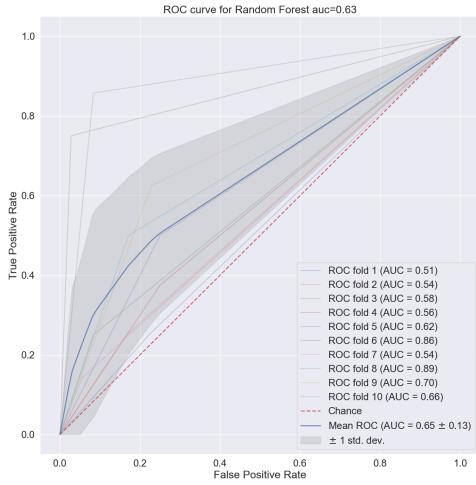
Figure 3.8: Performances of all the models

Even without looking at the ROC curves it's plain to see that the data at hand is proving to be difficult for this model. Much like before radiological features alone are useless. In evaluating these confusion matrices it should be kept in mind that the data is heavily unbalanced, since only $\sim 15\%$ of the patient died or were admitted in the ICU. Even when using SMOTE in the training phase to correct this problem it seems that the classifiers learn that it's optimal to guess that someone is alive. When it comes to the ROC curves Figure 3.9 and Table 3.11 summarises the results

Table 3.11: Recap table with the performance of the various families of features

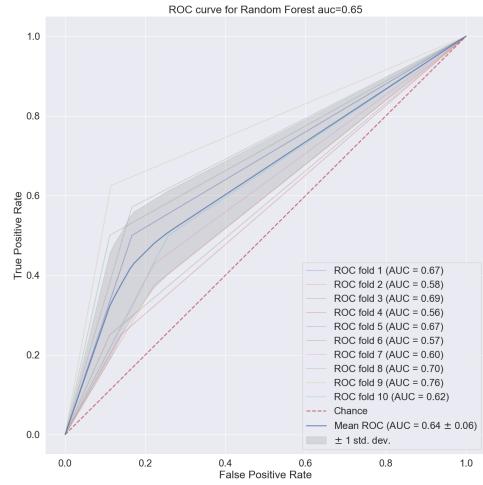
Features used	mean AUC \pm std
Radiomic	0.65 ± 0.13
Clinical	0.64 ± 0.06
Radiological	0.50 ± 0.07
All	0.66 ± 0.8

Smote corrected performance with rad features, predicting on Death



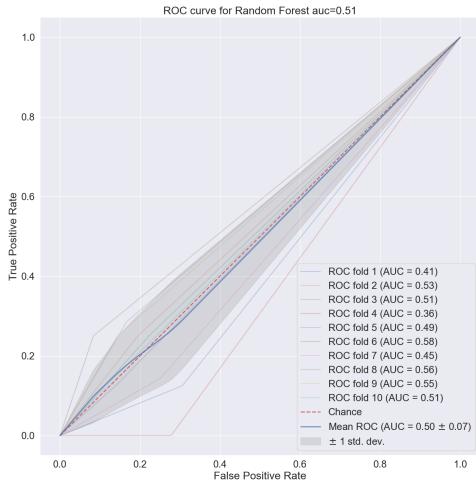
(a) Radiomic Features

Smote corrected performance with cli features, predicting on Death



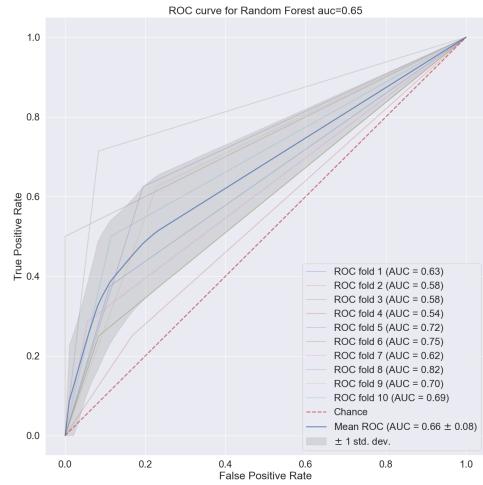
(b) Clinical Features

Smote corrected performance with radiologiche features, predicting on Death



(c) Radiological Features

Smote corrected performance with all features, predicting on Death



(d) All Features

Figure 3.9: Performances of all the models

Once again the curves are evidently not statistically different. This

1383 counterintuitive behaviour seems to be constant across the two im-
 1384 plemented methods and also across different labels tried. An attempt
 1385 to explain this phenomenon will be postponed to the end of the next
 1386 subsection

1387 3.2.2 ICU Admission

1388 Even when using the admission in the ICU the performance remains
 1389 pretty much the same, so the comments would still be the same as
 1390 before

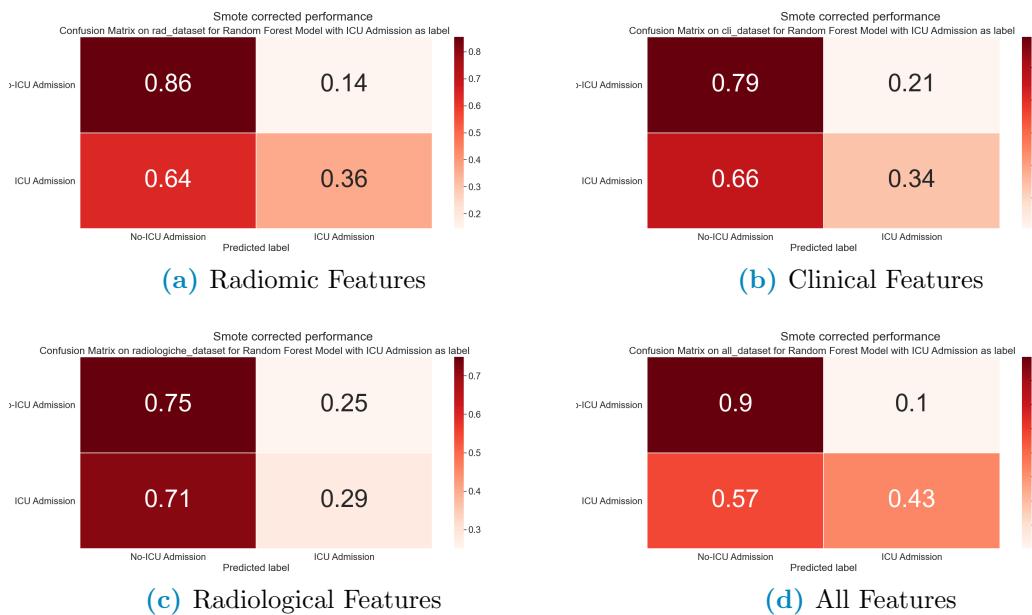


Figure 3.10: Performances of all the models

Table 3.12: Recap table with the performance of the various families of features

Features used	mean AUC \pm std
Radiomic	0.62 ± 0.08
Clinical	0.56 ± 0.08
Radiological	0.51 ± 0.11
All	0.64 ± 0.09

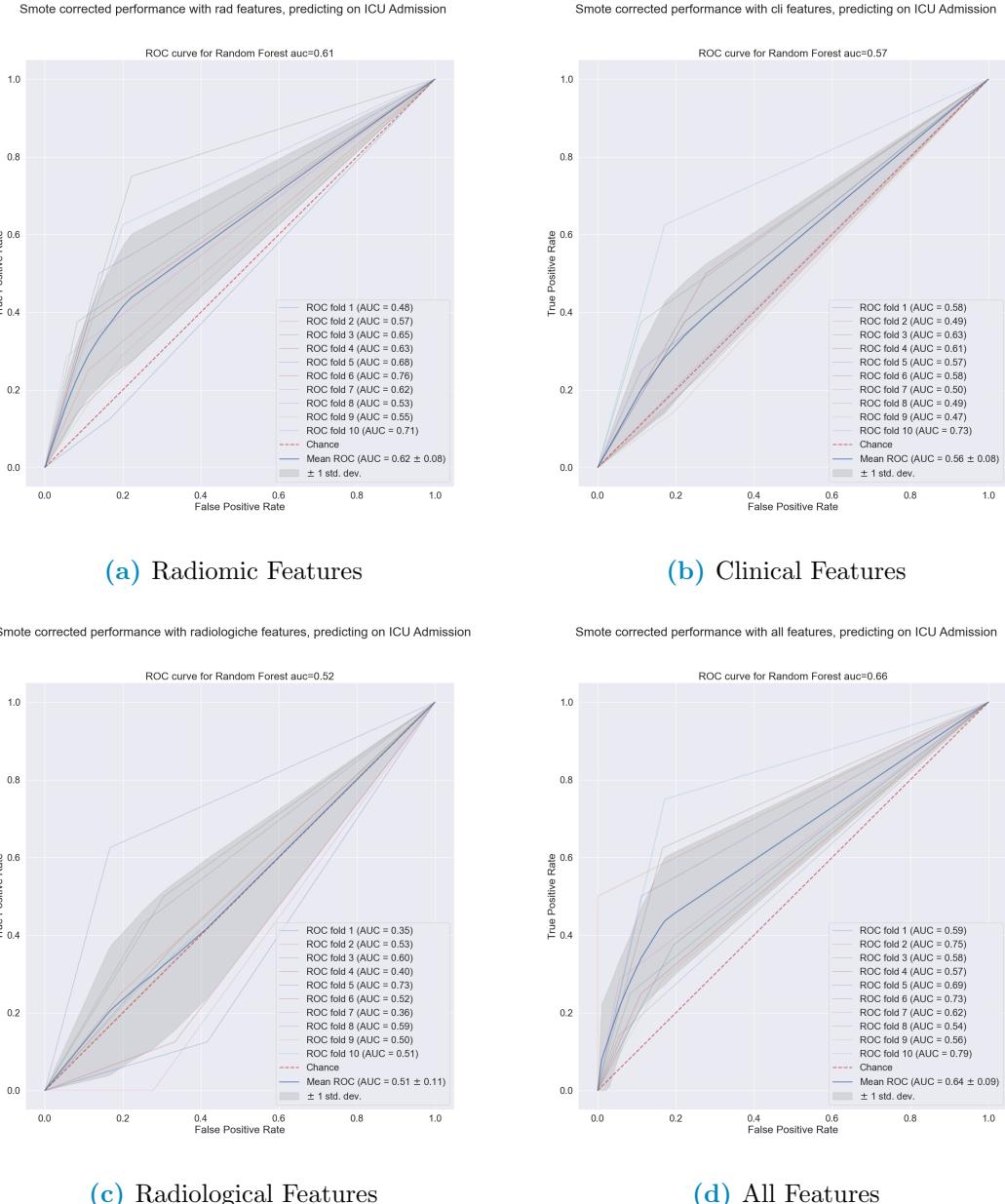


Figure 3.11: Performances of all the models

1391

It is very difficult to quantify what the expectations for each model

were a priori, this transposes to difficulties in trying to diagnose problems in the models when considered singularly. However when combining all of the available features it's very strange that no improvement in performance can be achieved. This would mean that 8 clinical variables provide the exact same amount of information as the total ~200 features most of which derived from images which, at least ideally, contain much more accurate information. One could surmise that the analysis methods have been implemented in an incorrect way, yet when two methods implemented differently and separately obtain the same result it reinforces the idea that the problem is somewhere before the analysis. All of these analyses started from the following hypotheses:

1. Clinical labels are informative of the final prognosis of the patient
2. Radiological images contain a lot of useful information
3. Radiomics can extract these information
4. This information is informative of the prognosis of the patient

It is very clear that the first three hypotheses are verified with a caveat, the performance of the radiomic pipeline hinges on the quality of the images and of the segmentation procedure performed on them. Since the images used in this thesis were, are and will be used by the hospital in routine processes it's close to impossible that all of them have problems, especially because of the preliminary screening done before segmentation. Another possibility is that the images are not really representative of the situation of the patient. Since one of the prevailing properties of *Sars-COVID19* is the speed with which the clinical picture of the patient can change it's possible that images taken at admission are not as informative of the final prognosis. This problem is, however, unavoidable in the setting of this thesis which has as aim the construction of a model that, exactly at admission, can discriminate between serious and easier cases. Another possibility is that there are two or more subgroups in the patient cohort and the performance on these is widely different, determining an average performance below the expectations. To diagnose if this is the case a few dimensionality reduction techniques have been used to visualize the data and to prepare for clustering in case of need, all of the results of these procedures will be presented in the following section.

1427 3.3 Dimensionality reduction

1428 For these analyses the data was always fed into a standard scaler
1429 before applying the technique of choice, furthermore a custom gravity
1430 score by classifying as 4 the dead individuals and then by assigning
1431 a progressive score from 1 to 3 by looking at the time of permanence
1432 was built as follows:

- 1433 1. Gravity 1: Survived individuals with permanence from 0th per-
1434 centile to 25th percentile
- 1435 2. Gravity 2: Survived individuals with permanence from 25th per-
1436 centile to 75th percentile
- 1437 3. Gravity 3: Survived individuals with permanence from 75th per-
1438 centile to 100th percentile
- 1439 4. Gravity 4: Dead individuals without regard for permanence in
1440 the hospital

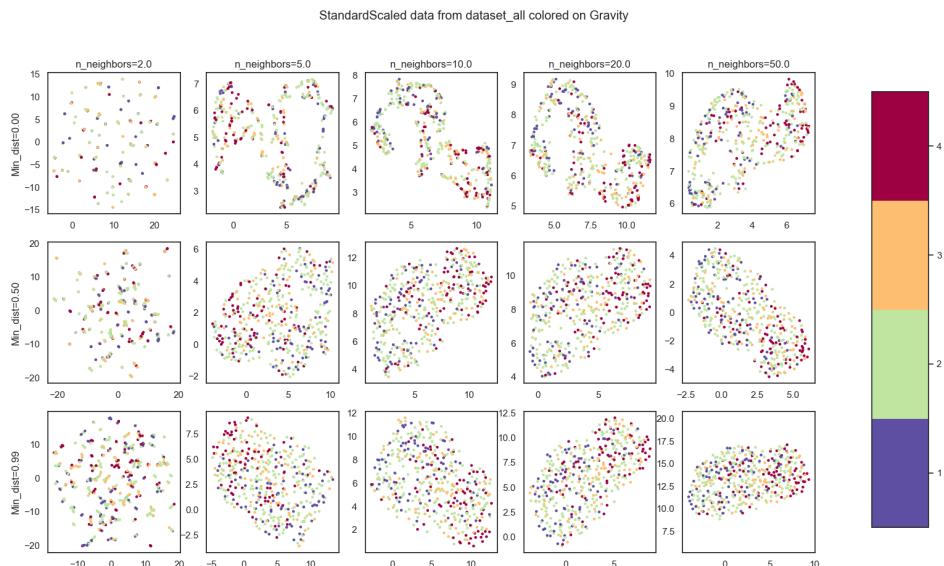


Figure 3.12: Possible combination for umap hyperparameters "number of neighbours" and "minimum distance". Color coding is done with aforementioned gravity score and all the features, i.e. clinical radiomic and radiological, were used.

1441 Also, before proceeding, the hyperparameter space for umap was
1442 explored since it's the method that allows the most control over rather

1443 intuitive parameters. Changing the value of the minimum distance of
1444 points in the final space from 0 to 0.99 changes how the structure is
1445 projected, while changing the number of neighbours changes how much
1446 the local or global structure of the data influences the final projection.

1447 3.3.1 PCA

1448 Starting from PCA, the data was reduced to either two or three di-
1449 mensions considering clinical and radiomic features, both separated
1450 and together. In this first example there seems to be a kind of left
1451 leaning polarization of the dead individuals, however there are no clear
1452 separations in the data.

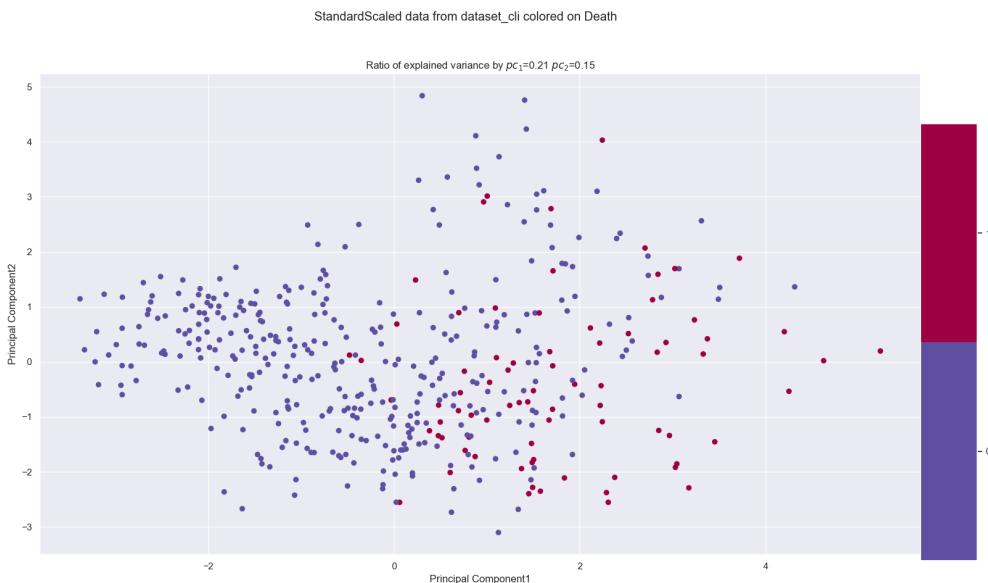


Figure 3.13: 2-Principal Component on clinical features.

1453 Working on the clinical dataset it can also be noted that the first
1454 two components of the PCA explain only 36% of the total variance.
1455 This leads to the conclusion that changes in the data cannot be ex-
1456 plained by a single, nor a few, features or linear combination thereof.
1457 The next approach was using the first three principal components
1458 using various labels available, most relevant of which being ICU ad-
1459 mission, Death and Gravity score.

1460 In all cases it seems like introducing the radiomic features causes
1461 the loss of the polarization structure that could be seen in the PCA

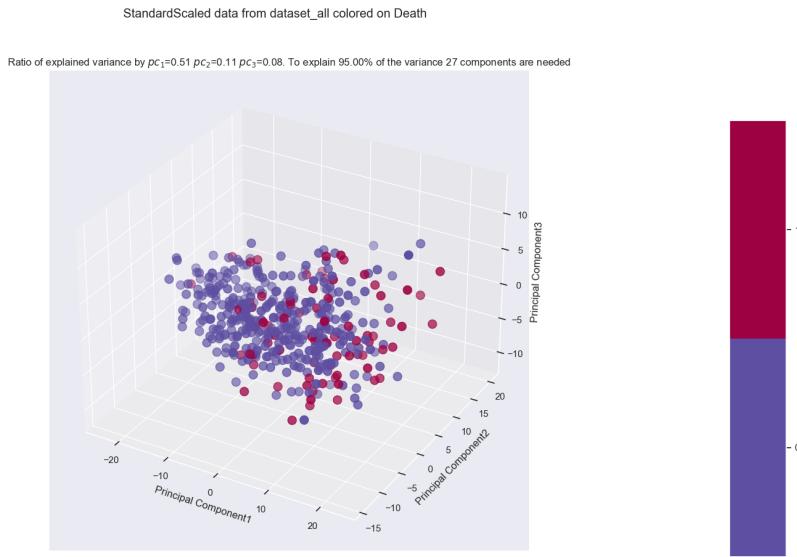


Figure 3.14: 3D PCA of whole dataset, colored with death label

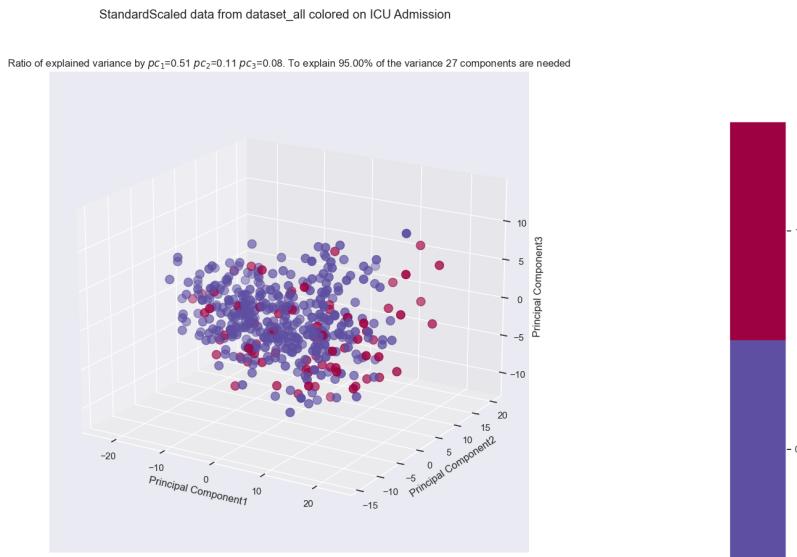


Figure 3.15: 3D PCA of whole dataset, colored with ICU Admission label

1462 on the clinical dataset alone in fig3.13.

1463 Since there are no visible clusters proceeding with cluster analysis
 1464 would mean incurring in the risk of finding non meaningful results so it
 1465 seemed appropriate to try other dimensionality reduction techniques.
 1466 The next technique tried was unsupervised Umap. Following the con-
 1467 clusions derived from fig:3.12 the number of neighbours was set to 10
 1468 and the minimum distance was set to 0. Once again the comparison

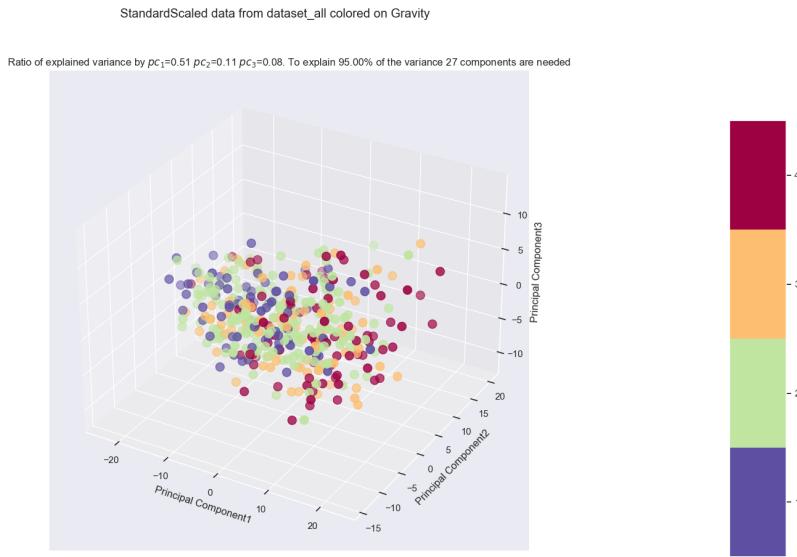


Figure 3.16: Comparison between various colour labels of the top 3 principal components for the entire dataset. Note that to explain 95% of the variance 27 components would be needed

were made between clinical and radiomic dataset as well as different possible labellings. Starting from the clinical dataset, without reporting all labels used, it's clear to see that the dataset seems to indicate very local well separated structures which don't seem correlated to gravity outcome

There are 9 well defined groups which don't seem to be correlated to any of the available labels. The dimension of these group is also very prohibitive if thinking of further analyses since groups of 35-50 people in a dataset with 15% mortality rate would mostly be very unbalanced if they were to be used for classification. However if the introduction of radiomic features were to unite some of these groups then this embedding could be meaningfully used for analysis. Looking at the 3D embedding for the whole dataset, the results are:

Once again the introduction of the radiomic feature seems to be a confounding factor in the seemingly clear-cut order present in the clinical dataset alone. There seems to be a well connected structure, which makes sense because umap sets out with the objective of preserving said structure. However since the variables of interest as label are Death, ICU admission or some kind of combination of them with

StandardScaled data from dataset_cli colored on Gravity
Umap parameters: N_Neighbours = 10, Min_distance =0, Distance metric=euclidean

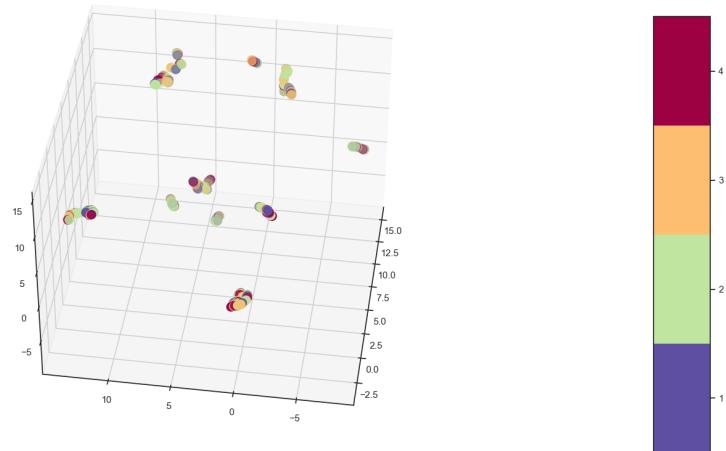


Figure 3.17: 3D umap of clinical dataset, colored based on gravity

StandardScaled data from dataset_all colored on Death
Umap parameters: N_Neighbours = 10, Min_distance =0, Distance metric=euclidean

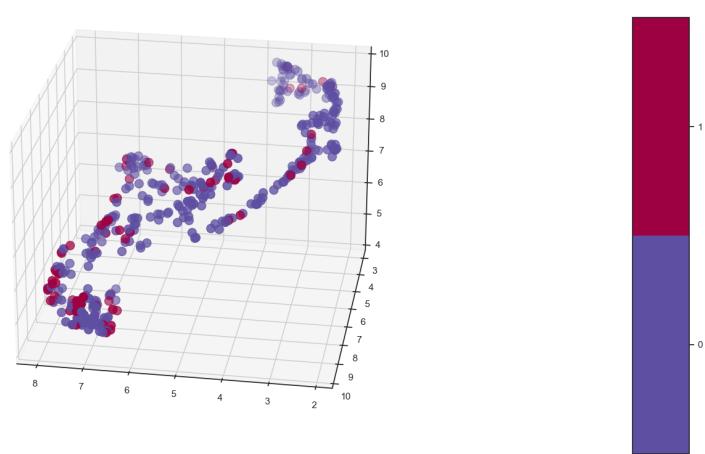


Figure 3.18: 3D umap of whole dataset, colored based on death

StandardScaled data from dataset_all colored on ICU Admission

Umap parameters: N_Neighbours = 10, Min_distance =0, Distance metric=euclidean

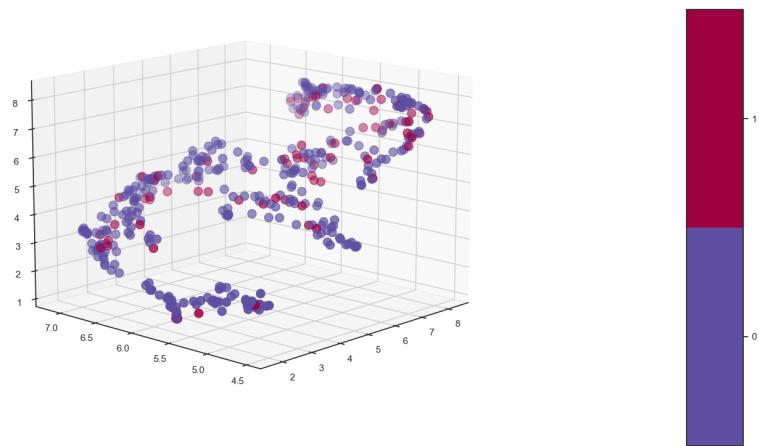


Figure 3.19: 3D umap of whole dataset, colored based on ICU admission

StandardScaled data from dataset_all colored on Gravity

Umap parameters: N_Neighbours = 10, Min_distance =0, Distance metric=euclidean

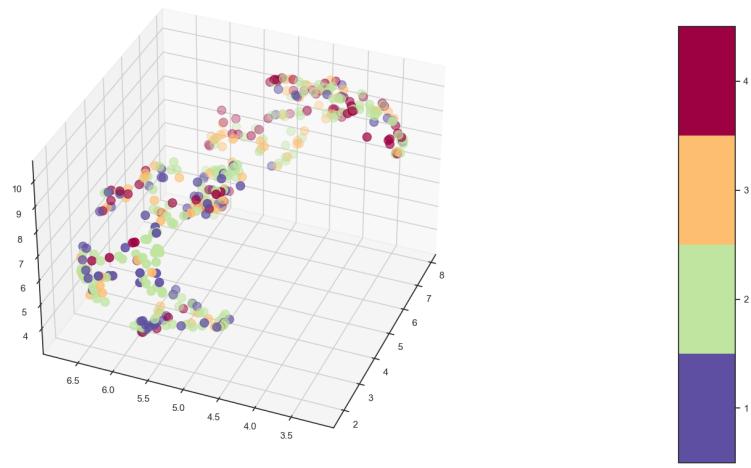


Figure 3.20: 3D umap of whole dataset, colored based on gravity

1488 hospital permanence there seems to be no visual correlation between
1489 structure and label. As such the next dimensionality reduction was
1490 tried to see if it yielded better results.

1491 Moving on from unsupervised methods to a supervised one, PLS-
1492 DA was used giving as label both death and ICU using both whole
1493 dataset, and singularly radiomic or clinical features. Starting from
1494 the clinical features alone, predicting on death Figure 3.21 can be
1495 obtained.

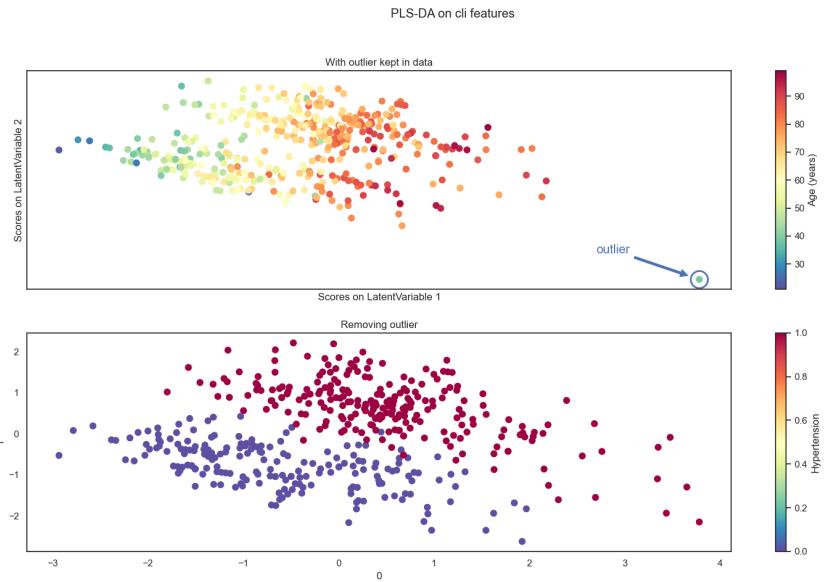


Figure 3.21: PLS-DA predicting on death coloured with age and hypertension

1496 In Figure 3.21 there are a few things to note. The first is the
1497 presence, in the top plot, of an outlier which, since PLS-DA is based
1498 on minimization of least squares, can ruin a lot the performance of the
1499 procedure. For this reason in the second plot the outlier was removed
1500 and the algorithm was run again on the cleaned data. The second
1501 thing to notice is the coloring used which, in the first plot, was used
1502 to highlight that along the one of the two latent variables the data is
1503 roughly distributed depending on age while, in the second plot, was
1504 used to highlight that the algorithm is able to perfectly separate the
1505 subjects with hypertension from those without it. However, by looking
1506 at the same embedding labelled with death and ICU admission fig 3.22
1507 can be obtained

1508 It's clear to see that, when predicting on death, the PLS-DA algo-

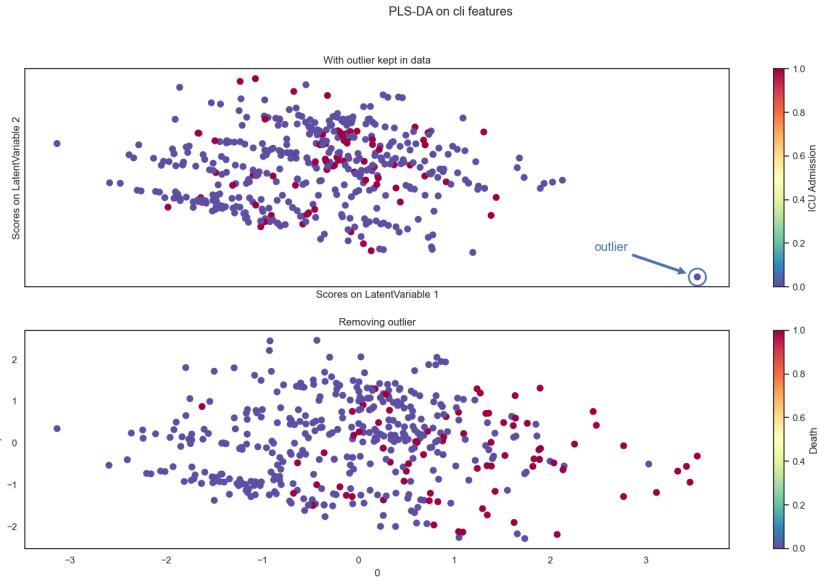


Figure 3.22: PLS-DA predicting on death coloured with death(bottom) and ICU Admission(top)

Table 3.13: PLS-DA feature weights in prediction on death using clinical features

Feature Name	Importance
Respiratory Rate	0.120206
Age (years)	0.116305
Obesity	0.004293
Hypertension	-0.004626
History of smoking	-0.012314
Febbre	-0.045431
Sex_bin	-0.054947

1509 rithm doesn't find any behaviour relevant for ICU admission. It's also
 1510 clear that there is at least a pattern of points labeled as dead being
 1511 towards the right of the image, this can be easily explained by looking
 1512 at how the ages are distributed in the first plot of fig: 3.21. From this
 1513 it's possible to deduce that older individuals tend to die more and that
 1514 hypertension does not seem to be relevant when considering death as
 1515 a clinical outcome. If necessary the PLS-DA algorithm allows also to
 1516 see the weights given to the features in predicting the label. At least
 1517 for the clinical dataset, which has a reasonable number of features, it's
 1518 interesting to report it ordering the coefficients by descending absolute
 1519 value:

1520 Doing the exact same procedure on the whole dataset, which means

1521 by including the radiomic features, Figure 3.23 can be obtained.

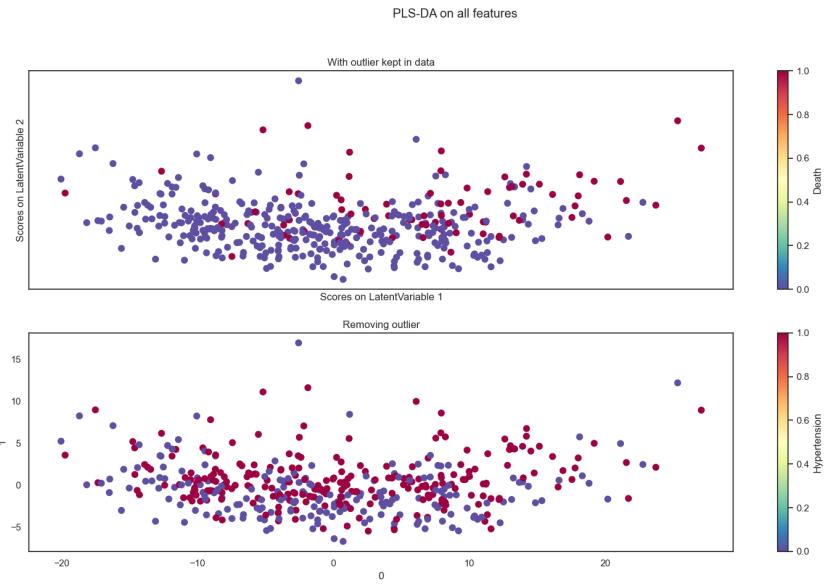


Figure 3.23: PLS-DA predicting on death coloured with death (top) and hypertension (bottom) on whole dataset

1522 Once again adding the radiomic features has evidently introduced
1523 noise in the system, which no longer displays any kind of behaviour,
1524 pattern nor separation. Doing the same analysis but using ICU Ad-
1525 mission as a label Figure 3.24 can be obtained.

1526 Now the colors have been chosen to highlight that respiratory rate and
1527 age have the main role in determining the latent variables. However,
1528 in this case, there doesn't appear to be a clear cut distinction as it
1529 happened before with hypertension. Looking at how the points scatter
1530 by coloring them according to the two interesting clinical labels the
1531 following figure can be obtained:

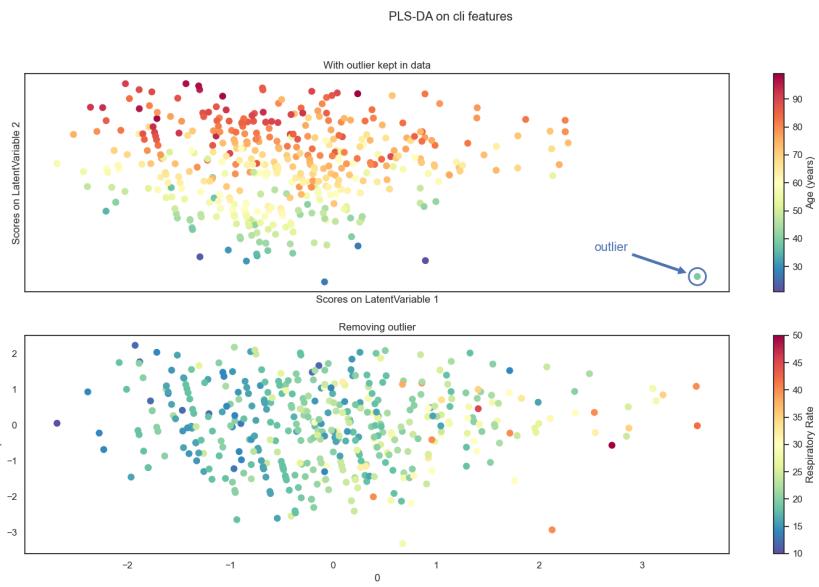


Figure 3.24: PLS-DA predicting on death coloured with Age and respiratory rate on clinical features

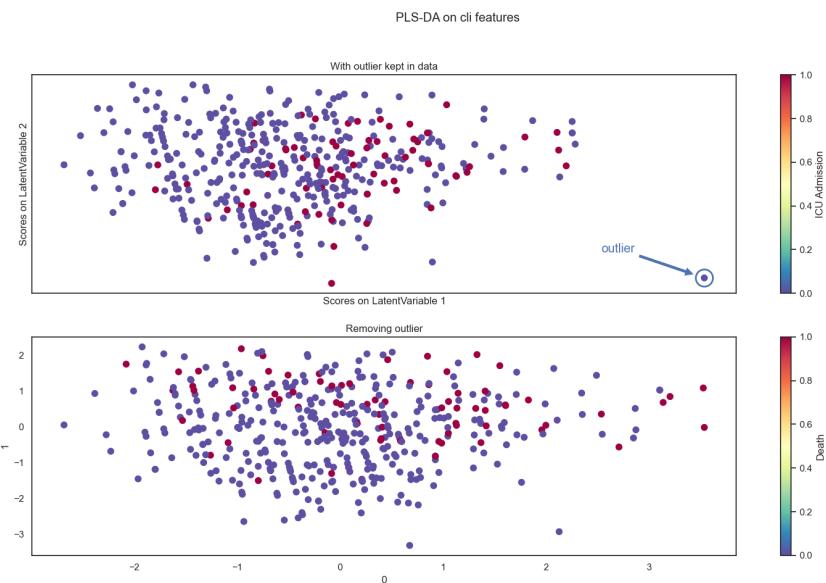


Figure 3.25: PLS-DA predicting on death(bottom) coloured with death and ICU admission(top) on clinical features

1532 Finally, introducing the radiomic features in the analysis the usual
 1533 effect of reducing separation can be seen can be seen in the figure
 1534 below:

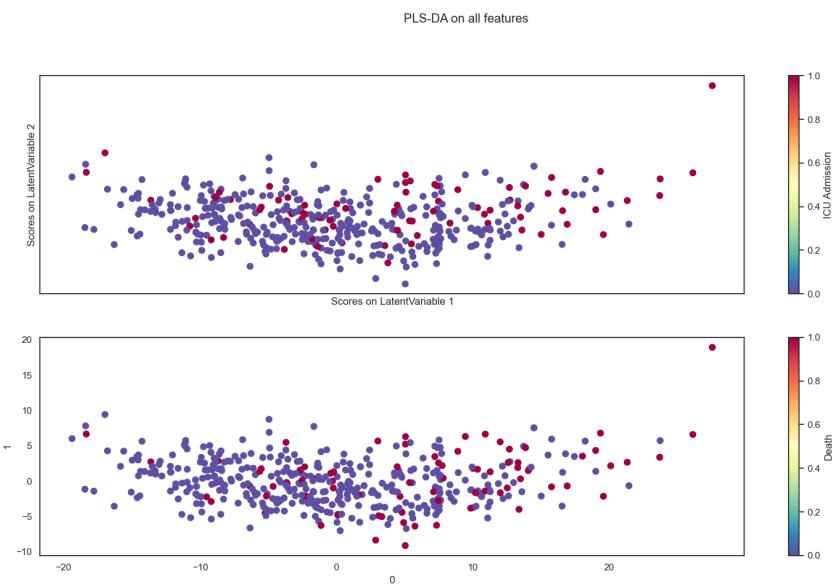


Figure 3.26: PLS-DA predicting on death coloured with death(bottom) and ICU admission(top) on all available features

¹⁵³⁵ Bibliography

- ¹⁵³⁶ [1] Nema ps3 / iso 12052, digital imaging and communications in
¹⁵³⁷ medicine (dicom) standard, national electrical manufacturers as-
¹⁵³⁸ sociation, rosslyn, va, usa. available free at, 2021.
- ¹⁵³⁹ [2] J. L. V. 1, R. Moreno, J. Takala, S. Willatts, A. D. Mendonça,
¹⁵⁴⁰ H. Bruining, C. K. Reinhart, P. M. Suter, and L. G. Thijs. The
¹⁵⁴¹ sofa (sepsis-related organ failure assessment) score to describe
¹⁵⁴² organ dysfunction/failure. on behalf of the working group on
¹⁵⁴³ sepsis-related problems of the european society of intensive care
¹⁵⁴⁴ medicine. *Intensive Care Medicine*, 1996.
- ¹⁵⁴⁵ [3] U. Attenberger and G. Langs. How does radiomics actually work?
¹⁵⁴⁶ – review. *RöFo - Fortschritte auf dem Gebiet der Röntgenstrahlen
und der bildgebenden Verfahren*, 193, 12 2020.
- ¹⁵⁴⁸ [4] M. Barker and W. Rayens. Partial least squares for discrimina-
¹⁵⁴⁹ tion, journal of chemometrics. *Journal of Chemometrics*, 17:166
¹⁵⁵⁰ – 173, 03 2003.
- ¹⁵⁵¹ [5] R. W. Brown, Y.-C. N. Cheng, E. M. Haacke, M. R. Thomp-
¹⁵⁵² son, and R. Venkatesan. *Magnetic Resonance Imaging: Physical
Principles and Sequence Design, 2nd Edition*. Wiley-Blackwell,
¹⁵⁵³ 2014.
- ¹⁵⁵⁵ [6] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer.
¹⁵⁵⁶ Smote: Synthetic minority over-sampling technique. *Journal of
Artificial Intelligence Research*, 16:321–357, Jun 2002.
- ¹⁵⁵⁸ [7] E. R. DeLong, D. M. DeLong, and D. L. Clarke-Pearson. Com-
¹⁵⁵⁹ paring the areas under two or more correlated receiver operating
¹⁵⁶⁰ characteristic curves: A nonparametric approach. *Biometrics*,
¹⁵⁶¹ 44(3):837–845, 1988.

- 1562 [8] K. P. F.R.S. Liii. on lines and planes of closest fit to systems of
1563 points in space. *The London, Edinburgh, and Dublin Philosophical*
1564 *Magazine and Journal of Science*, 2(11):559–572, 1901.
- 1565 [9] D. G. George H. Joblove. Color spaces for computer graphics.
1566 *ACM SIGGRAPH Computer Graphics*, (12(3), 20–25), 1978.
- 1567 [10] L. Guo. Clinical features predicting mortality risk in patients with
1568 viral pneumonia: The mulbsta score. *Frontiers in Microbiology*,
1569 2019.
- 1570 [11] G. N. Hounsfield. Computed medical imaging. nobel lecture.
1571 *Journal of Computer Assisted Tomography*, (4(5):665-74), 1980.
- 1572 [12] L. M. J. Cine computerized tomography. *The International Jour-*
1573 *nal of Cardiac Imaging*, 1987.
- 1574 [13] A. Kaka and M. Slaney. *Principles of Computerized Tomographic*
1575 *Imaging*. Society of Industrial and Applied Mathematics, 2001.
- 1576 [14] J. Kirby. Mosmeddata: dataset, 09/2021.
- 1577 [15] P. Lambin, R. T. H. Leijenaar, T. M. Deist, J. Peerlings, E. E. C.
1578 de Jong, J. van Timmeren, S. Sanduleanu, R. T. H. M. Larue,
1579 A. J. G. Even, A. Jochems, Y. van Wijk, H. Woodruff, J. van
1580 Soest, T. Lustberg, E. Roelofs, W. van Elmpt, A. Dekker, F. M.
1581 Mottaghy, J. E. Wildberger, and S. Walsh. Radiomics: the bridge
1582 between medical imaging and personalized medicine. *Nature re-*
1583 *views. Clinical oncology*, 14(12):749—762, December 2017.
- 1584 [16] L. Lee and C.-Y. Liong. Partial least squares-discriminant anal-
1585 ysis (pls-da) for classification of high-dimensional (hd) data: a
1586 review of contemporary practice strategies and knowledge gaps.
1587 *The Analyst*, 143, 06 2018.
- 1588 [17] G. Lemaître, F. Nogueira, and C. K. Aridas. Imbalanced-learn: A
1589 python toolbox to tackle the curse of imbalanced datasets in ma-
1590 chine learning. *Journal of Machine Learning Research*, 18(17):1–
1591 5, 2017.
- 1592 [18] J. McCarthy. What is artificial intelligence? 01 2004.
- 1593 [19] L. McInnes, J. Healy, and J. Melville. Umap: Uniform manifold
1594 approximation and projection for dimension reduction, 2020.

- 1595 [20] S. J. McMahon. The linear quadratic model: usage, interpretation
1596 and challenges. *Physics in medicine and biology*, 2018.
- 1597 [21] S. Morozov. Mosmeddata: Chest ct scans with covid-19 related
1598 findings dataset. *IAU Symp.*, (S227), 2020.
- 1599 [22] MosMed. Mosmeddata: dataset, 28/04/2020.
- 1600 [23] Y. Ohno. Cie fundamentals for color measurements. *International*
1601 *Conference on Digital Printing Technologies*, 01 2000.
- 1602 [24] D. L. Pham, C. Xu, and J. L. Prince. Current methods in medical
1603 image segmentation. *Annual Review of Biomedical Engineering*,
1604 2(1):315–337, 2000. PMID: 11701515.
- 1605 [25] P. C. Rajandeep Kaur. A review of image compression techniques.
1606 *International Journal of Computer Applications*, 2016.
- 1607 [26] W. C. Roentgen. On a new kind of rays. *Science*, 1986.
- 1608 [27] A. Saltelli. *Global Sensitivity Analysis. The Primer*. John Wiley
1609 and Sons, Ltd, 2007.
- 1610 [28] C. P. Subbe. Validation of a modified early warning score in
1611 medical admissions. *QJM: An international journal of medicine*,
1612 2001.
- 1613 [29] L. Tommy. Gray-level invariant haralick texture features. *PLOS*
1614 *ONE*, 2019.
- 1615 [30] S. Webb. The physics of medical imaging. 1988.
- 1616 [31] L. WS, van der Eerden MM, and e. a. Laing R. Defining commu-
1617 nity acquired pneumonia severity on presentation to hospital: an
1618 international derivation and validation study. *Thorax* 58(5):377-
1619 382, 2003.
- 1620 [32] A. Zwanenburg, M. Vallières, M. A. Abdalah, H. J. W. L. Aerts,
1621 V. Andrearczyk, A. Apte, S. Ashrafinia, S. Bakas, R. J. Beukinga,
1622 R. Boellaard, and et al. The image biomarker standardization init-
1623 iative: Standardized quantitative radiomics for high-throughput
1624 image-based phenotyping. *Radiology*, 295(2):328–338, May 2020.