bachelor's thesis Bachelorarbeit

Image enhancement for historical documents

David Spitzenberg

March 26, 2018

supervising examiner: Prof. Dr. Tomas Sauer

Chair of Digital Image Processing Faculty of Computer Science and Mathematics University of Passau

Betreuender Prüfer: Prof. Dr. Tomas Sauer

Lehrstuhl für Mathematik mit Schwerpunkt Digitale Bildverarbeitung Fakultät für Informatik und Mathematik Universität Passau co-supervisor: Prof. Dr. Malte Rehbein

Chair of Digital Humanities Faculty of Arts and Humanities University of Passau

Ko-Betreuer: Prof. Dr. Malte Rehbein

Lehrstuhl für Digital Humanities Philosophische Fakultät Universität Passau

This page is intentionally left blank.

Abstract

Historical manuscripts often feature passages illegible to the human viewer. Here, the illegibilities are not caused by unclear writing. Over the course of this thesis, we study illegibilities caused by covering writing with paper. In particular, we study the non-invasive reconstruction of ink-writing on paper covered by additional layers of paper. Based on our studies, we try to recover covered writing found in an early modern urbarium. We build our reconstruction on the non-invasive data acquisition tool of multispectral imaging. We analyze the obtained images using the framework of recursive principal component analysis, an adapted version of principal component analysis. Examining the interplay of multispectral imaging and covered writing, we show major limitations of the taken approach: the exposing radiation's wavelength, the covering layer of paper's thickness, the physical distance between writing and covering layer of paper as well as the writing's stroke width. Furthermore, we recover the title and the outline of the urbarium's examined text. The results obtained for the urbarium require further examination of the manuscript. Based on this work, the further examination can be planned.

Historische Manuskripten beinhalten häufig Passagen, welche für den menschlichen Betrachter nicht lesbar sind. Die Unleserlichkeiten liegen hierbei nicht in undeutlicher Schrift begründet. Im Rahmen dieser Arbeit untersuchen wir Unleserlichkeiten, welche durch das Verdecken von Schrift durch Papier bedingt werden. Insbesondere untersuchen wir die nicht-invasive Wiederherstellung von Schrift bestehend aus Tinte auf Papier, welche durch weitere Lagen Papier verdeckt wird. Aufbauend auf unseren Untersuchungen versuchen wir, die in einem frühneuzeitlichen Urbar enthaltene verdeckte Schrift wiederherzustellen. Die für die Wiederherstellung benötigten Daten nehmen wir nicht-invasiv durch Multispektralfotographie auf. Wir analysieren die aufgenommenen Bilder mittels rekursiver Primärkomponenten Analyse, einer Adaption der Primärkomponenten Analyse. Durch Untersuchung des Zusammenspiels aus Multispektralfotographie und verdecktem Text zeigen wir zentrale Grenzen dieses Ansatzes auf: die Wellenlänge der belichtenden Strahlung, die Dicke der verdeckenden Lage Papier, der physische Abstand von Schrift und verdeckender Lage Papier sowie die Strichdicke der Schrift. Darüber hinaus stellen wir den Titel und den Aufbau des am Urbar untersuchten Texts wieder her. Die für das Urbar erhaltenen Ergebnisse erfordern weitere Untersuchungen des Manuskripts. Diese weiteren Untersuchungen können aufbauen auf dieser Arbeit gestaltet werden.

This page is intentionally left blank.

Acknowledgments

This thesis could not have been realized without the support by a number of people. I thank the parochial archives Kößlarn and the "Staatliche Bibliothek Passau" for providing me with the testcase manuscripts examined throughout this thesis. Furthermore, I thank Mr. Mario Puhane and the town archives Schärding for providing me with the urbarium examined throughout this thesis. Regarding the urbarium, I further thank Mr. Georg Thuringer for his visual examination of the urbarium. Moreover, I thank Prof. Dr. Malte Rehbein and Prof. Dr. Tomas Sauer for providing me with both the hard- and the software used to implement the developed analysis as well as supportive feedback on the development of the analysis. By giving constructive feedback regarding my report, Prof. Dr. Tomas Sauer furthermore contributed to the eventual version of this thesis. Last but not least, I especially thank Dr. Ludger Drost for his support throughout the work on this thesis in all aforementioned domains.

This page is intentionally left blank.

Contents

A	bstra	\mathbf{ct}		3
A	cknov	wledgn	nents	5
Li	st of	Figure	es	8
Li	st of	Tables	5	10
In	\mathbf{trod}	uction		12
1	Fun	damen	ntals of digital image processing	14
	1.1	Monoo	chromatic images	14
		1.1.1	Image models	14
		1.1.2	Bit depth of digital images	15
		1.1.3	Gray level transforms	16
	1.2	Color	images	17
2	Fun	damen	tals of multispectral imaging	18
-	2.1	Funda	mentals of light	18
	2.2	Chara	cteristics of multispectral imaging	20
	2.3	Comp	arison of color images and multispectral series	21
	2.4	Image	acquisition setup	$\frac{-}{22}$
		2 4 1	Hardware	22
		242	Software	22
		243	Setun	23
		2.4.4	Acquisition of a multispectral series	23 23
3	Illes	rible m	anuscripts and PCA	26
	3.1	The ri	inning example: parchment manuscripts	26
	3.2	Pre-pr	ocessing of raw images	$\frac{-3}{28}$
	0	321	Histogram equalization	30
		3.2.2	Grav level range maximization	32
		323	Normalization	35
		324	Pre-processing pipeline	35
	33	Introd	uction of PCA	36
	0.0	331	Redundancy elimination	38
		332	Linear basis transform	40
		333	O and Singular Value Decomposition	40
		0.0.0		40

		3.3.4 Variances and restriction to relevant basis vectors	41
		B.3.5 Summary and algorithm overview	42
		B.3.6 Application of PCA to the toy example	42
	3.4	Application of PCA to the parchment images	43
	~		
4	Cov	red text and RPCA	50
	4.1	The running example: paper manuscripts	51
	4.2	visual examination of covered text	52
	4.3	Application of PCA	54
	4.4	ntroduction of recursive principal component analysis	55
		4.4.1 The basic idea of recursive principal component analysis	57
		1.4.2 Interpreting recursive PCA	59
		4.3 Improving recursive PCA	60
		4.4.4 Interpreting the result of RPCA and pseudo image construction	63
	4.5	Application of RPCA in practice	64
		1.5.1 The results of RPCA	64
		5.2 Examination of Taylor inspired pseudo image construction	66
	4.6	Post-processing	69
	47	Comparison to baseline results	73
	1.1		
5	Ana	vsis of the urbarium	76
	5.1	Description of the urbarium	76
	5.2	Analysis of the urbarium	. 0 77
	0.2	2.2.1 Analysis of Patch 5	$\frac{1}{77}$
		Analysis of the overall results	۲1 21
			51
6	Fut	re work	88
Ŭ	61	mage acquisition	88
	6.2	Jid_processing	88
	6.3	Post processing	80
	0.5	ost-processing	59
Co	nclu	on	an
UU	Mai	findings	an
	One	anding	01
	Ope		91
Re	fere		92
100			, ,
Α	RP	A mid-processing decisions	94
	A 1	stack 2	94
	Δ 2	Patch 1	96
	Λ 3		06 06
	л.5 Л 4	at 2	90 07
	A.4		91
	A.O		91
	A.6		98
	A.7	Patch 6	98
р	Dia	al appondix 10	00
D	Dig	ai appendix 10	10
\mathbf{C}	Eide	stattliche Erklärung 10	01
		-	

List of Figures

1.1	Black to white gradients for bit depths $d = 6$ through $d = 9$	16
$2.1 \\ 2.2 \\ 2.3$	Incident light versus transmitted light exposure	19 20 25
$3.1 \\ 3.2 \\ 3.3 \\ 3.4$	The recto and the verso of the parchment discussed throughout Chapter 3 The recto of the parchment shown in Figure 3.1a with a focus on legibility Multispectral series of images of the parchment's recto	27 28 29
3.5	contrast	31
3.6	applied to a real-world image (cf. Figure 3.3a)	33
97	applied to an artificial image	34
১.7 ২.২	Regults of application of PCA to the toy example presented in Figure 3.7	37
3.0 3.0	Pre-processed multispectral images of the parchment introduced in Figure 3.3	44
3.3 3.10	Principal components obtained by application of PCA to images shown in Figure 3.9	43 48
4.1	Document 1 used to set up paper stacks throughout Chapter 4	52
4.2	Document 2 used to set up paper stacks throughout Chapter 4	53
4.3	Document 3 used to set up paper stacks throughout Chapter 4	54
4.4	Document 1 with the parts considered throughout Section 4.2 highlighted in yellow	55
4.5	The pre-processed multispectral series obtained with Document 1	56
$4.6 \\ 4.7$	The view through frosted glass for increasing distances s between glass and sheet Documents 2 and 3 with the parts examined throughout Section 4.3 highlighted	57
	in yellow	57
4.8	Results of application of PCA to multispectral series depicting Stack 2	58
4.9	Charts outlining the data flow within a single application of PCA, recursive PCA and recursive PCA with mid-processing	60
4 10	Stack 2's means of order 0 through 10 obtained by application of BPCA	65
4.11	Pseudo image $F_{\mathbf{X}}^{\text{Taylor}}(x)$ constructed for various $x \in \mathbb{R}$ using the means depicted	00
4.12	In Figure 4.10	67
	in Figure 4.10	70
4.13	Histogram of the constructed pseudo image depicted in Figure 4.12i prior to gray	_
	level range maximization	72

$\begin{array}{c} 4.14 \\ 4.15 \end{array}$	Result of application of our first draft post-processing pipeline to Figure 4.12i Comparison of the baseline result and the pseudo image obtained by analyzing	72
	Stack 2's multispectral series	74
4.16	Reprint of Figure 4.15b with a focus on legibility	75
5.1	The closed urbarium imaged from in front and from behind	77
5.2	Sketch of the back cover's structure	78
5.3	Transmitted light images of the back cover	79
5.4	The urbarium's back cover opened	80
5.5	The decomposition of the back cover's recto into patches	81
5.6	Pre-processed multispectral series depicting Patch 5 of the back cover's recto	82
5.7	Means $\overline{\mathbf{w}}_k$ of orders $k = 0, \dots, 2$ obtained by application of RPCA to the multi-	
	spectral series depicting Patch 5 of the urbarium	85
5.8	$F_{\mathbf{x}}^{\text{Equal}}$ and the post processed $F_{\mathbf{x}}^{\text{Equal}}$ obtained by equally weighting pseudo image	
	construction using the means depicted in Figure 5.7	86
5.9	The final results obtained by application of the analysis described throughout this	
	thesis to the urbarium	87

List of Tables

2.1	Approximative correspondence between wavelength and perceived color according	
	to Smith	19
2.2	Primary colors and their corresponding radiations' wavelengths according to the	
	CIE	22
2.3	Filter types and respective approximative transmissivities used throughout the	
	thesis	22
2.4	Default exposure configurations used throughout a single multispectral series	25

Introduction

The archiving of historical manuscripts has superseded the mere storing of manuscripts for a long time. Nowadays, people in charge examine the kept items closely instead. One goal of such examination is the reconstruction of illegible passages. The reconstruction increases the contemporary knowledge not only about the manuscripts themselves but also about the ages and circumstances they originate from. Due to the manuscripts' ages and values, examinations are required to impede as little damage to the manuscripts as possible. In the ideal case, the examination is completely non-invasive.

This thesis analyses one specific type of illegibility: ink-writing on paper covered by additional layers of paper. Here, the covering layers of paper render the writing illegible to the human viewer. As instance of a real-world object featuring such illegibility, we examine a early modern urbarium. Urbaria were used to keep track of fiefs during early modern feudalism.

As non-invasive tool of data acquisition, we choose the tool of multispectral imaging. Multispectral imaging emerged from geo- and astrophysics. In the past, this technique has been adopted to the needs of scientists in the field of digital humanities. In the field of digital humanities, multispectral imaging has lead to astonishing discoveries. Hence, the technique has become a widespread tool in recent days.

We examine the obtained multispectral series using a computer-based analysis. Here, our analysis is mainly based on the framework of principal component analysis and its recursive counterpart.

We begin this thesis with a brief introduction to the bedrock of contemporary imaging in general: digital image processing. Building on this basis, we then study the fundamentals of multispectral imaging. To the end of computer-based examination of multispectral series, we further introduce the framework of principal component analysis. Given the tools of multispectral imaging and principal component analysis, we elaborate on how both tools are typically applied in the field of digital humanities. To this end, we study the reconstruction of aged inkwriting on parchment as a typical use-case. Based on the insights from this first use-case, we analyze the interplay of paper-covered writing and the tools of multispectral imaging and principal component analysis afterwards. Here, we examine two test-cases featuring paper-covered writing. The writing we try to recover in both test-cases is known a priori. The comparison of the a priori knowledge and the results of our recovery permits us to tune the computer-based analysis of multispectral series. To the end of tuning, we transmute principal component analysis into its recursive counterpart by recursive application. Concluding this thesis, we apply the recursive principal component analysis to above mentioned urbarium. This concluding application establishes our first attempt to recover a long lost passage from the urbarium. Throughout this thesis, we base our reconstructions on the solely use of the multispectral imaging hardware already present at the Chair of Digital Humanities at the University of Passau. In this regard, this thesis establishes a first feasibility study of what reconstruction quality is possible with the present hardware.

INTRODUCTION

Based on the findings of this thesis, a suitable further examination of the urbarium can be developed. First, we can decide whether to keep or to replace multispectral imaging as primary data acquisition tool. Second, different measures to improve the recovery can be ranked. Such measures include the development of a more sophisticated computer-based analysis and the acquisition of further imaging hardware.

The rest of this thesis is structured as follows: In Chapters 1 and 2, we give a brief introduction to the fundamentals of digital image processing and multispectral imaging respectively. Here, Chapter 2 also describes the experimental setup used throughout this thesis. We then introduce the framework of principal component analysis in Chapter 3. This introduction is accompanied by the study of the use-case of parchment examination. Chapter 4 is dedicated to the application of multispectral imaging to covered text. Based on our findings, we transmute principal component analysis into recursive principal component analysis. We conclude our analytical work in Chapter 5 with the application of recursive principal component analysis to the urbarium's back cover. Chapter 6 concludes our overall work with a discussion of future work resulting from this thesis.

Chapter 1

Fundamentals of digital image processing

While color images have mostly driven monochromatic images out of daily life, monochromatic images still play an important role in science. This thesis is based on the analysis of monochromatic images obtained using the multispectral approach. This section is dedicated to giving a basic introduction to the bedrock of digital image processing. In Section 1.1, we give some details on monochromatic images. In Section 1.2, we discuss the basics of color images.

1.1 Monochromatic images

Digital image processing describes the science of acquiring, modeling, manipulating and displaying images using digital technology. Throughout the field, there exist various models of images of which we introduce some on the following pages. Here, we focus mainly on monochromatic images and give only a brief outlook to color images afterwards.

Due to the scope of this thesis, this introductory section remains rather brief. To get insights on the details, the interested reader is referred to [6], by which the following is heavily inspired.

1.1.1 Image models

Definition 1.1 specifies our most general understanding of a monochromatic image.

Definition 1.1. Let $f: \mathbb{R}^2 \to \mathbb{R}$ denote a mapping with

$$f(x,y) \ge 0 \quad \forall (x,y) \in \mathbb{R}^2$$

and

$$\sup_{(x,y)\in\mathbb{R}^2}f\left(x,y\right)<\infty$$

Then f is called a monochromatic, continuous image. Here, \mathbb{R}^2 is called the spatial space, \mathbb{R} is referred to as gray level space and f(x, y) is called gray level at $(x, y) \in \mathbb{R}^2$.

Continuous images appear convenient in the mathematical perspective. They allow - after imposing needed additional assumptions - the application of the full mathematical tool-set such as differentiation and integration. However, the modeling of continuous images using a computer is not so easy. Restricting the domain of an image, we obtain an image model more suitable for computer-based processing: **Definition 1.2.** Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow [0, 1]$ with $m, n \in \mathbb{N}$ denote a mapping. Then f is called sampled image. We further call the transition of domains from \mathbb{R}^2 to $\{0, \ldots, m-1\} \times \{0, \ldots, n-1\}$ sampling. We can interpret f as matrix $\mathbf{F} \in \text{Mat}([0, 1], m, n)$ given by

$$\mathbf{F} := \begin{pmatrix} f(0,0) & \dots & f(0,n-1) \\ \vdots & \ddots & \vdots \\ f(m-1,0) & \dots & f(m-1,n-1) \end{pmatrix}$$

The matrix' elements are called pixels. Analogously, a matrix $\mathbf{H} \in Mat([0,1], m, n)$ can be interpreted as image $h: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow [0,1]$.

Sampling an image eases the efficient storing of the image. Nonetheless, the real range [0, 1] still imposes some challenges to the computer-based modeling of a sampled image. Usually, the modeling is based on approximative floating point arithmetic of sufficient numerical accuracy. To get rid of the real range, we restrict the range of a sampled image even further:

Definition 1.3. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ with $m, n, g \in \mathbb{N}$ denote a mapping. Then f is called a digital image. The transition of ranges from $\mathbb{I} \subset \mathbb{R}$ to $\{0, \ldots, g\}$ is called quantization.

Both sampled and digital images feature their own advantages and disadvantages. While storage of a digital image is simple¹, its discrete range imposes difficulties on the processing of the image using divisions and multiplications. Divisions and multiplications are easily performed for sampled images, but the real range [0, 1] can only be modeled approximatively. Fortunately, there is no need to restrict ourselves to one of the two models. We can easily convert a sampled image into a digital one and vice versa using below mappings:

$$\begin{array}{rcccc} D_g \colon & [0,1] & \to & \{0,\dots,g\} \\ & x & \mapsto & \lfloor g \cdot x \rfloor \end{array}$$
(1.1)

$$\begin{array}{rccc} S_g \colon & \{0, \dots, g\} & \to & [0, 1] \\ & x & \mapsto & \frac{x}{q} \end{array} \tag{1.2}$$

 D_g converts a sampled image into a digital image. S_g performs the conversion in the opposite direction.

1.1.2 Bit depth of digital images

As stated in Definition 1.3, digital images feature a discrete range of gray levels $\{0, \ldots, g\}$ with $g \in \mathbb{N}$. Due to the low-level representation of images, g is often given by

$$g = 2^d - 1$$

with $d \in \mathbb{N}$. This way, a single pixel may be stored in and represented by d bits. This naturally leads to the following definition:

Definition 1.4. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ denote a digital image with $g = 2^d - 1$ as well as $m, n, d \in \mathbb{N}$. Then d is called bit depth of the image f.

¹In fact, modern imaging hardware usually provides the user with digital images.



Figure 1.1: Black to white gradients for bit depths d = 6 through d = 9. Notice the vertical edges in Subfigures (a) and (b). In contrast, Subfigures (c) and (d) appear smooth to the human eye.

The attentive reader might wonder what choices for bit depths d are suitable for monochromatic images. Before deciding on a suitable bit depth, we have to define what we mean by the term "suitable". At this point, we have to trade off between two conflicting goals. On the one hand, we want to be able to model smooth images. That is, a black to white gradient contained in an image should not feature noticeable "edges". Greater bit depths allow more subdivisions of the black to white range. Therefore, greater bit depths allow the modeling of smoother images. Furthermore, greater bit depths allow higher computational accuracies when processing digital images. On the other hand, we do not want to waste memory when storing images digitally. Since the memory consumption of an image scales quadratically in its bit depth, we want to keep the bit depth as small as possible. To conclude, we have to trade off smoothness and computational accuracy against memory consumption.

Let $d_{min} \in \mathbb{N}$ denote the minimal bit depth which allows representation of smooth black to white gradients. We decide on d_{min} by plotting black to white gradients for various bit depths as in Figure 1.1. We notice the vertical edges Figures 1.1a and 1.1b feature. Since the edges disappear in case of Figure 1.1c, we define:

Proposition 1.1. Let $d_{min} \in \mathbb{N}$ denote the minimal bit depth which allows representation of smooth black to white gradients. We set

$$d_{min} := 8$$

Definition 1.1 is crucial to the display of images. The black to white gradient appears smooth to the human eye when plotted using a bit depth of $d_{min} = 8$. Hence, there is no need to aim for greater bit depths when merely showing an image to the human viewer. The human eye is not capable of recognizing all the fine grained variation in images of greater bit depth. Even if the data we process allows greater bit depth, we may eventually safely fall back to d_{min} when presenting our results to the human eye. Following Definition 1.1, during this fall back we do not lose any information recognizable by the human eye.

1.1.3 Gray level transforms

So far, we have discussed how to model images in the field of digital image processing. We have not yet discussed how to actually alter or enhance images. One of the most basic techniques is to alter an image's gray levels pixel-wise. **Definition 1.5.** Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \to \mathbb{G}$ denote a sampled $(\mathbb{G} = [0,1])$ or digital $(\mathbb{G} = \{0, \ldots, g\}$ with $g \in \mathbb{N})$ image. Let further $T: \mathbb{G} \to \mathbb{G}$ denote a mapping. Then the composition

$$\begin{array}{ccc} T \circ f \colon & \{0, \dots, m-1\} \times \{0, \dots, n-1\} & \to & \mathbb{G} \\ & & (i, j) & \mapsto & T\left(f\left(i, j\right)\right) \end{array}$$

is called application of the gray level transform T to the image f.

1.2 Color images

Most of this thesis is based on the processing and examination of monochromatic images. Nonetheless, we will encounter situations, where the display of information using a color image outperforms the display of the same information using a monochromatic image. Let us therefore give a brief explanation on how to model a color image.

A color image is simply a stack of monochromatic images. Within this stack, every image represents a color channel. A widespread encoding of color images is a stack of three images representing the colors red, green and blue.

Definition 1.6. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \to \mathbb{G}^3$ denote a stack of three sampled or digital images. Let further $f_k: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \to \mathbb{G}$ with $k = 1 \ldots 3$ denote the mappings given by

$$\begin{array}{rcl} f_k: & \{0,\ldots,m-1\} \times \{0,\ldots,n-1\} & \to & \mathbb{G} \\ & & (i,j) & \mapsto & f(i,j)_k \end{array}$$

Then f is called a color image. If f_1, f_2 and f_3 represent the red, blue and green color channel of the color image f, then f is called a color image in RGB-encoding.

Chapter 2

Fundamentals of multispectral imaging

At this point, we know how to model images. We have not yet discussed how to obtain images though. This section is dedicated to the process of image acquisition. Put in simple words, multispectral imaging is a subtype of photography. Photography, in turn, is based on capturing light. Consequently, in Section 2.1, we discuss the physical foundations of light in general. We discuss the characteristics of multispectral imaging in contrast to usual photography in Section 2.2. Section 2.3 gives a brief link between color images and multispectral imaging. Section 2.4 concludes this section with a description the experimental setup used throughout this thesis.

2.1 Fundamentals of light

Photography captures the light received from an object using light-sensitive hardware. We distinguish two types of photography in general:

- (1) In case of incident light exposure, the sensitive hardware captures the light reflected by the object of interest (cf. Figure 2.1a). Thus, the resulting image depicts mainly the surface of the object.
- (2) In case of transmitted light exposure, the sensitive hardware captures the light transmitted by the object of interest (cf. Figure 2.1b). Thus, the resulting image depicts the integral through the inner of the object.

The term "light" usually refers to the visible range of electromagnetic radiation. Here, the radiation's wavelength λ is of crucial importance for two reasons.

First, the wavelength determines, whether a specific type of radiation is visible to the human eye or not. The visible spectrum of electromagnetic radiation ranges from wavelengths of approximately 400 nm to 700 nm (cf. [6]). At $\lambda = 400$ nm and $\lambda = 700$ nm there is no sharp "cut" in terms of visibility of radiation. Instead, the visible light smoothly blends into radiation called ultraviolet radiation (approximately 300 nm to 400 nm) and infrared radiation (approximately 700 nm to 1500 nm).

Second, the wavelength determines the color a human perceives visible radiation in. White light consists of a mixture of radiation with wavelengths evenly distributed over the whole visible spectrum. "Colored" light either features radiation of wavelengths within some predominant,



Figure 2.1: Incident light versus transmitted light exposure

Table 2.1: Approximative correspondence between wavelength and perceived color according to Smith

wavelength [nm], approximatively	perceived color
390 to 420	near ultraviolet to violet
420 to 490	blue
490 to 510	blue-green
510 to 550	green
550 to 575	yellow-green
575 to 590	yellow
590 to 650	orange
650 to 750	red to near infrared

narrow sub-range of the visible spectrum. Or it features a lack of radiation of wavelengths corresponding to its complementary color. For example, red light either consists of radiation featuring wavelengths around $\lambda = 700 \text{ nm}$ ("red wavelengths") or it consists of radiation featuring all possible wavelengths except for wavelengths around $\lambda = 520 \text{ nm}$ ("green wavelengths"). Table 2.1 approximatively maps wavelengths to colors.

Both for transmitted and reflected light, the color an object appears to the human viewer depends on the type of radiation received from the object. For the sake of simplicity, we focus our explanation on reflected light. Nonetheless, the very same explanation applies for transmitted light as well. Objects can either reflect or absorb incident radiation. The reflectance and absorbance can be either leveled over the whole visible spectrum or favor radiation of specific wavelengths.

Gray objects feature leveled reflectance and absorbance over the whole visible spectrum. The amount of absorbance determines how dark objects appear. White and black objects establish the extreme cases at opposite ends of the scale. They reflect and absorb all incident radiation respectively.

In contrast, colored objects favor some radiation of specific wavelengths. Equivalent to "colored" light, there are two possibilities for an object to appear of a certain color. First, the object can reflect radiation of a specific wavelength only. Second, the object can reflect all radiation except for radiation "of complementary color".

The interrelation of absorbance and reflectance of an object and the characteristics of incident light enable us to both hide and emphasize details of an imaged object. To emphasize this interrelation, we briefly discuss the example depicted in Figure 2.2. The Figure depicts colored pencil of various colors on white paper (cf. Figure 2.2a).



(g) 590 nm

(h) 627 nm

(i) 655 nm

Figure 2.2: The influence of exposing radiation to the contrast of imaged objects. The images depict colored pencil on white paper. Due to different reflectance and absorbance, different colors appear of different contrast in different images.

We first consider the line featuring yellow colored pencil. The pencil receives its color by reflecting all but the "complementary, blueish light". When exposed to white light, the writing appears as simple writing on white background (cf. Figure 2.2b). When exposed to blueish light, the writing absorbs most incident light (cf. Figures 2.2c and 2.2d). Because of being white, the paper reflects all incident blueish light. Therefore, the writing appears with maximal contrast in the obtained images. When exposed to light of greater wavelengths, the writing reflects most incident light (cf. Figures 2.2f through 2.2i). Because of being white, the paper reflects all incident light of greater wavelengths as well. Consequently, the writing disappears in the obtained images.

We second consider the line featuring red colored pencil. In contrast to the yellow colored pencil, the red colored pencil receives its color by absorbing all but "reddish light". Thus, when exposed to light of short wavelengths, the pencil appears with maximal contrast in the obtained images (cf. Figures 2.2c through 2.2e). The greater the exposing light's wavelength gets, the more light the writing reflects. Hence, the greater the wavelength, the more the writing disappears from the obtained images (cf. Figures 2.2g through 2.2i).

2.2Characteristics of multispectral imaging

For usual photography, the hardware is sensible to the whole visible spectrum. Hardware sensible to infrared or ultraviolet radiation is rather uncommon. The object of interest is exposed with light consisting of a broad range of radiation of differing wavelengths. All the radiation of differing wavelengths is captured within a single image.

Multispectral imaging reverts the ratio of the number of radiation wavelengths to images. Per image, radiation of only either a single wavelength or a narrow band of wavelengths is captured. To scan the whole visible spectrum, multiple images are taken. As discussed in Section 2.1, this approach allows the analysis of different details in contrast at different wavelengths. In the following, we call a set of images of the same object obtained using all supported exposure radiation wavelengths a *multispectral series*.

To the end of image acquisition, specialized hardware is used. Usually, as with usual photography, the camera is sensitive to radiation in the whole visible spectrum. Beyond the visible spectrum, the sensitivity often extends far into the infrared and ultraviolet spectrum. To obtain multiple images of differing exposure, two approaches can be taken. First, filters attached to the camera can restrict the detectable radiation to a narrow range in terms of wavelengths. Second, the exposing hardware may be specialized. Here, the exposing hardware emits radiation of a single or few wavelengths.

Both infrared and ultraviolet exposure serve their very unique purposes. First, according to Smith, infrared radiation features little to no scatter with minute structures of matter. Thus, infrared radiation penetrates shallow surfaces. This enables infrared radiation to scan regions beneath the surface even for reflected light exposure. Second, ultraviolet radiation is valuable for the examination of parchment. According to Hollaus et al., parchment is fluorescent. Therefore, writing of ink appears of great contrast when exposed using ultraviolet radiation.

The range of possible applications of multispectral imaging features a vast set of different objects. In [3], [5] and [11], the multispectral approach was used to recover several types of writing. The three publications focus on the examination of palimpsests, the enhancement of faint text and the recovery of censored text respectively. In [1] and [2], the authors examined paintings. Using the multispectral approach, pigment composition was analyzed and underdrawings were revealed. All mentioned applications feature two similarities and emphasize two main advantages of the multispectral approach:

- (1) Due to the historical value of the examined objects, invasive physical or chemical examinations were not worthy of consideration. In contrast, the multispectral approach establishes a non-invasive tool applicable even to objects of poor physical constitution.
- (2) The multispectral approach is capable of capturing fine details detectable only with exposure of discrete or narrow-ranged wavelengths. While the details remain physically present with exposure of white light, they get lost in the superposition of the various types of reflected radiation.

2.3 Comparison of color images and multispectral series

In a sense, the color images from Section 1.2 were the very first multispectral data we encountered throughout this thesis. According to the Commission Internationale de l'Eclairage (CIE, cf. [6]), red, green and blue establish the so-called *primary colors*. The primary colors correspond with radiation of specific wavelength λ . Table 2.2 lists these correspondences. Red, green and blue are chosen as primary colors, because the human eye features distinct receptors for the three of them.

The color channels of a RGB-encoded image directly map to above primary colors. Other colors can not be modeled directly. Instead, other colors are modeled by superposition of the Table 2.2: Primary colors and their corresponding radiations' wavelengths according to the CIE

primary color	wavelength λ [nm]
red	700
green	546.1
blue	435.8

Table 2.3: Filter types and respective approximative transmissivities used throughout the thesis

filter type	transmissivity
ultraviolet-block (UVB)	$unknown^1$
ultraviolet-pass (UVP)	unknown
yellow (Y22)	Wratten number 22

colors red, green and blue. The acquisition of a general purpose color image is usually implemented in a similar manner. The visible spectrum is divided into three wavelength bands which are simultaneously captured independently from each other.

Dedicated multispectral imaging tries to avoid the superposition of different colors. Consequently, a much greater number of distinct wavelengths is considered upon image acquisition. To model dedicated multispectral series as a single color image, the image has to feature a distinct color channel for each wavelength. The set of wavelengths considered upon multispectral imaging is not standardized. Hence, there is no dedicated encoding for multispectral series as single color image. Instead, multispectral series are usually represented as set of monochrome images.

2.4 Image acquisition setup

2.4.1 Hardware

As indicated by Section 2.2, specialized hardware eases the acquisition of multispectral data. This thesis relies on hardware manufactured by MegaVision.

Throughout the thesis, the "E6-mono" camera back was used. The "E6-mono" is a full frame camera back providing images of $5412 \times 7216 \,\mathrm{px}$. Images feature a bit depth of $12 \,\mathrm{bit}$. The "E6-mono"'s sensor is a CCD sensor and thus sensitive to radiation from $350 \,\mathrm{nm}$ to $1050 \,\mathrm{nm}$. Here, the sensor's pixels capture mere intensities. In contrast to general purpose sensors, the pixels do not feature subpixels discriminating radiation of different wavelengths.

The discrimination of different wavelengths is implemented by two supporting MegaVision LED-panels. The panels' supported wavelengths are 365 nm, 450 nm, 470 nm, 505 nm, 530 nm, 590 nm, 627 nm, 655 nm and 850 nm. Here, each wavelength can be emitted independently from others. The panels further support white light emission.

Attached in front of the lens, MegaVision's filter wheel provides even more exposure possibilities. The wheel features three types of filters or may be disabled at all. The filters and their respective transmissivities are given in Table 2.3.

2.4.2 Software

Early imaging hardware forced the photographer to manually exchange filters in order to obtain images of differing exposure. This manual exchange caused small shifts in the position of the

 $^{^{1}}$ We asked MegaVision's sales TechVision for clarification of the ultraviolet filters' transmissivities. Unfortunately, the answer is still pending as of the time of writing.

imaging hardware. Therefore, images obtained during a single multispectral series had to be registered before analyzing them. This registration process assured the invariance of the mapping of pixel position to physical position on the object.

The hardware used throughout this thesis interoperates with external input hardware. I.e. the shutter release, the wavelength selection and the filter selection may be performed using software shipped by MegaVision. This eliminates the need for further physical interaction with the imaging hardware once the hardware has been calibrated. Thus, no shift of the object relative to the imaging hardware occurs between the acquisition of several images of differing exposure. Consequently, no further image registration-process is needed.

The software further supports two types of correction. "Flat field calibration" corrects for variations within a single image. The corrected variations include inhomogeneous lighting and lens fall-off. "Light balancing" corrects for variations within a single series of multispectral images. The balancing accounts for varying brightness due to differing exposing radiation. After applying this balancing, the obtained images feature roughly the same brightness.

2.4.3 Setup

Figure 2.3a depicts the overall imaging setup for incident light exposure. The camera is mounted onto a copy stand facing vertically downwards. The filter wheel is attached to the camera's lens.

The object of interest is positioned on the base. The object is placed on top of black cardboard to prevent reflections from underlying surfaces. If present and possible, parts of the object not imaged throughout a multispectral series are covered by black cardboard as well.

Some sheets of parchment and paper imaged throughout this thesis feature noticeable bending. In the extreme case, the bending exceeds the camera's depth of field. To cope with this bending, the upper layer of black cardboard is covered with a plate of anti-reflexive glass where needed.

The two lighting panels are positioned "to the left" and "to the right" of the copy stand. Both panels are positioned facing towards the center of the focal plane. The panels' direction and the focal plane enclose an angle of 45°. As depicted in Figure 2.3a, both panels feature a diffuser.

Due to external restrictions, the copy stand is positioned in front of a white wall. Throughout the early stage of this thesis, obtained images featured inhomogeneous exposure due to radiation reflected by the wall. To cope with these reflections, a black curtain is placed between the copy stand and the wall.

In case of transmitted light exposure, the two lighting panels are replaced by a single "Dörr LP-1218 LED" lighting panel. This time, the object is placed directly on top of the panel without a layer of cardboard beneath it. The further setup matches the setup of incident light exposure and is depicted in Figure 2.3b.

In case of both incident and transmitted light exposure within the same multispectral series, both aforementioned setups are fused into a single setup (cf. Figure 2.3c). This fused setup permits the acquisition of both incident and transmitted light exposure images. Thus, the setup removes the need for further physical interaction throughout a single series. Again, the avoidance of physical interaction serves the elimination of an image registration step.

2.4.4 Acquisition of a multispectral series

The acquisition of a multispectral series is similar to a general purpose photo shoot. First, the object of interest is placed beneath the camera. Second, the camera is focused on the object.

Third, the multispectral series is shot. Unless otherwise stated, a series consists of a total of 13 images. Table 2.4 lists the different exposure configurations.

Some objects including the urbarium are imaged in several patches depicting different parts of the object. This spread over several patches allows us to increase the number of pixels per physical area. Consequently, the camera captures finer details of the imaged object. In other words, the spread over several patches allows us to image the object in greater detail.



Table 2.4: Default exposure configurations used throughout a single multispectral series

exposure type	exposing radiation's wavelength [nm]	filter
monochromatically exposed	365	-
monochromatically exposed	450	-
monochromatically exposed	470	-
monochromatically exposed	505	-
monochromatically exposed	530	-
monochromatically exposed	590	-
monochromatically exposed	627	-
monochromatically exposed	655	-
monochromatically exposed	850	-
white light exposed	-	-
white light exposed	-	"UVB"
white light exposed	-	"UVP"
white light exposed	-	"Y22"

Chapter 3

Recovery of illegible parts of historical manuscripts and Principal Component Analysis

Now that we know the theoretical foundations of multispectral imaging, let us apply our knowledge in practice. To this end, we need an object to image. This thesis aims at developing a *robust* analysis for the examination of the urbarium. We base our development on

- (1) a sound understanding of the tools used throughout the analysis and
- (2) the application of our analysis to objects with known outcome.

(1) eases the actual development of the analysis. (2) lets us compare our results with some baseline and therefore enables us to critically reflect our results. While (1) is possible by working exclusively with the urbarium, (2) is certainly easier by gradually increasing the complexity of the imaged objects. Before examining the urbarium, we hence focus on instructive toy examples first.

This section is dedicated to a typical task in the field of historical document examination: the recovery of illegible parts of a manuscript. The rest of this section is structured as follows. In Section 3.1, we present the running example of this section. Section 3.2 is dedicated to the pre-processing steps of our analysis. We introduce the concepts of Principal Component Analysis in Section 3.3. To conclude, we apply the Principal Component Analysis to the running example and discuss our findings in Section 3.4.

3.1 The running example: parchment manuscripts

A typical task in the field of historical document examination is the recovery of illegible parts of manuscripts. According to [4], historical manuscripts often suffer from various types of damages. The damages are caused by a number of reasons like mistakes during production, improper storage and improper use. Most mentioned types of damage decrease the legibility of manuscripts. Therefore, the recovery of illegible parts of manuscripts has become an omnipresent task in the field of historical document examination.

Throughout this chapter, we focus on a parchment manuscript kindly provided by the "Staatliche Bibliothek Passau". Figure 3.1 shows both the front and the back of the parchment. The parchment dates back to the first half of the second millennium. Both the front and the back feature



(a) Recto

(b) Verso

Figure 3.1: The recto and the verso of the parchment discussed throughout Chapter 3. Remains of glue frame the verso's text in Figure (b).

writing consisting of ink. Presumably, the ink present on the parchment is iron-gall-ink. Notice that we avoid the technically more correct terms "recto" and "verso" instead of "front" and "back" up to this point. Since the parchment is not embedded in a composed object like a book, it is not clear which side establishes the recto at first glance.

According to the "Staatliche Bibliothek Passau", the parchment was used as cover of a book in the past. Later on, the person responsible for the book decided to separate the parchment from the inner book. Due to the use as cover, the side we defined as back features remains of glue. Motivated by these remains, we conjecture that the side we defined as back in fact originally faced towards the verso of the book cover. Thus, the side we defined as front originally established the recto of the book cover.

Protected by the rest of the book cover, the parchment's verso sustained aging in best possible manner. The verso's text, an extract of the chapter "De sancto Martino episcopo" from "Speculum ecclesiae" by the Honorius of Autun, features high legibility. In contrast, the parchment's recto features characteristics of high manual use and wear off. Use and wear have rendered parts of the recto's text illegible. Figure 3.2 highlights legible and illegible parts of the recto's text.

In order to reconstruct the damaged parts of the recto's text, we image the recto using the multispectral approach discussed in Chapter 2. The acquired images are shown in Figure 3.3.

Before we have a close look at the images shown in Figure 3.3, let us introduce a term we will use throughout our discussion.

Definition 3.1. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ denote a digital image. Then the histogram H_f of f is defined as

$$\begin{array}{rcl} H_f \colon & \{0, \dots, g\} & \to & \mathbb{N} \\ & x & \mapsto & |\{(i, j) \in \{1, \dots, m\} \times \{1, \dots, n\} \mid f(i, j) = x\}| \end{array}$$

 H_f simply counts the absolute occurrences of a specific gray level x in the image f.



Figure 3.2: The recto of the parchment shown in Figure 3.1a with a focus on legibility. Legible and illegible parts are highlighted in yellow and blue respectively.

Several features of Figure 3.3 immediately catch our attention:

- As discussed in Section 2.1, different parts of the recto's text feature different contrast for different exposure wavelengths. Figure 3.4 emphasizes the contrast for the capital letter "Q" located in the upper right quarter of the parchment.
- (2) The multispectral approach itself already increases the legibility of the recto's text.
- (3) The images are not exposed homogeneously across the whole series. In other words, different images appear of differing overall brightness. This difference is present despite the use of "Light balancing" (cf. Subsection 2.4.2).
- (4) The images make use of only a small subspace of the available gray level space. This misuse of the available gray level space becomes even more apparent when looking at the histogram shown in Figure 3.5b.

On the one hand, (1) and (2) indicate that the multispectral approach is a good measure to tackle the reconstruction of the recto's text.

On the other hand, (3) and (4) hinder both the further visual and the further computer-based examination of the multispectral images.

3.2 Pre-processing of raw images

Fortunately, there is a number of measures to overcome the Shortcomings (3) and (4). We apply these measures before proceeding to the actual analysis of the multispectral series. In the following, we introduce the measures and compose them into a sequence of processing steps.



Figure 3.3: Multispectral series of images of the parchment's recto. The shown images represent the raw images provided by the used camera without further processing. The shown images feature a pixel density of approximately 910 dpi. Continued on Page 30.

CHAPTER 3. ILLEGIBLE MANUSCRIPTS AND PCA



(j) White, Filter none





(k) White, Filter UVB



(l) White, Filter Y22

Figure 3.3 continued: Multispectral series of images of the parchment's recto. The shown images represent the raw images provided by the used camera without further processing. The shown images feature a pixel density of approximately 910 dpi. Continuation from Page 29.

(m) White, Filter UVP

3.2.1 Histogram equalization

As depicted in Figure 3.3, the raw multispectral images of the parchment feature arguably poor contrast. This poor contrast is significantly caused by the poor use of the available gray level space. E.g. Figure 3.4a mainly features gray levels in $\{0, \ldots, 12\,000\}$. At the same time Figure 3.4a offers an available gray level space of $\{0, \ldots, 2^{16} - 1\}$. As a consequence, the "lightest" pixels of the discussed Figure appear of dark gray on black background to the human viewer.

In contrast, an image of high contrast makes use of the whole available gray level space. Furthermore, there are no predominant gray levels that appear significantly more often than other gray levels. Thus, the absolute frequencies of occurrences of the different gray levels do not vary much. In other words, for a digital image $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ of high contrast, we expect the histogram H_f to be approximately given by

$$\begin{array}{cccc} H_f \colon \{0, \dots, g\} & \to & \mathbb{N} \\ & x & \mapsto & \left| \frac{m \cdot n}{g + 1} \right| \end{array}$$

$$(3.1)$$

We now aim at finding a gray level transform. This gray level transform shall transform a real-world digital image $h: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ into a digital image $\tilde{h}: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ with $H_{\tilde{h}} \stackrel{!}{=} H_f$. To this end, let us introduce another term closely related to Definition 3.1.



Figure 3.4: Multispectral series of images of the parchment's recto with a focus on differing contrast. The shown images are patches taken from the images shown in Figure 3.3. All patches show the same area around the capital letter "Q" located in the upper right quarter of the parchment's recto.

(m) White, Filter UVP

Definition 3.2. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ denote a digital image. Then the cumulative histogram A_f of f is defined as

$$\begin{array}{rcl} A_{f}: & \{0, \dots, g\} & \to & \mathbb{N} \\ & x & \mapsto & \sum_{k=0}^{x} H_{f}\left(k\right) = \sum_{k=0}^{x} |\{(i,j) \mid f\left(i,j\right) = k\}| \\ & = |\{(i,j) \mid f\left(i,j\right) \leq x\}| \end{array}$$

 A_f counts the absolute occurrences of gray levels below a specific threshold x in the image f.

Equation 3.1 can be almost¹ equivalently expressed in terms of A_f as

$$\begin{array}{cccc} A_f: & \{0, \dots, g\} & \to & \mathbb{N} \\ & x & \mapsto & \left\lfloor \frac{m \cdot n}{g+1} \cdot (x+1) \right\rfloor \end{array} \tag{3.2}$$

In other words, given a gray level x we consider the ratio of absolute occurrences of gray levels $x' \leq x$ and the total area of the image $m \cdot n$. We further consider the ratio of x + 1 and g + 1. For an ideal image, we expect both ratios to be roughly equal. That is:

$$\frac{A_f(x)}{m \cdot n} \approx \frac{x+1}{g+1} \tag{3.3}$$

We can force a real-world image into satisfying Equation 3.3 by mapping a given gray level x onto the product of g and the relative frequency of occurrence of gray levels less than or equal to x. The resulting gray level transform is defined in Definition 3.3.

Definition 3.3. Let $f: \{0, ..., m-1\} \times \{0, ..., n-1\} \rightarrow \{0, ..., g\}$ denote a digital image. f induces a gray level transform called histogram equalization based on f defined as

$$\begin{array}{rccc} T_f^{\texttt{equ}} \colon & \{0, \dots, g\} & \to & \{0, \dots, g\} \\ & x & \mapsto & \left\lfloor \frac{A_f(x)}{m \cdot n} \cdot g \right\rfloor \end{array}$$

The application of the gray level transform T_f^{equ} to the image f is called Equalization of f's histogram.

Figure 3.5c depicts the equalized version of the image shown in Figure 3.3a alongside the equalized histogram. Figure 3.6c shows a further, artificial example of the effects of histogram equalization applied to Figure 3.6a. Clearly, histogram equalization improves the images' contrast.

3.2.2 Gray level range maximization

Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ denote a digital image and \tilde{f} its histogram equalized counterpart. \tilde{f} features one last caveat. By definition, we find:

$$\max_{\substack{i \in \{0,...,m-1\}\\j \in \{0,...,n-1\}}} \hat{f}(i,j) = \max_{\substack{x \in \{0,...,g\}\\}} T_f^{equ}(x)$$
$$= T_f^{equ}(g)$$
$$= g$$

as well as

$$\min_{\substack{i \in \{0,\dots,m-1\}\\j \in \{0,\dots,n-1\}}} \tilde{f}(i,j) = \min_{\substack{x \in \{0,\dots,g\}}} T_f^{\mathsf{equ}}(x)$$
$$= T_f^{\mathsf{equ}}(0) \\= \left\lfloor \frac{A_f(0)}{m \cdot n} \cdot g \right\rfloor$$

Assured by definition of T_f^{equ} , \tilde{f} makes use of gray levels near the upper bound g. In contrast, \tilde{f} will not make use of gray levels near the lower bound 0, if f features a significant portion of pixels of gray level 0. Figure 3.6c shows an illustrative example of this behavior. While \tilde{f} features gray levels $x \in \{32, \ldots, 255\}$, gray levels $x \in \{0, \ldots, 31\}$ do not occur at all. To force \tilde{f} 's minimal and maximal gray levels into matching 0 and 255 respectively, we apply the linear mapping introduced in Definition 3.4.

Definition 3.4. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow [0, 1]$ denote a sampled image. Let further

$$g_{\min} := \min_{\substack{i \in \{0, \dots, m-1\}\\ j \in \{0, \dots, n-1\}}} f(i, j)$$
$$g_{\max} := \max_{\substack{i \in \{0, \dots, m-1\}\\ j \in \{0, \dots, n-1\}}} f(i, j)$$

Then the gray level transform T_{f}^{range} defined as

$$\begin{array}{rccc} T_f^{\texttt{range}} \colon & [0,1] & \to & [0,1] \\ & x & \mapsto & \frac{x - g_{\min}}{g_{\max} - g_{\min}} \end{array}$$

is called gray level range maximization based on f.

¹ If we were discussing images with a range $\mathbb{I} \subset \mathbb{R}$, Equations 3.1 and 3.2 would in fact be equivalent. As we are discussing images with a range $\mathbb{I} \subset \mathbb{N}$, the application of $\lfloor \cdot \rfloor$ introduces slight inaccuracies which prevent equivalence.



(a) Original image; reprint from Figure 3.3a.



(c) Histogram equalized version of (a). Clearly, the image's contrast profits from the application of histogram equalization.



(b) Histogram of (a). The predominant gray levels are contained in $\{0,\ldots,12\,000\}$ although (a) features an available gray level space of $\{0,\ldots,2^{16}-1\}.$



(d) Histogram of (c)



(e) Result of subsequent application of histogram equalization and gray level range maximization to (a)



Figure 3.5: Effects and interplay of histogram equalization and gray level range maximization applied to a real-world image (cf. Figure 3.3a)



histogram equalization and gray level range maximization to (a)

Figure 3.6: Effects and interplay of histogram equalization and gray level range maximization applied to an artificial image

3.2. PRE-PROCESSING OF RAW IMAGES

Concerning the gray level spaces used throughout the definition of T_f^{equ} and T_f^{range} we note: Divisions are far easier to perform over the interval [0,1] instead of the discrete set $\{0, \ldots, g\}$. Hence, we use the gray level space [0,1] in Definition 3.4. Nonetheless, we can not use [0,1]as gray level space in T_f^{equ} 's defining equation, as [0,1] does not induce any histogram bins to base the definition on. In contrast to [0,1], $\{0,\ldots,g\}$ does induce bins. As we have seen in Subsection 1.1.1, there is no need to restrict ourselves to one of the two mentioned gray level spaces. Instead, we can easily change an image's representation as needed.

Figures 3.5e and 3.6e depict the effect of gray level range maximization. As with histogram equalization, the images' contrast profit from the application of gray level range maximization.

3.2.3 Normalization

The interplay of histogram equalization and gray level range maximization is targeted at enhancing an image's contrast (cf. item (4) from the introduction of Section 3.1). In opposition, both techniques are not dedicated to the differing exposure across the series shown in Figure 3.3 (cf. item (3)). During the analysis of the images, we do not want to take into account the differing exposure. We therefore try to get rid of the differences. To this end, we introduce a quantity measuring the "energy" contained in an image. Here, the "energy" depends on the exposure of an image. Given this "energy", we can then fix the "energies" of all images under consideration to 1. As quantity measuring the "energy", we simply use the Euclidean norm as defined in Definition 3.5.

Definition 3.5. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ denote a digital image. f's Euclidean norm $||f||_2$ is defined as

$$\| f \|_{2} := \sqrt{\sum_{i=1}^{m} \sum_{j=1}^{n} f(i,j)^{2}}$$

To fix the energy of a specific image to 1, we apply the gray level transform given in Definition 3.6.

Definition 3.6. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow [0, 1]$ denote a sampled image with $||f||_2 > 0$. The gray level transform

$$\begin{array}{cccc} T_f^{\operatorname{norm}} \colon & [0,1] & \to & [0,1] \\ & x & \mapsto & \frac{x}{\|f\|_2} \end{array}$$

is called normalization based on f.

Application of normalization usually leads to images which appear all black to the human viewer. Consequently, we omit the depiction of the effect of normalization in a similar manner to Figure 3.5. Nonetheless, the appearance does not hinder the computer-based further analysis of the images.

3.2.4 Pre-processing pipeline

With the tools of histogram equalization, gray level range maximization and normalization we overcome issues (3) and (4) discussed in the introduction of Section 3.1. As we re-encounter the mentioned issues several times throughout the rest of this thesis, we give the tools a summarizing name. Definition 3.7 gives the details.

Definition 3.7. When analyzing a set of images, we often encounter several transformations targeted at preparing the raw images for the actual analysis. All transformations performed before the application of the actual analysis are thus called pre-processing of the raw images. A specific sequence of pre-processing transformations is referred to as pre-processing pipeline.

At this point, our pre-processing pipeline consists of the following steps:

Pre-processing:	
(1) Histogram equalization	
(2) Gray level range maximization	
(3) Normalization	

As most pre-processing steps alter an image's Euclidean norm, we keep an eye on applying the Normalization as very last step of a pre-processing pipeline.

We are now ready to perform the actual analysis of the images shown in Figure 3.3.

3.3 Introduction of PCA

When imaging historical documents using a multispectral camera, we end up with a vast set of images as shown in Figure 3.3. The images represent a tremendous amount of information. In order to gain insights on the imaged object, we have to carefully examine the vast set of images. This analysis can be a tedious and difficult task. To simplify the task, we use a dimension reduction tool to condense all the information into a single image revealing the information we aim at. By the term "dimension reduction tool" we refer to the fact that data obtained from an arbitrary experiment usually does not feature independent variations in all possible dimensions. Instead, variations in one dimension correlate with variations in another dimension. Thus, the data may be uniquely described using a smaller, minimal set of dimensions instead of the initial set. Dimension reduction tools help us find a minimal set of dimensions needed to describe our data.

To make our point more clear, let us consider the following toy example. Imagine we image the period at the end of a sentence. The multispectral camera we use for the imaging provides us with the 4×4 -images shown in Figure 3.7. Given an arbitrary but fixed pixel in one of the images, the pixel may vary its gray level throughout the series of images. This variation is independent from variations of the other pixels. Thus, every pixel in the 4×4 -images establishes a dimension. Leaving apart the restrictions imposed by the available gray level space, the images can be represented as vectors in the vector space $\mathbb{R}^{4\cdot 4} = \mathbb{R}^{16}$. Carefully analyzing the images we recognize two properties:

- The four inner pixels always feature roughly the same gray level. Thus, to describe an image uniquely, we need to know the gray level of only a single of these four pixels. The gray levels of the remaining three pixels may be reconstructed from the fourth pixel's gray level.
- The outer 12 pixels share a similar property. Again, knowledge about one of the pixels' gray level lets us reconstruct the gray levels of all remaining pixels.

In conclusion, the data we encounter in Figure 3.7 features a lot of redundancy. The data can be embedded in an at most two-dimensional subspace of \mathbb{R}^{16} . Dimension reduction tools serve as
3.3. INTRODUCTION OF PCA



general purpose measure to extract the relevant dimensions along which a set of high-dimensional data varies. Here, reduction is not lossy. The extracted dimensions still suffice to reconstruct the original data. In other words, we can think of the dimension reduction as a basis transform. The transform maps the possibly inadequate basis in which we conduct our experiment to a basis more suitable for describing the obtained data.

Above example emphasizes our need for a dimension reduction tool when analyzing multispectral data. Before we continue our study of this example, let us decide on a tool. As stated by Giacometti et al. and Hollaus et al., there are several reasonable dimension reduction tools. These include Principal Component Analysis (PCA), Independent Components Analysis (ICA), Linear Spectral Mixture Analysis (LSMA) and Linear Discriminant Analysis (LDA). While PCA, ICA and LSMA are unsupervised dimension reduction tools, LDA is a supervised tool. LDA therefore needs to be trained on a set of labeled images prior to application on previously unknown images. Benefiting from the information contained in the training data, LDA outperforms the mentioned unsupervised tools according to the findings reported in [7]. The setup of a database of labeled images is beyond the scope of this thesis. Hence, despite the superior performance of LDA, we focus on unsupervised dimension reduction tools. Giacometti et al. provide a general purpose recommendation concerning dimension reduction tools. According to their findings in [4], PCA outperforms the other mentioned unsupervised tools in the average case. Furthermore, PCA reaches competitive performance in all other cases. Throughout this thesis, we thus focus on Principal Component Analysis as dimension reduction tool.

Let us now come back to above toy example. As we noted earlier, the quantities we monitor during our experiment are given by the pixels contained in a 4×4 -image. Every quantity we monitor, that is every pixel, induces a dimension of the space we embed our observations in. Independent of our methodology of experimenting, our observations may well be contained in a much smaller subspace. PCA as dimension reduction tool aims at finding a suitable basis of

the smaller subspace. The rest of this subsection is dedicated to an introduction to the concepts of PCA. The following is heavily inspired by [13] and [10]. Nonetheless, PCA is a well-studied topic and the interested reader is pointed to literature such as [9].

3.3.1 Redundancy elimination

Let us recall our main observation of our first study of the images in Figure 3.7. The shown images feature a high degree of redundancy. In order to get rid of it, we first have to introduce a measure of the redundancy contained in our data. PCA chooses the quantity of covariance as measure. It is defined as follows.

Definition 3.8. Let $n \in \mathbb{N}$ denote the number of trials we obtain during an experiment. Let $\mathbf{x} \in \mathbb{R}^n$ denote a quantity we observe. Let the subscript index $j = 1 \dots n$ further denote the specific trial we refer to in a given context. Then the mean \overline{x} of \mathbf{x} is defined as

$$\overline{x} := \frac{1}{n} \sum_{j=1}^{n} x_j$$

Unless otherwise stated, $\mathbf{x}, \mathbf{y}, \mathbf{x}_1, \ldots, \mathbf{x}_m$ for a given $m \in \mathbb{N}$ denote quantities as in Definition 3.8 throughout the rest of this section. Here, x_{ij} refers to the *j*-th trial of the *i*-th quantity. **Definition 3.9.** The variance $\sigma_{\mathbf{x}}^2$ of \mathbf{x} and the covariance $\operatorname{cov}(\mathbf{x}, \mathbf{y})$ between \mathbf{x} and \mathbf{y} are defined as

$$\sigma_{\mathbf{x}}^2 := \frac{1}{n} \sum_{j=1}^n (x_j - \overline{x})^2$$
$$\operatorname{cov}(\mathbf{x}, \mathbf{y}) := \frac{1}{n} \sum_{j=1}^n (x_j - \overline{x}) (y_j - \overline{y})$$

The elements of \mathbf{x} are located in a range around the mean \overline{x} . The variance $\sigma_{\mathbf{x}}^2$ is a measure for the narrowness of this range. The greater the variance, the broader the range.

Let us now come back to eliminating redundancies. To this end, we study the linear prediction of one quantity based on another.

Proposition 3.1. Let $\sigma_{\mathbf{x}}^2 > 0$. Let further

$$\mathbf{b} := \begin{pmatrix} b \\ \vdots \\ b \end{pmatrix} \in \mathbb{R}^n$$
$$a^* := \frac{\operatorname{cov}\left(\mathbf{x}, \mathbf{y}\right)}{\sigma_{\mathbf{x}}^2}$$
$$b^* := \overline{y} - a^* \cdot \overline{x}$$
$$\tilde{\mathbf{y}}_{ab} := a\mathbf{x} + \mathbf{b}$$

where $a, b \in \mathbb{R}$. $\tilde{\mathbf{y}}_{ab}$ is called linear prediction of \mathbf{y} based on \mathbf{x} . Then

$$\min_{a,b\in\mathbb{R}} \frac{1}{n} \sum_{j=1}^{n} (\tilde{y}_{ab,j} - y_j)^2 = \frac{1}{n} \sum_{j=1}^{n} (\tilde{y}_{a^*b^*,j} - y_j)^2$$

3.3. INTRODUCTION OF PCA

Instead of stating the proof, let us reformulate Proposition 3.1 in prose:

- (1) The covariance establishes a measure of linear relationship between \mathbf{x} and \mathbf{y} . cov $(\mathbf{x}, \mathbf{y}) > 0$ indicates that linearly increasing values of \mathbf{x} tend to coincide with linearly increasing values of \mathbf{y} . cov $(\mathbf{x}, \mathbf{y}) < 0$ indicates that linearly increasing values of \mathbf{x} tend to coincide with linearly decreasing values of \mathbf{y} . Nonetheless, strong coincidence indicated by the covariance does not imply any causal relationship between \mathbf{x} and \mathbf{y} .
- (2) In order to minimize the mean squared error of the linear prediction of \mathbf{y} based on \mathbf{x} , we have to take into account their covariance. In case $cov(\mathbf{x}, \mathbf{y})$ vanishes, we obtain

$$\tilde{\mathbf{y}}_{a^*,b^*} = \overline{\mathbf{y}} \quad \text{where} \quad \overline{\mathbf{y}} := \begin{pmatrix} \overline{y} \\ \vdots \\ \overline{y} \end{pmatrix} \in \mathbb{R}^n$$

Thus, \mathbf{x} is of no use to predict \mathbf{y} . To put it in the context of dimension reduction: non-zero covariance indicates some degree of redundancy between \mathbf{x} and \mathbf{y} . In contrast, vanishing covariance corresponds with vanishing redundancy between \mathbf{x} and \mathbf{y} .

(2) establishes the overall task of PCA. In order to select a redundancy-free, minimal set of dimensions, we calculate the covariances of the initial quantities corresponding to the initial dimensions. We then try to base our selection on these covariances. For a small number of dimensions this task might be performed by hand. But our camera provides us with images of resolution $5412 \times 7216 \,\mathrm{px}$, resulting in $39\,052\,992$ dimensions. Clearly, we need a proper approach to target the above task. To this end, we make use of some basic linear algebra.

Definition 3.10. Let $\sigma_{\mathbf{x}_1}^2, \ldots, \sigma_{\mathbf{x}_m}^2 = 0$. The covariance matrix $\mathbf{C}_{\mathbf{X}}$ of $\mathbf{x}_1, \ldots, \mathbf{x}_m$ is defined as

$$\mathbf{C}_{\mathbf{X}} := \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

where

$$\mathbf{X} := (x_{ij})_{\substack{i=1...m\\j=1...n}}$$

For the rest of this section, we assume $\sigma_{\mathbf{x}_1}^2, \ldots, \sigma_{\mathbf{x}_m}^2 = 0$. Even for real-world data, we can easily assure $\sigma_{\mathbf{x}_1}^2, \ldots, \sigma_{\mathbf{x}_m}^2 = 0$ by subtraction of the per dimension mean.

Definition 3.10 arranges the quantities $\mathbf{x}_1, \ldots, \mathbf{x}_m$ in an $m \times n$ -matrix². Here, the rows of **X** correspond to the set of quantities. The columns of **X** correspond to the set of trials obtained during the experiment. The terms $C_{\mathbf{X},ij}$ then yield the covariances $\operatorname{cov}(\mathbf{x}_i, \mathbf{x}_j)$ for $i, j \in \{1, \ldots, m\}, i \neq j$. Further, the terms $C_{\mathbf{X},ii}$ yield the variances $\sigma_{\mathbf{x}_i}^2$ for $i \in \{1, \ldots, m\}$. The introduction of the covariance matrix allows us to reformulate the overall task of the PCA:

Find a basis transform, such that $\mathbf{C}_{\mathbf{X}}$ becomes a diagonal matrix.

 $^{^{2}}$ We note that this arrangement is always possible, even if the monitored quantities relate to each other spatially. In our case, we can simply linearize the obtained images into a one-dimensional column vector. At this step, we remember the mapping from pixel position to vector element. Later, we recover the image by delinearizing the column vector into a matrix using the mapping above.

3.3.2 Linear basis transform

So far, we have not specified the type of basis transform we take into consideration. In order to precisely define the type of transform, let us introduce some more terms from the field of linear algebra.

Definition 3.11. (1) $\mathbf{A} \in Mat(\mathbb{R}, m, m)$ is called orthogonal, iff $\mathbf{A}^T \mathbf{A} = 1$.

- (2) A linear mapping $f : \mathbb{R}^m \to \mathbb{R}^m$ is called orthogonal, iff its corresponding transformation matrix is orthogonal.
- (3) $\mathbf{A} \in Mat(\mathbb{R}, m, m)$ is called symmetric, iff $\mathbf{A}^T = \mathbf{A}$.
- (4) A set $\{\mathbf{b}_1, \ldots, \mathbf{b}_m\}$ is called an orthonormal basis of \mathbb{R}^m iff

(a)
$$\{\mathbf{b}_1, \dots, \mathbf{b}_m\}$$
 is a basis of \mathbb{R}^m and
(b) $\sum_{i=1}^m b_{ij}^2 = 1$ $\forall i \in \{1, \dots, m\}$

In the framework of PCA, we choose to focus on a rather basic type of basis transforms. In fact, we choose to focus on orthogonal transforms.

The attentive reader might wonder whether the diagonalization of C_X using an orthogonal transform is possible at all. Fortunately, C_X is symmetric by definition. Proposition 3.2 assures that C_X can be orthogonally diagonalized.

Proposition 3.2. Let $\mathbf{A} \in \text{Mat}(\mathbb{R}, m, m)$ denote a symmetric matrix. Then there exists an orthogonal matrix $\mathbf{Q} \in \text{Mat}(\mathbb{R}, m, m)$ such that QAQ^T is diagonal.

Let us now examine how the trials we obtained during our experiment behave under orthogonal transformation. To this end, let $\mathbf{Q} \in \text{Mat}(\mathbb{R}, m, m)$ with \mathbf{Q} orthogonal. We further set $\mathbf{Y} := \mathbf{Q}\mathbf{X}$. Y_{ij} is the projection of the *j*-th trial \mathbf{x}_j onto the *i*-th row of \mathbf{Q} . Since \mathbf{Q} is orthogonal, the rows of \mathbf{Q} form an orthonormal basis of \mathbb{R}^m . Hence, every column of \mathbf{Y} yields a trial expressed with respect to the basis induced by \mathbf{Q} 's rows. In other words, if $\mathbf{C}_{\mathbf{Y}}$ is diagonal, then the rows of \mathbf{Q} induce the basis we are looking for. Definition 3.12 concludes our observation.

Definition 3.12. Let $\mathbf{Q} \in Mat(\mathbb{R}, m, m)$ with \mathbf{Q} orthogonal and $\mathbf{C}_{\mathbf{QX}}$ diagonal. Then the rows of \mathbf{Q} are called the principal components of \mathbf{X} .

With Definition 3.12 the overall task of PCA becomes:

Find a matrix $\mathbf{Q} \in Mat(\mathbb{R}, m, m)$ with \mathbf{Q} orthogonal, such that $\mathbf{C}_{\mathbf{QX}}$ is diagonal.

3.3.3 Q and Singular Value Decomposition

So far, we have restated the overall task of PCA over and over again. Each time, we introduced new assumptions helping us reach the final goal of an improved basis. In the following, we tackle the final version of the overall task.

For reasons that will become clear afterwards, let us introduce a matrix factorization known as *Singular Value Decomposition*.

Theorem 1. Let $\mathbf{A} \in Mat(\mathbb{R}, m, n)$ denote a matrix of rank $r_{\mathbf{A}}$. Then there exist matrices

- $\mathbf{U} \in Mat(\mathbb{R}, m, m)$ with U orthogonal,
- $\mathbf{V} \in Mat(\mathbb{R}, m, m)$ with V orthogonal and

- $\Sigma \in Mat(\mathbb{R}, m, n)$ with Σ diagonal. Furthermore, the elements Σ_{ii} for i = 1...min(m, n) are ordered descendingly and are non-negative. The number of strictly positive diagonal elements Σ_{ii} equals the rank $r_{\mathbf{A}}$.
- $\mathbf{A}, \mathbf{U}, \mathbf{\Sigma}$ and \mathbf{V} satisfy

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$

This factorization of \mathbf{A} is known as \mathbf{A} 's Singular Value Decomposition. Σ 's diagonal elements are called \mathbf{A} 's singular values.

In Subsection 3.3.2 we examined the behavior of trials under orthogonal transformations. Theorem 1 allows us to extend our examination to the corresponding covariance matrix. To this end, we set $\mathbf{Z} := \frac{1}{\sqrt{n}} \mathbf{X}^T$. By construction, we find

$$\mathbf{Z}^T \mathbf{Z} = \mathbf{C}_{\mathbf{X}}$$

We now factorize \mathbf{Z} as in Theorem 1 such that $\mathbf{Z} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$. This allows us to introduce $\mathbf{Y} := \mathbf{V}^T \mathbf{X}$. With above definitions, we finally find

$$C_{\mathbf{Y}} = \frac{1}{n} \mathbf{Y} \mathbf{Y}^{T}$$

$$= \frac{1}{n} \mathbf{V}^{T} \mathbf{X} \mathbf{X}^{T} \mathbf{V}$$

$$= \mathbf{V}^{T} C_{\mathbf{X}} \mathbf{V}$$

$$= \mathbf{V}^{T} \mathbf{Z}^{T} \mathbf{Z} \mathbf{V}$$

$$= \mathbf{V}^{T} \mathbf{V} \mathbf{\Sigma}^{T} \mathbf{U}^{T} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{T} \mathbf{V}$$

$$= \mathbf{\Sigma}^{T} \mathbf{\Sigma}$$

Since Σ is diagonal, the product $\Sigma^T \Sigma$ is diagonal as well. In conclusion, the rows of \mathbf{V}^T - that is the columns of \mathbf{V} - are the principal components of \mathbf{X} defined in Definition 3.12.

The attentive reader might wonder whether \mathbf{Y} equals \mathbf{X} 's independent component analysis. \mathbf{Y} 's rows are uncorrelated by construction. However, the rows are not independent in the general case. Hence, \mathbf{Y} does not necessarily equal \mathbf{X} 's independent component analysis.

3.3.4 Variances and restriction to relevant basis vectors

The attentive reader might further wonder what happened to the initial goal of actually *reducing* the number of dimensions. So far, we changed the basis with respect to which we represent our data. But the number of dimensions remained invariant.

At this point, the variances on the diagonal of C_Y come into play again. We recall two facts:

- (1) The variance $\sigma_{\mathbf{x}}^2$ of a quantity \mathbf{x} is a measure for the narrowness of the range in which its trials are located.
- (2) The diagonal elements of $\mathbf{C}_{\mathbf{Y}}$ are the variances of our data expressed with respect to the principal components. In other words, \mathbf{Z} 's singular values are the variances of our data expressed with respect to the principal components.

(1) results in our final assumption. Since we obtain our data through an experiment, our data will certainly feature some noise. As the PCA is not aware of the experimental setup, the PCA can not remove this noise. Instead, we will likely obtain principal components encoding the noise in

the data. But as we developed the experimental setup, we certainly designed it to maximize the signal-to-noise ratio. For this reason, we can expect the variance of our data to be mainly caused by the underlying links and laws we are studying. In other words, principal components with great variance encode important structure whereas principal components with small variance encode noise. Principal components featuring zero variance are not induced by the data but by the PCA-algorithm itself.

To summarize, Definition 3.13 concludes the ordering of principal components based on their respective variances.

Definition 3.13. The principal component of \mathbf{X} featuring greatest variance is called its first principal component. Its second, third, ... principal components are given by the components featuring the second-greatest, third-greatest, ... variance.

In order to focus our examination on the relevant variance of our data, we restrict our considerations to the principal components of greatest variance. Following (2) and Theorem 1, these components are given by the first columns of \mathbf{V} . Applying this restriction, we finally reduce the number of dimensions.

3.3.5 Summary and algorithm overview

To sum up our insights, let us restate the three main assumptions underlying the PCA. Let us further give a short outline of an algorithm performing PCA.

Assumptions:

- Covariance encodes redundancy.
- The basis transform is an orthogonal transform.
- Greater variances indicate important structure. Smaller variances indicate noise.

Algorithm:

- (1) Subtract the per dimension mean of and from the data.
- (2) Arrange the data in a "dimensions \times trials"-matrix **X**.
- (3) Calculate the Singular Value Decomposition of $\mathbf{Z} := \frac{1}{\sqrt{\text{trials}}} \mathbf{X}^T$.

3.3.6 Application of PCA to the toy example

Before we apply PCA to the running example presented in Section 3.1, let us briefly come back to the toy example presented in Figure 3.7. I.e. let us briefly check whether PCA is capable of recognizing the two dimensions of variance discussed in this section's motivation.

To this end, we first have to decide on how to encode principal components in images. As we have seen in Subsection 1.1.1, sampled and digital images feature a range $\mathbb{G} \subset \mathbb{R}_{\geq 0}$. Principal components do not represent images themselves. Instead, the components represent deviations from the mean of the data under consideration. As such, the components are given by vectors $\mathbf{q}_1, \ldots, \mathbf{q}_m \in \mathbb{R}^m$. Therefore, the components will almost certainly feature negative elements. These negative elements can not be encoded in digital images directly. We consider two approaches to overcome this limitation:

3.4. APPLICATION OF PCA TO THE PARCHMENT IMAGES

(1) We can map the principal component \mathbf{q}_i onto the interval [0,1] by element-wise application of below mapping t_i :

$$\begin{array}{rccc} t_i \colon & \mathbb{R} & \to & [0,1] \\ & x & \mapsto & 1 + \frac{x}{2\max\limits_{j \in \{1,\dots,m\}} |q_{ij}|} \end{array}$$

The drawback of this approach is that small deviations will be hard to distinguish from vanishing elements. Furthermore, the sign of small deviations will be hard to deduce from the image.

(2) We generate a pseudo color image for each principal component. In these images, yellow and blue pixels encode positive and negative deviations respectively. The brighter a pixel is, the greater the represented deviation is in magnitude.

To this end, we first scale every principal component such that all its elements are contained in the interval [-1, 1]. This scaling is performed by simple division by the maximum absolute element. Second, we map the absolute values to the set $\{0, \ldots, 255\}$. The image of the absolute values is then histogram equalized. Last but not least, the equalized image is mapped to a color-image using the signs of the original elements of the principal component.

Processing the absolute values instead of the signed values in step two assures that deviations of equal magnitude are encoded by pixels of equal brightness. Otherwise, the brightness of blue and yellow pixels might differ for deviations of equal magnitude.

Since (2) eases the distinction between negative and positive deviations, we choose (2) as the encoding of principal components.

The results of PCA applied to the toy example presented in Figure 3.7 are shown in Figure 3.8. As expected (cf. the motivation of this Section), the first and the second principal component encode the period as well as the background. Both principal components feature positive variances. Since components three through five feature vanishing variances³, subsequent components feature vanishing variances as well. Thus, components three through 16 are not needed to express the data shown in Figure 3.7. In conclusion, the original images are contained in a two-dimensional subspace of \mathbb{R}^{16} and PCA is capable of detecting the relevant dimensions.

3.4 Application of PCA to the parchment images

At this point, we are ready to apply PCA to the parchment images presented in Figure 3.3. Throughout this section, we introduce one final pre-processing step. Afterwards, we apply PCA and review the results.

As shown in Figure 3.3, the images currently feature dark writing on a light background. This depiction matches the usual human perception of writing on parchment and paper. As discussed in Section 3.3, the images represent elements of $\mathbb{R}^{5412\cdot7216}$. Since dark and white pixels feature gray levels close to and at a distance from zero respectively, the elements are currently characterized by the background pixels. Obviously, we are not interested in the background. Instead, we are interested in the writing. We therefore try to pre-process the images such that their equivalent elements in $\mathbb{R}^{5412\cdot7216}$ are characterized by the writing instead of the background.

³The computational accuracy of the processing hardware equals 2.2204×10^{-16} . Since $\sigma_{\mathbf{p}_3}^2$ is far less than 2.2204×10^{-16} , we consider the variances of principal component three through 16 as numerically zero.



Figure 3.8: Results of application of PCA to the toy example presented in Figure 3.7. To save space, only the first five principal components are presented. The pseudo color images depicted in Subfigures (b) through (f) were generated as described in Subsection 3.3.6.

Negating the images, we achieve the desired characterization. As negation alters the images' norms, we negate the images *before* normalizing them. In conclusion, our pre-processing pipeline from Subsection 3.2.4 is adapted as follows:

Pre-processing:

- (1) Histogram equalization
- (2) Gray level range maximization
- (3) Negation
- (4) Normalization

Figure 3.9 shows the pre-processed images. To emphasize the images' content, the images are not normalized. Apart from the inverted gray levels, Figure 3.9a does not differ from the earlier obtained Figure 3.5e. In other words, the negation of the images does not disturb the effects of the remaining pre-processing steps.

Figure 3.10 shows the principal components obtained by application of PCA to the preprocessed images. Compared to the images shown in Figure 3.3, the mean shown in Figure 3.10a improves the legibility of the depicted writing. Nonetheless, parts of the text remain illegible. In contrast, the first principal component shown in Figure 3.10b reveals nearly all the applied



Figure 3.9: Pre-processed multispectral images of the parchment introduced in Figure 3.3. Again, to emphasize the images' content, the images are not normalized. Continued on Page 46.

CHAPTER 3. ILLEGIBLE MANUSCRIPTS AND PCA



(j) White, Filter none



(m) White, Filter UVP

the second principal component (cf. Figure 3.10c).





(l) White, Filter Y22

Figure 3.9 continued: Preprocessed multispectral images of the parchment introduced in Figure 3.3. Again, to emphasize the images' content, the images are not normalized. Continuation from Page 45.

writing. Only a few letters contained in the text in the upper left corner remain illegible. The missing letters can be deduced from the enclosing context as well as the equivalent portions of

Furthermore, some principal components separate different colors of ink. For example, the fourth component separates the capital letters "Q" and "C" from the surrounding text (cf. Figure 3.10e).

Principal component 13 shown in Figure 3.10n appears all black. Manual examination of the component reveals non-zero elements. These elements are distributed sparsely over the whole component. Due to the high number of elements of the component, the infrequent non-zero elements are not recognizable for the human eye. In addition, principal component 13 features a numerically vanishing variance of 4.075596×10^{-30} . In the toy example discussed throughout Section 3.3, the internal structure of the examined data gave the reason why some of the encountered components feature vanishing variances. The structure depicted by the multispectral images examined throughout this section is unlikely to feature only twelve dimensions of variance. Thus, the conclusions drawn in Section 3.3 do not apply here.

Nonetheless, we can give a theoretical reason why only the first twelve principal components feature positive variance. We recall that the present multispectral series features 13 images. Hence, the matrix \mathbf{X} 's rank $r_{\mathbf{X}}$ is less than or equal to 13. The per dimension mean subtraction during PCA decrements $r_{\mathbf{X}}$ by one. Consequently, \mathbf{X} 's ranks are bounded by 12. Due

to Theorem 1, the number of principal components featuring positive variances is bounded by 12 as well.

To put it figuratively, imagine three arbitrarily chosen elements in \mathbb{R}^3 . No matter how these elements are chosen, we are always able to find a single plane which passes through all three elements. Given a "starting point" in the fitting plane, we can reach any other point in the plane with only two dimensions of variance. In other words, the plane features two dimensions of variance. Trivially, the same fit is possible for only one or two elements in \mathbb{R}^3 . In contrast, a set of elements, which can not be passed through by a plane, consists of at least four elements. Similarly, a set of $n \in \mathbb{N}$ elements of $\mathbb{R}^{5412 \cdot 7216}$ with $n \leq 5412 \cdot 7216$ can always be passed through by a hyperplane of at most n - 1 dimensions. In the framework of PCA, the mean of the examined data and the principal components define above starting point and dimensions of variance respectively. Thus, the analysis of n images will always provide at most n - 1 principal components featuring positive variances.



Figure 3.10: Principal components obtained by application of PCA to images shown in Figure 3.9. Components 13ff feature vanishing variances. To save space, further principal components are not shown here. Continued on Page 49.



(f) Principal component 5, $\sigma^2 = 5.669588 \times 10^{-4}$



(i) Principal component 8, $\sigma^2 = 1.641\,541\times 10^{-4}$



(l) Principal component 11, $\sigma^2 = 9.369368 \times 10^{-5}$



(g) Principal component 6, $\sigma^2 = 2.806\,472\times 10^{-4}$



(j) Principal component 9, $\sigma^2 = 1.247\,888 \times 10^{-4}$



(m) Principal component 12, $\sigma^2 = 8.080\,487\times 10^{-5}$



(h) Principal component 7, $\sigma^2 = 2.346\,302\times 10^{-4}$



(k) Principal component 10, $\sigma^2 = 1.054\,668 \times 10^{-4}$



(n) Principal component 13, $\sigma^2 = 4.075\,596\times 10^{-30}$

Figure 3.10 continued: Principal components obtained by application of PCA to images shown in Figure 3.9. Components 13ff feature vanishing variances. To save space, further principal components are not shown here. Continuation from Page 49.

Chapter 4

Recovery of covered text and recursive principal component analysis

Let us recall the two principles of this thesis introduced on Page 26 in Chapter 3. This thesis aims at

- (1) a sound understanding of the used tools and
- (2) the comparison of the results obtained by our pipeline to some baseline.

While Chapter 3 targets (1), we have not yet addressed (2). To this end, the initial plan of this thesis was to closely analyze the accessible characteristics of the urbarium. Such characteristics include the type of materials present, the number of layers of parchment and paper, the number of layers of writing and so on. With this knowledge, the plan was to build a fake urbarium featuring approximately the same characteristics. For this fake urbarium, all contained writing would have been known a priori to the application of the pipeline developed throughout this thesis. As such, the contained writing would have served as baseline for the rating of the pipeline.

Unfortunately, materials of same characteristics as present within the urbarium are rare. In order to build a fake urbarium, the materials would have to be used in destructive manner. Especially the intended destructive use prevented the acquisition of suitable materials. These difficulties of acquisition eventually rendered the initial plan unfeasible.

The crucial type of "damage" present in the urbarium is writing covered by additional layers of paper. To permit the examination of this type of "damage" despite above difficulties, we take a different approach. Throughout this section, we set up two stacks of loose sheets of paper. Here, the sheet at the bottom features writing covered by an additional sheet of paper. Within these stacks, the sheets of papers remain undamaged. Thus, the use of materials similar to those of the urbarium is eased. Furthermore, the stacks can be easily composed and decomposed to allow the examination of the covered writing. This way, the usage of loosely set up stacks allows to pursue (2).

There are two major drawbacks of the approach taken:

(1) The stacks of paper do not feature any glue in between the sheets of paper. Glue might alter the crucial characteristics of the imaged object. Hence, the appearance of glued and not glued objects throughout a multispectral series might deviate. Therefore, the

4.1. THE RUNNING EXAMPLE: PAPER MANUSCRIPTS

results obtained with the paper stacks might deviate from those eventually obtained for the urbarium.

(2) The stacks of paper hinder the quantitative comparison of results obtained with the pipeline and baseline results. To obtain baseline results, the uncovered writing has to be imaged. To obtain a suitable test object, at least one additional layer of paper has to be introduced afterwards. This addition has to be done before acquiring a multispectral series of the resulting stack. The resulting physical interaction with the initial sheet of paper causes tiny movements of the writing relative to the imaging hardware. Thus, a registration as described in Section 2.4.2 between the images depicting the uncovered and the covered writing is inevitable. Again, such a registration is beyond the scope of this thesis. Consequently, the comparison is performed on a qualitative basis only.

The rest of this section is structured as follows. Section 4.1 gives a more detailed description of the documents used and composed to paper stacks. Section 4.2 examines a multispectral series depicting covered text visually. In Section 4.3, we apply the principal component analysis to the multispectral series depicting covered text. Based on the obtained results, Sections 4.4 through 4.6 are dedicated to improving the current analysis. Section 4.7 concludes this section with a description of the results obtained by the improved analysis.

4.1 The running example: paper manuscripts

To simulate the characteristics of the paper present within the urbarium, three paper documents provided by the parochial archives Kößlarn are used. The documents date back to the second half of the 18-th century and feature weights of approximately 80 g m^{-2} to 90 g m^{-2} . Figures 4.1 through 4.3 depict all three documents.

Document 1 consists of a bunch of sheets of paper. Only the last two of those sheets are of interest. Figure 4.1a depicts the verso of the first sheet. As can be seen, the verso features writing in most parts of the page. The recto of Sheet 2 appears completely blank (cf. Figure 4.1b). Depicted in Figure 4.1c, Sheet 2's verso features writing only in the lower right quarter. Apart from this writing, the verso appears completely blank. The writing consists of brown ink applied in clear, thick strokes. The sheets appear relatively thin. Thus, Sheet 1 verso's writing shines through Sheet 2 in Figure 4.1c. Similarly, Sheet 1 recto's writing shines through in Figure 4.1a. Both sheets feature two strong creases dividing the sheets into four.

Document 2 consists of a single sheet of paper. The sheet is folded into two halves, of which only the one depicted in Figure 4.2 is of interest. The half features several further, light creases. Both recto and verso feature writing of brown ink. Here, the writing appears as rather thin, clear strokes. Additionally, parts of the recto's text have been underlined using brown ink, purple fineliner and gray pencil (cf. Figure 4.2a). This time, the sheet is thick enough to not let writing shine through.

Document 3 consists of a single sheet of paper again. The sheet is folded into two halves, of which only the one depicted in Figure 4.3 is of interest. The half features two more strong creases, dividing the relevant part of the sheet into four. Similar to Document 1, the sheet's verso features little writing of brown ink in the lower right part (cf. Figure 4.3b). Below the ink-writing, some words written with gray pencil are present. The rest of the verso, as well as the recto, appear blank. The sheet features a watermark in the center of the page. This watermark is not visible in Figure 4.3

In the following, we reference specific pages of Documents 1 through 3 by strings following the format "D α S $\beta\gamma$ ". Here, $\alpha \in \{1, ..., 3\}, \beta \in \{1, 2\}$ and $\gamma \in \{v, r\}$ state the document, the sheet and the page - recto or verso - respectively.



Figure 4.1: Document 1 used to set up paper stacks throughout Chapter 4

4.2 Visual examination of multispectral series depicting covered text

With the documents described in Section 4.1 at hand, we set up stacks of paper. In a first step, let us visually examine the influence of covering writing with paper to the obtained multispectral series. We then apply PCA to the series in a second step.

To the end of visual examination, we image the closed Document 1 "from behind". Here, the versos face towards the camera. Our goal is to recover the text written on D1S1v. During the acquisition of the multispectral series, the writing is covered by D1S2. As depicted in Figure 4.1c, the text we try to reconstruct shines through the covering sheet of paper. We focus our analysis on the upper half of the document. Figure 4.4 highlights the considered part.

Figure 4.5 depicts the preprocessed multispectral series. Here, preprocessing is performed as described in Section 3.4. Again, the images in Figure 4.5 are not normalized. We remark two observations:

- (1) The smaller the wavelength of the exposing radiation, the more noise the images feature. Due to the appearance of the noise, we conjecture that it is caused by scatter of the exposing radiation with the minute fibers of the covering sheet of paper. Obviously, the scattering effect highly depends on the radiation's wavelength. The smaller the wavelength of the exposing radiation, the more scatter occurs. Figures 4.5a and 4.5c show the extreme cases in more detail. The image shown in Figure 4.5a was exposed using ultraviolet radiation. Due to the high amount of scatter, no writing is visible at all. In contrast, Figure 4.5c features negligible scatter. The exposing infrared radiation penetrates the covering sheet of paper and reveals the underlying writing. Our observation matches the findings reported by Smith (cf. [14]). As the underlying physics are beyond the scope of this thesis, we do not study the effect in greater detail.
- (2) In Figure 4.5b, the text appears of differing sharpness. The closer the letters get to the image's rim, the sharper they appear. In the center of the image, the letters appear blurry compared to the outermost letters. We conjecture that the differences in sharpness result from an interplay of the covering and the covered sheet of paper. Both sheets bend



Figure 4.2: Document 2 used to set up paper stacks throughout Chapter 4

towards and away from the camera. Here, the bending is caused by the creases mentioned in Section 4.1. To minimize bending, a weighting plate of glass was placed on top of the sheets during the imaging process. Despite this glass plate, the bending causes air pockets in between the sheets of paper. Due to these pockets, the covering sheet of paper develops an effect similar to the view through frosted glass. We conjecture that the pockets are mainly located in the center of the image. In areas where the covering sheet is pressed firmly onto the covered sheet of paper the covered writing appears sharp (cf. the rim of Figure 4.5b). Figure 4.6a depicts the equivalent case for frosted glass. In contrast, areas where the covering and covered sheet of paper enclose air feature blurry writing (cf. the center of Figure 4.5b). Again, Figure 4.6f depicts the equivalent case for frosted glass.

We re-encounter (1) and (2) several times throughout the rest of this thesis. Both (1) and (2) establish crucial limitations of the recovery of covered text by the means of multispectral imaging. Here, the limitations are caused by the underlying physics - not by our experimental setup. To summarize both limitations, let us restate them as rules of thumb:

(Limitation 1)	The greater the wavelength of an image's exposing radiation, the more information about covered writing is encoded in the image. Similarly, the greater the wavelength, the less noise caused by scatter occurs.
(Limitation 2)	The greater the size of air pockets enclosed in between covered writing and covering sheet of paper, the blurrier the writing appears in images.

Due to the high quality of the recovery of the writing in Figure 4.5c, there is no need to apply PCA to the multispectral series depicted in Figure 4.5. We therefore examine a second stack of paper.



Figure 4.3: Document 3 used to set up paper stacks throughout Chapter 4

4.3 Application of PCA

This second stack is formed by Documents 2 and 3. Document 3 is used to cover Document 2. Here, D2S1r and D3S1v face towards the camera. We focus on the parts highlighted in Figure 4.7. In the following, we refer this stack of paper as *Stack 2*.

To save space, we print neither the original nor the pre-processed multispectral series corresponding to Stack 2. Instead, we directly jump to the analysis of the results obtained by application of PCA.

The results are shown in Figure 4.8. By careful analysis of the results we notice:

- (1) The mean depicted in Figure 4.8a reveals some of the covered writing. This characteristic of the mean indicates that several images of the multispectral series encode writing as well.
- (2) A high number of the first principal components encode information on the covered writing. The parchment examined throughout Section 3.4 features a similar spread of information. Nonetheless, in case of the parchment, all the information we are aiming to obtain is condensed into the first few principal components. In contrast, every component depicted in Figure 4.8 focuses on different parts of Stack 2. While Components 1 and 2 focus mainly on the left half Component 7 yields details about the center and Components 4, 9 and 11 encode information about the right half of the stack (cf. Figures 4.8b, 4.8c, 4.8e, 4.8h, 4.8j and 4.8l). This separation of areas hinders the analysis of the principal components.
- (3) Not all principal components encode information about the covered writing. Especially Components 13 ff are likely to not encode any insights at all since their corresponding variances vanish. Furthermore, there is no general guarantee that all of the first components do encode the desired information as in Figure 4.8.
- (4) Different components encode information in different ways. While Components 1 and 2 feature yellowish writing on blueish background, Component 9 features blackish writing



Figure 4.4: Document 1 with the parts considered throughout Section 4.2 highlighted in yellow

on yellowish background (cf. Figures 4.8b, 4.8c and 4.8j). Further combinations like blueish writing on yellowish background are also possible.

We recall that we initially introduced PCA to overcome the spread of information. The eventual reason for applying PCA is to merge the desired information into a small number of principal components. Judging by this eventual reason, we have to admit that PCA does not meet its purpose in the current example.

Projections establish a common approach to merge the information encoded in the relevant principal components. Here, the original data is projected onto the affine space induced by the per dimension mean $\overline{\mathbf{x}}$ of the data and the relevant components. Let P denote this projection. Consider the matrix $\mathbf{Q} \in \text{Mat}(\mathbb{R}, m, n)$ featuring the relevant components as columns. Then P is given by

$$P: \mathbb{R}^m \to \mathbb{R}^m$$
$$\mathbf{x} \mapsto \overline{\mathbf{x}} + \mathbf{Q}\mathbf{Q}^T (\mathbf{x} - \overline{\mathbf{x}})$$

Throughout the work on this thesis, we experimented with above projection. Unfortunately, the obtained results featured very little enhancement of the information about the covered writing. However, the principal components enhance the legibility of the covered writing compared to the original data. Hence, we do not want to discard the components and PCA in general. Instead, we try to come up with a different approach to merge the information encoded in the relevant principal components.

4.4 Introduction of recursive principal component analysis

Throughout the work on this thesis, we realized that two subsequent applications of PCA can improve the obtained principal components. Here, put in simple words, a set of slightly altered components obtained by the first application is used as input to the second application of PCA. Before we elaborate on "slightly altering" principal components, let us focus on the overall idea of recursive application of PCA first.



Figure 4.5: The pre-processed multispectral series obtained with Document 1. The closed Document 1 is imaged "from behind" with the versos facing the camera. The parts of Document 1 depicted in this figure are highlighted in Figure 4.4.







Figure 4.7: Documents 2 and 3 with the parts examined throughout Section 4.3 highlighted in yellow

4.4.1 The basic idea of recursive principal component analysis

Applying PCA recursively is based on the idea that - after per dimension mean subtraction the previously obtained principal components might be easier to reduce in dimensionality than the initial input data. Thus, the second application of PCA might condense desired information more easily. We do not have to restrict ourselves to two applications of PCA. Instead, we can easily generalize the recursive application of PCA to arbitrarily many applications.

Figure 4.9a depicts a single application of PCA as a basic flow chart. Figure 4.9b outlines how a single application of PCA transforms into arbitrarily many recursive applications. Here, we separate the obtained means for further use. Hence, the recursive applications receive the previous principal components only. Definition 4.1 gives the recursive application of PCA a summarizing name.

Definition 4.1. Let us analyze some data using PCA. Let us further analyze the obtained principal components and their variances using a recursive application of PCA. Here, we do not focus on two recursive applications but allow additional applications. Then we call the repeated recur-







(i) Principal component 8, $\sigma^2 = 8.088448 \times 10^{-4}$



(l) Principal component 11, $\sigma^2 = 4.853352 \times 10^{-4}$



(J) Principal component 9, $\sigma^2 = 6.639\,230 \times 10^{-4}$



(m) Principal component 12, $\sigma^2 = 4.718\,354 \times 10^{-4}$



(k) Principal component 10, $\sigma^2 = 5.937110 \times 10^{-4}$

Figure 4.8 continued: Results of application of PCA to multispectral series depicting Stack 2. The mean and almost all components encode some information on the covered text. Again, components featuring vanishing variance are not depicted. Continuation from Page 58.

sive application of PCA recursive principal component analysis (RPCA). Figure 4.9b outlines the data flow within RPCA.

4.4.2 Interpreting recursive PCA

With the flow depicted in Figure 4.9b, let us analyze the results of the various applications of PCA. With the notation introduced in Section 3.3, the first application provides us with the per dimension mean $\overline{\mathbf{x}} \in \mathbb{R}^m$ of the input data $\mathbf{X} \in \text{Mat}(\mathbb{R}, m, n)$ alongside the principal components and their respective variances $\sigma_1^2, \ldots, \sigma_m^2$. To ease notation, we let $\mathbf{Q}_1 \in \text{Mat}(\mathbb{R}, m, m)$ denote the principal components arranged as matrix. In slightly misleading notation, we let further $\sigma_1^2 \in \mathbb{R}^m$ denote the variances arranged as element of \mathbb{R}^m . Here, the subscript index 1 refers to the first application of PCA. As stated in Section 3.3, we can interpret the principal components as deviations of the input data \mathbf{X} from its mean $\overline{\mathbf{x}}$. The variances encode the components' respective importance.

Analogously, the second application provides us with the mean $\overline{\mathbf{w}}_1$ of the initially obtained principal components. We further obtain a new set of principal components $\mathbf{Q}_2 \in \text{Mat}(\mathbb{R}, m, m)$ and corresponding variances $\sigma_2^2 \in \mathbb{R}^m$. This time, the $\overline{\mathbf{w}}_1$ represents the mean of initial deviations. Similarly, the components \mathbf{Q}_2 represent \mathbf{Q}_1 's deviations of their mean $\overline{\mathbf{w}}_1$. To ease further analysis, Definition 4.2 introduces a set of terms.

Definition 4.2. Let $\mathbf{X} \in Mat(\mathbb{R}, m, n)$ denote some input data to the recursive application of *PCA* as in Figure 4.9b. For $k \in \mathbb{N}, k \geq 1$ let further denote

- $\overline{\mathbf{x}} \in \mathbb{R}^m$ the per dimension mean of the input data,
- $\mathbf{Q}_0 := \mathbf{X}$,
- $\overline{\mathbf{w}}_0 := \overline{\mathbf{x}},$



Figure 4.9: Charts outlining the data flow within a single application of PCA, recursive PCA and recursive PCA with mid-processing

- $\mathbf{Q}_k \in \text{Mat}(\mathbb{R}, m, m)$ the principal components obtained by the k-th recursive application of *PCA* arranged as matrix (\mathbf{Q}_0 is not orthogonal in contrast to \mathbf{Q}_k),
- $\sigma_k^2 \in \mathbb{R}^m$ the respective variances arranged as element of \mathbb{R}^m and
- $\overline{\mathbf{w}}_k$ the per dimension mean of \mathbf{Q}_{k-1} .

Then we call $\overline{\mathbf{w}}_k$, the principal components encoded in \mathbf{Q}_k and the variances contained in σ_k^2 the k-th order mean, principal components or deviations and variances respectively.

In general with the notation as in Definition 4.2, the k-th application of PCA provides us with the mean of the principal components of order k - 1, the principal components of order k and the respective variances of order k. The mean $\overline{\mathbf{w}}_k$ can be interpreted as average deviation of order k. The principal components \mathbf{Q}_k can be interpreted as deviations of next higher order from this average. While such interpretation of the mean and the principal component is easy for small k's, the interpretation loses meaning for great k's.

This interpretation might remind the attentive reader of repeated differentiation of a mapping $f: \mathbb{R} \to \mathbb{R}$. Here, the derivatives correspond to principal components of increasing order. We will, in fact, come back to this analogy later.

4.4.3 Improving recursive PCA

Let us now come back to "slightly altering" a set of principal components upon recursive application of PCA. There are several issues with RPCA as depicted in Figure 4.9b:

(1) The termination of RPCA is not clear yet. Without a proper termination condition, we can not decide whether RPCA terminates at all.

- (2) Item (4) from Section 4.3 ails the calculation of means. In case of conflicting encoding of writing, the information about the writing cancels out in the principal components' mean.
- (3) The result obtained by RPCA is not clear yet. So far, RPCA merely provides us with several means of deviations of different orders.

To overcome issues (1) and (2), we introduce additional processing steps to RPCA. Definition 4.3 gives these processing steps a summarizing name.

Definition 4.3. Let us analyze some data using RPCA as in Definition 4.1. In contrast to Definition 4.1, we apply further processing steps in between the applications of PCA. Here, the sequence of processing steps remains unchanged throughout the whole analysis. We refer to the processing steps in between two applications of PCA as mid-processing. Throughout this thesis, mid-processing gets a set of principal components alongside their respective variances as input. The output of mid-processing is a possibly smaller set of elements in \mathbb{R}^m . Figure 4.9c outlines the data flow within this adapted version of RPCA. Unless otherwise stated, the term RPCA refers to recursive principal component analysis with interleaved mid-processing.

Throughout the rest of this section, we consider the k + 1-th recursive application of PCA during application of RPCA with mid-processing. To this end, $\mathbf{Q}_k \in \text{Mat}(\mathbb{R}, m, m)$ denotes the principal components of order k arranged in a matrix. Further, $\sigma_k^2 \in \mathbb{R}^m$ and $r_k \in \mathbb{N}$ denote the corresponding variances arranged as element of \mathbb{R}^m and the number of non-zero variances respectively. In other words, \mathbf{Q}_k and σ_k^2 serve as input to the k-th application of mid-processing. The mid-processing result is then used as input to above k + 1-th recursive application of PCA.

Let us first address issue (1). To this end, we recall a property of PCA from Section 3.4: given n_{k-1} trials of input, at most $n_{k-1} - 1$ principal components obtained by the k-th application of PCA feature strictly positive variances. Further components of vanishing variances are not needed to reconstruct the input data. In other words, during mid-processing we can safely discard principal components featuring vanishing variances. This way, n_k is bounded by

$$n_k < n_{k-1} < \dots < n_1 < n_0$$

Once all variances provided by an application of PCA vanish, we end the execution of RPCA. This way, we apply PCA at most n_0 times. To accelerate termination, we make use of item (3) from Section 4.3: we discard unneeded principal components during mid-processing by hand.

Definition 4.4. We refer to discarding principal components during mid-processing as described above as variance based pruning and manual pruning. Variance based and manual pruning can be expressed as

$$\mathbf{Q}_k \mathbf{\Delta}_k \mathbf{\Lambda}_k$$

Here,

$$\mathbf{\Delta}_k := (\delta_{ij})_{\substack{i=1\dots m\\j=1\dots r_k}}$$

represents variance based pruning¹. $\Lambda_k \in Mat(\mathbb{R}, r_k, n_k)$ denotes a matrix with exactly one element equal to 1 per column and at most one element equal to 1 per row. All remaining elements of Λ_k equal zero. Λ represents manual pruning.

Second, we address issue (2). To this end, we recall how we depict principal components in pseudo color images. Based on the sign of a component's element, the equivalent pixel is either

 $^{^{1}\}delta_{ij}$ denotes the Kronecker-delta.

colored yellow or blue. In the following, we consider two principal components of which one encodes desired information about the writing in blue on yellow background. The other component encodes information the other way round. In other words, the former component encodes the writing with "negative elements on positive background". The latter component encodes the writing with "positive elements on negative background". If we compute the components," mean directly, the information about the writing will cancel out. Thus, the mean will not yield any information about the writing. To avoid the cancellation, we negate one of the components. After negation, the information propagates into the components' mean. Here, we may not arbitrarily choose one component to negate. Instead, the encoding of writing in the pre-processed multispectral series determines which type of encoding we prefer. After pre-processing, the series' images encode writing using light pixels on dark background (cf. Section 3.4). In other words, writing corresponds to gray levels of great magnitude. To emphasize the writing, we thus force the components to encode writing using "positive elements on negative background". To the end of negation, we introduce another manual mid-processing step. During this step, we negate individual principal components to match our preferred way of encoding writing.

Definition 4.5. We refer to negating individual principal components during mid-processing as described above as manual negation. Assuming no pruning, manual negation can be expressed as

$$\mathbf{Q}_k \mathbf{\Xi}_k$$

where

$$\Xi_k := (\xi_{ij})_{\substack{i=1...m\\ j=1...m}}^{i=1...m}$$

and
$$\xi_{ij} := \begin{cases} \pm 1 & i=j\\ 0 & otherwise \end{cases}$$

The multiplication in Definition 4.5 can be easily extended to matrices $\Phi_k \in Mat(\mathbb{R}, m, m)$ where Φ_k is given by

$$\Phi_k := (\varphi_{ij})_{\substack{i=1...m\\j=1...m}}$$

and
$$\varphi_{ij} := \begin{cases} x \in (0,1] & i=j\\ 0 & \text{otherwise} \end{cases}$$

If we further demand

$$\sum_{i=1}^{m} \varphi_{ii} = 1$$

we can even give a proper interpretation: the product $\mathbf{Q}_k \mathbf{\Phi}_k$ weights the components according to their relevance. Here, we decide on the relevance of a given component by hand based on the amount of writing it encodes. The mean calculation during the subsequent application of PCA becomes a weighted mean calculation.

Definition 4.6. We refer to the manual weighting of principal components as manual relevance weighting. With above notation, manual relevance weighting can be expressed as

$$\mathbf{Q}_k \mathbf{\Phi}_k$$

With above definitions, we can give our mid-processing pipeline:

	Mid-processing:
	(1) Variance based pruning
	(2) Manual pruning
	(3) Manual negation
	(4) Manual relevance weighting
	 (2) Manual pruning (3) Manual negation (4) Manual relevance weighting

Due to the manual processing steps during mid-processing, the results of RPCA are highly user-depending. To make our results reproducible, we outline our decisions taken during the manual processing steps in Appendix A.

4.4.4 Interpreting the result of RPCA and pseudo image construction

We still lack a proper interpretation of the results of RPCA. Before giving an interpretation, let us recall (3) from Subsection 4.4.3. The application of RPCA provides us a set of means of principal components of increasing order. The first order mean equals to the mean of the input data. The second order mean can be interpreted as average deviation of the input data from its mean. Further means can be interpreted as average deviations of higher order. As stated before, the increasing order of deviations reminds us of the increasing order of derivatives of a smooth mapping $f \colon \mathbb{R} \to \mathbb{R}$.

In the following, we describe a pseudo image construction based on the means obtained by RPCA. To motivate this construction, we state a well-known approximation of smooth functions:

Definition 4.7. Let $f : \mathbb{R} \to \mathbb{R}$ and $x_0 \in \mathbb{R}$ denote a smooth mapping and a fixed, real element respectively. Then T_{f,x_0}^K given by

$$T_{f,x_0}^K \colon \mathbb{R} \to \mathbb{R}$$
$$x \mapsto \sum_{k=0}^K \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k$$

is called K-th order Taylor polynomial of f at point x_0 . Here, k! and $f^{(k)}$ denote the factorial of k and the k-th derivative of f respectively.

One can show that T_{f,x_0}^K approximates f in some neighbourhood of x_0 . Identifying the derivatives $f^{(k)}$ with the means $\overline{\mathbf{w}}_k$, we aim at defining a pseudo image construction. To this end, let $\overline{\mathbf{w}}_k$ for $k \in \{0, \ldots, K\}$ denote the means of different order obtained by application of RPCA to some input data $\mathbf{X} \in \text{Mat}(\mathbb{R}, m, n_0)$. Consider the mapping

$$\tilde{F}_{\mathbf{X}}^{\text{Taylor}} \colon \mathbb{R} \to \mathbb{R}^m \\ x \mapsto \sum_{k=0}^K \frac{x^k}{k!} \overline{\mathbf{w}}_k$$

For reasonable choices of $x \in \mathbb{R}$, we find $\tilde{F}_{\mathbf{X}}^{\text{Taylor}}(x) \in [0,1]^m$. After delinearizing $\tilde{F}_{\mathbf{X}}^{\text{Taylor}}(x)$ into a matrix, we can thus interpret $\tilde{F}_{\mathbf{X}}^{\text{Taylor}}(x)$ as an image. In the general case, we can force the

mapping's image to be contained in $[0, 1]^m$ by capping its elements. To this end, we consider the mappings

$$\begin{array}{rccc} \max \colon & \mathbb{R}^m & \to & \mathbb{R}^m \\ & & (\mathbf{x}, \mathbf{y}) & \mapsto & (\max{(x_i, y_i)})_{i=1, \dots, m} \end{array}$$

and

$$\begin{array}{rccc} \min \colon & \mathbb{R}^m & \to & \mathbb{R}^m \\ & & (\mathbf{x}, \mathbf{y}) & \mapsto & (\min \left(x_i, y_i \right) \right)_{i=1, \dots, m} \end{array}$$

as well as the elements

$$\mathbf{1} := (1)_{i=1,...,m} \in [0,1]^m$$
$$\mathbf{0} := (0)_{i=1,...,m} \in [0,1]^m$$

In above setting, we can define our pseudo image construction.

Definition 4.8. Let $\overline{\mathbf{w}}_k$ for $k \in \{0, \ldots, K\}$ denote the means of different order obtained by application of RPCA to some input data $\mathbf{X} \in Mat(\mathbb{R}, m, n_0)$. Then we call

$$F_{\mathbf{X}}^{Taylor}: \quad \mathbb{R} \quad \to \quad [0,1]^{m}$$
$$x \quad \mapsto \quad \max\left(\mathbf{0}, \min\left(\mathbf{1}, \sum_{k=0}^{K} \frac{x^{k}}{k!} \overline{\mathbf{w}}_{k}\right)\right)$$

Taylor inspired pseudo image construction based on **X**. Due to the fact that the mid-processing introduced throughout this section depends highly on the user, there is a variety of constructions $F_{\mathbf{X}}^{Taylor}$ for the same input data **X**. For a given $x \in \mathbb{R}$, we call $F_{\mathbf{X}}^{Taylor}(x)$ the Taylor inspired pseudo image based on **X** at point x. Rearranging $F_{\mathbf{X}}^{Taylor}(x)$ into a matrix, we interpret $F_{\mathbf{X}}^{Taylor}(x)$ as sampled image.

The Taylor inspired pseudo image establishes a crucial difference between the type of results of PCA and RPCA. On the one hand, PCA provides us a mean alongside a set of principal components. Upon analysis of the results, we thus have to consider a set of components. As we have seen, the information we are aiming at may be spread over several components. The spread of information can render the analysis confusing. On the other hand, RPCA in conjunction with Taylor inspired pseudo image construction provides us a single image. Due to the various mid-processing steps, this single image contains all the desired, condensed information.

4.5 Application of RPCA in practice

4.5.1 The results of RPCA

Let us now examine the results of applying RPCA in practice. To this end, we apply RPCA to the multispectral series depicting Stack 2. Figure 4.10 shows the obtained means of increasing order. By definition of RPCA, Figure 4.10a depicts the multispectral series' mean. All depicted means feature some information about the covered writing. While the information is rather coarse in case of Figures 4.10a through 4.10d, the means of higher order focus on comparably fine details. Furthermore, the means of higher order focus mainly on specific parts of the overall image.

64



(a) Order 0 mean $\overline{\mathbf{w}}_0$



(d) Order 3 mean $\overline{\mathbf{w}}_3$



(g) Order 6 mean $\overline{\mathbf{w}}_6$



(j) Order 9 mean $\overline{\mathbf{w}}_9$



(b) Order 1 mean $\overline{\mathbf{w}}_1$



(e) Order 4 mean $\overline{\mathbf{w}}_4$



(h) Order 7 mean $\overline{\mathbf{w}}_7$



(k) Order 10 mean $\overline{\mathbf{w}}_{10}$



(c) Order 2 mean $\overline{\mathbf{w}}_2$



(f) Order 5 mean $\overline{\mathbf{w}}_5$



(i) Order 8 mean $\overline{\mathbf{w}}_8$

Figure 4.10: Stack 2's means of order 0 through 10 obtained by application of RPCA. Figures 4.10a and 4.8a equal by definition. All shown means feature some information about the covered writing. Again, blue and yellow pixels encode negative and positive deviations respectively. The brighter a pixel, the greater the deviation in magnitude.

65

4.5.2 Examination of Taylor inspired pseudo image construction

With the means from Subsection 4.5.1, we can now evaluate how the Taylor inspired pseudo image construction performs. As pointed out in Subsection 4.4.4, the images constructed using $F_{\mathbf{X}}^{\text{Taylor}}$ depend on the choice of $x \in \mathbb{R}$. Consequently, we examine the influence of the choice of x. To this end, the left column of Figure 4.11 depicts several pseudo images constructed using the means from Figure 4.10 for $x \in \{0, 0.125, 0.25, 0.5, 1, 2, 4, 8\}$.

We notice: the information about the covered writing increases for increasing x up to x = 1. Around x = 2, the contained information saturates. For greater x, the contained information rapidly degrades.

By definition, $F_{\mathbf{X}}^{\text{Taylor}}$ caps its image's elements of to the interval [0,1]. The right column of Figure 4.11 shows adapted versions of the pseudo images. These adapted versions highlight the capping performed upon pseudo image construction. Blue pixels encode negative elements capped to gray level 0. Analogously, yellow pixels encode elements greater than 1 capped to gray level 1. Not surprisingly, the amount of capping performed increases for increasing x. Unexpectedly, capping is performed much more often "at the lower gray level bound" than "at the upper". We conjecture that the degradation of information for x > 2 is caused by excessive capping.

There is one last issue with Taylor inspired pseudo image construction, especially with the images of maximal desired information $F_{\mathbf{X}}^{\text{Taylor}}(1)$ and $F_{\mathbf{X}}^{\text{Taylor}}(2)$. Although the high order means feature fine-grained details on the writing, the writing appears as rather thick strokes. Presumably, this appearance is caused by the decreasing weighting of $\overline{\mathbf{w}}_k$ for increasing k. By definition of $F_{\mathbf{X}}^{\text{Taylor}}$, the decreasing weighting is based on the factorials k!. To reduce the bold appearance of the writing, we equalize the weighting of the $\overline{\mathbf{w}}_k$ s. Here, we leave the weight of $\overline{\mathbf{w}}_0$ invariant. Furthermore, we leave the total weight of means $\overline{\mathbf{w}}_k$ with $k \geq 1$ invariant as well. We remark that our definition $F_{\mathbf{X}}^{\text{Taylor}}$ was motivated by the Taylor polynomial for smooth mappings $f: \mathbb{R} \to \mathbb{R}$. However, the definition of $F_{\mathbf{X}}^{\text{Taylor}}$ does not feature any mathematically sound derivation. While this lack of derivation establishes a gap to be filled on the one hand, it permits us to tweak the definition for our liking at the other hand.

Definition 4.9. Let $\overline{\mathbf{w}}_k$ for $k \in \{0, \ldots, K\}$ denote the means of different order obtained by application of RPCA to some input data $\mathbf{X} \in Mat(\mathbb{R}, m, n_0)$. Let further

$$a := \frac{\sum\limits_{k=1}^{K} \frac{1}{k!}}{K-1}$$

Then we call

$$\begin{aligned} F_{\mathbf{X}}^{Equal} : & \mathbb{R} & \to & \left[0,1\right]^m \\ & x & \mapsto & \max\left(\mathbf{0},\min\left(\mathbf{1},\overline{\mathbf{w}}_0 + a\sum_{k=1}^K x^k \overline{\mathbf{w}}_k\right)\right) \end{aligned}$$

equally weighting pseudo image construction based on \mathbf{X} . Further remarks as in Definition 4.8 apply analogously to this definition.

The right column of Figure 4.12 depicts the pseudo images constructed using $F_{\mathbf{X}}^{\text{Equal}}$. Again, the pseudo images are based on the means shown in Figure 4.10. Moreover, the amount of desired information increases for increasing x up to x = 1. For greater x, the weighting of high order means $\overline{\mathbf{w}}_k$ overshoots. Thus, the constructed pseudo images are of no use for $x \geq 2$. The overshooting becomes even more apparent when looking at the right column of Figure 4.12. This right column emphasizes the capping performed by $F_{\mathbf{X}}^{\text{Equal}}$ similar to the right column



Figure 4.11: Pseudo image $F_{\mathbf{X}}^{\text{Taylor}}(x)$ constructed for various $x \in \mathbb{R}$ using the means depicted in Figure 4.10. The left column depicts the constructed pseudo images. The right column highlights capping performed upon pseudo image construction. Blue and yellow pixels correspond to negative values capped to gray level 0 and values greater 1 capped to gray level 1 respectively. Continued on Page 68.



Figure 4.11 continued: Pseudo image $F_{\mathbf{X}}^{\text{Taylor}}(x)$ constructed for various $x \in \mathbb{R}$ using the means depicted in Figure 4.10. The left column depicts the constructed pseudo images. The right column highlights capping performed upon pseudo image construction. Blue and yellow pixels correspond to negative values capped to gray level 0 and values greater 1 capped to gray level 1 respectively. Continuation from Page 67.

of Figure 4.11. Again, the amount of capping "at the lower gray level bound" increases for increasing x. In contrast to $F_{\mathbf{X}}^{\text{Taylor}}$, we encounter an unarguable amount of capping performed "at the upper gray level bound" for $F_{\mathbf{X}}^{\text{Equal}}$.

Based on Figures 4.11 and 4.12, we fix x on x = 1 for later pseudo image constructions. Figures 4.11i and 4.12i both capture most of the information contained in the principal components depicted in Figure 4.8. Since Figure 4.12i tends to capture finer details compared to Figure 4.11i, we decide to go with $F_{\mathbf{X}}^{\text{Equal}}$.

4.6 Post-processing

At this point, we are almost ready to finally compare the results of our analysis to the baseline results. Before we cover the comparison, let us briefly analyze the histogram of the constructed pseudo image. The histogram is shown in Figure 4.13. The shown histogram shares some disadvantageous characteristics with the histograms of the initial raw images discussed in Subsection 3.2.4. Namely, the pseudo image features only few gray levels and the histogram is not leveled. To overcome these characteristics, we introduce additional processing steps. These processing steps are performed after application of main analysis. Definition 4.10 gives these steps a summarizing name.

Definition 4.10. When analyzing a set of images, we often encounter several transformations targeted at enhancing the raw results of the main analysis. All transformations performed after the application of the main analysis are thus called post-processing of the raw results. A specific sequence of post-processing transformations is referred to as post-processing pipeline.

Our first draft post-processing pipeline consists of a combination of tools introduced throughout Subsection 3.2.4. Namely, we apply gray level range maximization followed by an application of histogram equalization. Afterwards, we apply gray level range maximization again. Currently, the writing depicted by the pseudo image is encoded in "light pixels on dark background". To revert this encoding, we conclude the post-processing pipeline with a negation of the pseudo image. The post-processed counterpart of Figure 4.12i is depicted in Figure 4.14.

Clearly, the post-processed pseudo image features a lot of fine grained noise. Especially in parts depicting light background, the noise appears similar to so called *salt-and-pepper-noise*. Here, salt-and-pepper-noise refers to pitch black and pure white pixels spread randomly over the whole image. Salt-and-pepper-noise is usually reduced using a median filter. Before we apply this type of filter as well, let us briefly introduce it.

Definition 4.11. Let $f: \{0, \ldots, m-1\} \times \{0, \ldots, n-1\} \rightarrow \{0, \ldots, g\}$ denote a digital image. Given an odd $k \in \mathbb{N}$ and $z \in \mathbb{Z}$, we set

$$\begin{split} \mathbb{S}_{z} &:= \left\{ z - \frac{k-1}{2}, \dots, z + \frac{k-1}{2} \right\} \\ f_{\text{pad}} : \quad \mathbb{Z}^{2} \quad \rightarrow \quad \{0, \dots, g\} \\ (x, y) \quad \mapsto \quad \left\{ \begin{array}{cc} f(x, y) & (x, y) \in \{0, \dots, m-1\} \times \{0, \dots, n-1\} \\ 0 & otherwise \end{array} \right. \end{split}$$

as well as

$$\begin{array}{rcl} f_{\text{med}}: & \{0,\ldots,m-1\} \times \{0,\ldots,n-1\} & \to & \{0,\ldots,g\} \\ & (x,y) & \mapsto & \text{median} \left\{ f_{\text{pad}}\left(s,t\right) \, | \, (s,t) \in \mathbb{S}_x \times \mathbb{S}_y \right\} \end{array}$$

Then we call f_{med} the image obtained by median filtering of filter size k applied to f.



Figure 4.12: Pseudo image $F_{\mathbf{X}}^{\mathrm{Equal}}(x)$ constructed for various $x \in \mathbb{R}$ using the means depicted in Figure 4.10. The left column depicts the constructed pseudo images. The right column highlights capping performed upon pseudo image construction. Blue and yellow pixels correspond to negative values capped to gray level 0 and values greater 1 capped to gray level 1 respectively. Continued on Page 71.



Figure 4.12 continued: Pseudo image $F_{\mathbf{X}}^{\mathrm{Equal}}(x)$ constructed for various $x \in \mathbb{R}$ using the means depicted in Figure 4.10. The left column depicts the constructed pseudo images. The right column highlights capping performed upon pseudo image construction. Blue and yellow pixels correspond to negative values capped to gray level 0 and values greater 1 capped to gray level 1 respectively. Continuation from Page 71.



Figure 4.13: Histogram of the constructed pseudo image depicted in Figure 4.12i prior to gray level range maximization \mathbf{F}_{i}



Figure 4.14: Result of application of our first draft post-processing pipeline to Figure 4.12i

Due to the high resolution of the pseudo image, we have to fix k on a comparably high value. Throughout our work on this thesis, k = 19 has proven to be an adequate choice. With the tool of median filtering, we can give our final version of the post processing pipeline:

Post-processing:

- (1) Gray level range maximization
- (2) Median filtering with filter size k = 19
- (3) Histogram equalization
- (4) Gray level range maximization
- (5) Negation

Figure 4.15b shows the post-processed counterpart of Figure 4.12i. In contrast to Figure 4.14, Figure 4.15b depicts the result of above final post-processing pipeline.
4.7 Comparison to baseline results

Apart from depicting the post processed pseudo image, Figure 4.15 also depicts the baseline result. We obtained the baseline result by imaging the uncovered text. Here, white light exposed the stack. Comparing both images, we remark that the legibility of the recovered text falls behind the baseline result. It is possible to decipher the recovered text using knowledge from the baseline result. In contrast, reading the recovered text on its own is not possible in major parts of the image. Nonetheless, individual parts of the text are recovered clearly. Some characteristics which influence the recovery's legibility are:

- (1) The contours of the writing are recovered pretty well in Figure 4.15b. Based on their contour, especially letters exceeding the center part of a line are easily recognizable. Such letters are for example f, g, h, l, etc.
- (2) Broad strokes are recovered more clear than thin strokes. E.g. the recovery of the capital letter "U" in the upper left corner is of better quality than the recovery of the second line's center part.
- (3) Some areas feature better overall legibility compared to others. Figure 4.16 highlights areas of exceptionally clear and poor legibility in yellow and blue respectively. Again, we encounter blue areas mainly in the neighborhood of the center crease of the covering sheet of paper. We thus conjecture that this spatially local decrease in legibility is again caused by air enclosed in between the two layers of paper in the area of the crease. Furthermore, the covering sheet of paper may vary spatially in thickness. "Thin" and "thick" areas then correspond to areas of clear legibility and poor legibility respectively.

Compared to the principal components depicted in Figure 4.5, the writing is clearly enhanced in Figure 4.15.

To summarize, our approach is capable of recovering the writing clearly in areas of good physical preconditions. In areas of poor physical preconditions, further processing steps have to be taken. These further processing steps are beyond the scope of this thesis. Possible processing steps are outlined in Chapter 6. Since we can not know the conditions we find in case of the urbarium, we bravely apply our analysis as presented up to this point in Chapter 5.

We re-encounter (2) throughout the analysis of the urbarium. Thus, let us extend the list of limitations of the recovery of covered writing by the means of multispectral imaging. Additionally to (Limitation 1) and (Limitation 2) from Section 4.2 we state

(Limitation 3) The thinner the strokes of covered writing, the poorer the writing's recovery.

usto ! flyam, griplif. qual. ch 72- 04 go hallow in got you from & Corlefor artiner! Al Soflow for ha hi qualige ? chicks late mas any ato & fred Supelly Rinfind of guileta an Sifemp din gr -low y dif Korge fore filit yator fory ars Per prous la ityd "Jefro 2.00 oll uto day g · 9 · 14 H. 1. in fin wolf grille for sins preleftin guller origing For 17 yur diaguela , dergying hingro bogafter. So son Som all. fin quelle falloude falerefse Problichoud And fingigne per 1. lobisindary exotol farings fins breauft roler; china febru chil for goton lif the or fast lendra Coin interfining Dalaforo no Patoston bringt, in ofin Q reflective have so Withtand in tenpor guardig & forfi Hyp gor Jelow dalangon Ab Mofo es. (a) Baseline result



(b) Post processed pseudo image

Figure 4.15: Comparison of the baseline result and the pseudo image obtained by analyzing Stack 2's multispectral series $\frac{1}{2}$



Figure 4.16: Reprint of Figure 4.15b with a focus on legibility. Areas of exceptionally clear and poor legibility are highlighted in yellow and blue respectively.

Chapter 5

Analysis of the urbarium

As stated throughout the introduction on Pages 12f, this thesis aims at the analysis of an early modern urbarium. With the tool-set developed throughout Chapters 1 through 4, we are now ready to perform the analysis. To this end, Section 5.1 gives a detailed description of the urbarium itself. In Section 5.2, we analyze the multispectral series obtained by imaging the urbarium.

5.1 Description of the urbarium

The urbarium dates back to 1609. It originates from the city of Schärding, Austria. The urbarium shares the main characteristics of a modern book. The inner book is protected by a front and a back cover. The front cover's recto and the back cover's verso are depicted in Figure 5.1. The urbarium is of good general condition. Thus, the content of the inner book is easily accessible and readable.

Our analysis focuses on the urbarium's back cover. The cover features an exceptionally high number of layers. Three layers of paper are protected from the outside by a layer of parchment. All four layers are glued together. While the three layers of paper are relatively evenly glued, the bunch of paper and the parchment layer enclose noticeable air pockets. The used glue has rendered all layers quite rigid. Thus, the layers withstand even high pressure. Therefore, the air pockets can not be reduced in size by pressing the layers together without harming the back cover's physical integrity. Figure 5.2 depicts the cover's structure.

Figure 5.3 depicts a part of the back cover as transmitted light image. Here, the cover's recto faces towards the camera. As can be deduced from Figure 5.3b, the cover features three layers of writing. Two layers of writing stem from the parchment's recto and verso. In Figure 5.3b, these layers are highlighted blue and yellow respectively. We also notice the recto's writing shine through in Figure 5.1b. Visual examination indicates that the writing consists of iron-gall-ink.

As depicted in Figure 5.4, the back cover's recto is blank. Hence, we conjecture that the additional third layer of writing is located somewhere in between the layers of paper. Figure 5.1b highlights this layer of writing in red. We can only guess the type of ink of which this writing consists. We assume either iron-gall- or soot-ink based on the visual examination of Mr. Thuringer.

The writing on the parchment's verso is easily accessible from the outside. Due to the enclosing layers of paper and parchment, the parchment's recto's writing and the third layer of writing are illegible.



(a) The front cover's recto



(b) The back cover's verso

Figure 5.1: The closed urbarium imaged from in front and from behind. Apart from the title on the front cover, the writing in fact appears upside down.

5.2 Analysis of the urbarium

To recover the writing, we apply the analysis described throughout Chapters 1 through 4. In addition to the images mentioned in Section 2.4.4, we add another image to the multispectral series. This additional image is a transmitted light image as depicted in Figure 5.3. Throughout the series, the back cover's recto faces towards the camera.

The transmitted light image possibly contains additional information about the layer of writing we aim at. Since the transmitted light image depicts all three layers of writing, the additional image might cause unwanted artifacts in the analysis' results. As we will see, our analysis is mostly capable of suppressing artifacts resulting from transmitted light images.

We image the back cover's recto in six patches. This decomposition of the cover into patches allows us to capture finer details. The patches are numbered as shown in Figure 5.5.

5.2.1 Analysis of Patch 5

The examination of the obtained multispectral series is similar for all patches. Therefore, we focus our report on a single patch. We arbitrarily choose Patch 5. For the remaining patches, we merely state the results from Page 81 onwards.

Figure 5.6 depicts the pre-processed multispectral series obtained by imaging Patch 5. Several depicted images feature information about the writing we aim at.



Figure 5.2: Sketch of the back cover's structure. The cover consists of three layers of paper and one layer of parchment. The parchment features writing on both its recto (depicted in blue) and its verso (depicted in yellow). Furthermore, a layer of writing is enclosed somewhere between the layers of paper (depicted in red; only one of both items applies).

(1) The incident light images depicted in Figures 5.6f through 5.6i feature information. In terms of the exposing radiation's wavelength λ , incident light images feature information for roughly $\lambda \geq 600$ nm. This information is not always visible at first sight. Close examination and comparison of Figures 5.6f through 5.6h to Figure 5.6i reveals that especially the depicted bottom parts feature faint information. Again the rule of thumb from Section 4.2 applies: the greater the exposing radiation's wavelength, the more information the obtained image features.

However, not all incident light images feature information about the covered writing. Figures 5.6a through 5.6e do not depict the covered writing. Depending on the depicted object, images exposed with radiation of the same wavelength may differ with respect to the information about covered writing. For a fixed object, let λ_{min} denote the minimal wavelength which allows the depiction of covered writing. Our findings in Chapter 4 indicate that λ_{min} satisfies $365 \text{ nm} < \lambda_{min} < 450 \text{ nm}$ in case of Stack 2. In contrast, we find $\lambda_{min} \approx 590 \text{ nm}$ in case of the urbarium. We conjecture that λ_{min} depends on the thickness of the covering layer of paper. The thicker the layer, the greater λ_{min} gets. In other words, the thicker a layer of paper is, the greater the incident radiation's wavelength has to be in order to penetrate the paper.

- (2) We find desired information in the white light exposed, UVP filtered image depicted in Figure 5.6m. This finding is unexpected. Due to our discussion in Section 4.2, we do not expect desired information in images exposed with ultraviolet radiation. In conformity with Section 4.2, Figure 5.6a does not feature any desired information. In contrast, Figure 5.6m does feature desired information. The use of the UVP filter during the acquisition Figure 5.6m suggests of that the information is encoded by the reflected radiation of "short" wavelength. Comparison of Figure 5.6m to Figures 5.6j through 5.6l confirms that the information is encoded exclusively by the reflected radiation of "short" wavelength.
- (3) The transmitted light image depicted in Figure 5.6n features desired information. The transmitted light image highlights the writing we aim at especially in regions where the parchment's writing appears blurry. Furthermore, the transmitted light image reveals desired information located in the page's margin. Here, the page's margin refers to the black vertical bar in Figure 5.6n. This additional information is mostly not captured by the incident light images in Figure 5.6.

(1) establishes another crucial characteristic of the interplay of multispectral imaging and covered writing. Again, we extend our list of limitations of the recovery of covered writing by the means of multispectral imaging. For the sake of completeness, we restate previously stated limitations in this final version of our list.



(a) Patch 5 of the urbarium's back cover depicted as transmitted light image. The part highlighted in yellow is shown in more detail in Subfigure (b).



(b) A more detailed view of Patch 5 of the urbarium's back cover. The image depicts three layers of writing highlighted in yellow, blue and red. Yellow, blue and red writing correspond with writing on the parchment's verso, on the parchment's recto and in between the paper layers respectively. The coloring corresponds to the coloring in Figure 5.2.

Figure 5.3: Transmitted light images of the back cover



Figure 5.4: The urbarium's back cover opened. The left and the right parts of the image depict the inner book's verso and the back cover's recto respectively.

(Limitation 1)	The greater the wavelength of an image's exposing radiation, the more information about covered writing is encoded in the image. Similarly, the greater the wavelength, the less noise caused by scatter occurs.
(Limitation 2)	The greater the size of air pockets enclosed in between covered writing and covering sheet of paper, the blurrier the writing appears in images.
(Limitation 3)	The thinner the strokes of covered writing, the poorer the writing's recovery.
(Limitation 4)	The thicker a covering layer of paper, the greater the pen- etration threshold λ_{min} gets. Hence, we expect the covered writing's recovery quality to decrease for increasing thickness of the covering layer of paper.

Figure 5.7 depicts the means $\overline{\mathbf{w}}_k$ of orders $k = 0, \ldots, 2$. The means have been obtained by application of RPCA to the multispectral series shown in Figure 5.6. $\overline{\mathbf{w}}_0$ does not feature any information about the writing we aim at. Given the pre-processed multispectral series in Figure 5.6, this lack of information is not surprising. In contrast, $\overline{\mathbf{w}}_1$ and $\overline{\mathbf{w}}_2$ do feature information about the writing we aim at. In themselves, neither $\overline{\mathbf{w}}_1$ nor $\overline{\mathbf{w}}_2$ suffice to decipher



Figure 5.5: The decomposition of the back cover's recto into patches. The writing we aim at appears upside down relative to the urbarium's inner book. Thus, this image appears upside down relative to Figure 5.4. The depicted frames serve as guideline only. The actual patches feature significant overlap.

the writing. Positive deviations (in other words yellowish pixels) in $\overline{\mathbf{w}}_1$ and $\overline{\mathbf{w}}_2$ mostly correspond with information about the writing. Negative deviations (in other words blueish pixels) in $\overline{\mathbf{w}}_1$ and $\overline{\mathbf{w}}_2$ mostly correspond with the darkening of coarse regions of $\overline{\mathbf{w}}_0$. Here, the coarse regions mainly correspond to characteristics like stains or creases of the covering layer of paper.

mainly correspond to characteristics like stains or creases of the covering layer of paper. Figure 5.8 shows both $F_{\mathbf{X}}^{\text{Equal}}(1)$ and its post processed counter part. Here, $F_{\mathbf{X}}^{\text{Equal}}(1)$ is the result of applying the equally weighting pseudo image construction to the means depicted in Figure 5.7. $F_{\mathbf{X}}^{\text{Equal}}(1)$ fuses the information contained in the means $\overline{\mathbf{w}}_0$ through $\overline{\mathbf{w}}_2$. Hence, the writing depicted by $F_{\mathbf{X}}^{\text{Equal}}(1)$ features improved legibility compared to $\overline{\mathbf{w}}_1$ and $\overline{\mathbf{w}}_2$. Despite the use of an additional transmitted light image as input to our analysis, $F_{\mathbf{X}}^{\text{Equal}}(1)$ does not show any artifacts caused by the parchment's writing. We conclude that our analysis is capable of suppressing such artifacts. In other words, our analysis is capable of separating the layers of writing contained in Figure 5.6n. The post processed counterpart of $F_{\mathbf{X}}^{\text{Equal}}(1)$ improves the legibility further. Nonetheless, the legibility varies spatially as with the example discussed throughout Subsection 4.7. Compared to the result from Figure 4.15, $F_{\mathbf{X}}^{\text{Equal}}(1)$ features a lot more background noise and distortions. We conjecture that the noise is caused by the covering layer of paper. The urbarium consists of paper of much rougher characteristics than the paper used to set up Stack 2 from Subsection 4.3. Consequently, the results obtained for the urbarium feature more noise than the results obtained for Stack 2.

Throughout the analysis of Patch 5, we developed a workflow of two steps. In both steps, we apply RPCA and the aforementioned pseudo image construction. However, the result obtained by the first application does not yet establish the final result. Instead, it serves as a reference throughout the second application of RPCA. This reference eases the mid processing steps of manual pruning and manual relevance weighting.

5.2.2 Analysis of the overall results

Figure 5.9 depicts the post processed pseudo images $F_{\mathbf{X}}^{\text{Equal}}(1)$ obtained for all Patches 1 through 6. With respect to these pseudo images we remark:

• All patches share similar noise. This noise appears the same both in amount and structure for all patches.



Figure 5.6: Pre-processed multispectral series depicting Patch 5 of the back cover's recto. Again, the depicted images are not normalized. The depicted images feature roughly 1292 dpi. Continued on Page 83.

5.2. ANALYSIS OF THE URBARIUM



(m) White, Filter UVP

(n) White, Filter none, transmitted light image

- Only Patches 2 and 6 feature transmitted light artifacts. Generalizing from Patch 5, we restate that our analysis is robust against transmitted light artifacts in the average case. The artifacts are mostly suppressed. Again, our analysis separates the different layers of writing in the average case.
- The legibility of the recovered writing varies spatially. In general, the writing on the back cover's recto's inner half features better legibility than the outer half. Here, the inner and the outer half correspond to the right and the left column of Figure 5.9 respectively. There are at least two guesses on why the inner half is recovered better than the outer half.

First, the ratio ρ of "area depicting writing" to "total depicted area" differs a lot for the inner and the outer half. In case of Patches 5 and 6, both areas roughly equal. Thus, ρ is approximately 1. In contrast, in case of Patches 1 and 2, ρ equals roughly 0.5. In case of Patches 3 and 4, ρ drops to approximately 0.25. ρ can reduce the performance of histogram equalization during pre-processing if the image to equalize is exposed heterogeneous. Thus, ρ can reduce the quality of our analysis' results. To check this guess, we applied our analysis to a cropped version of the multispectral series depicting Patch 2. Here, cropping pushes ρ to reach a value of roughly 1 again. The quality of the analysis' result obtained with the cropped series did not exceed the original quality. In conclusion, the difference in quality

between the recto's inner and outer half is not caused by differing ρ . This conclusion is supported by the recovery's quality of Patches 3 and 4. As stated above, in case of Patches 3 and 4, ρ is even lower than in case of Patches 1 and 2. Nonetheless, the recovery's quality of Patches 3 and 4 exceeds the quality of Patches 1 and 2.

Second, the drop of the recovery's quality might be caused by the major crease depicted by Patches 1 through 3. The crease aligns vertically with the line indent of the text we aim at. As we have seen throughout Chapter 4, creases tend to cause air pockets. The air pockets in turn hinder the recovery of covered text.

- The legibility of the recovered writing varies depending on the stroke width of the writing. We have encountered this dependence already throughout Chapter 4. In case of the urbarium, the title depicted by Patches 1 and 6 consists of much broader strokes than the remaining text. Hence, the title is recovered of greater quality than the remaining text.
- The inexperienced viewer of the results depicted in Figure 5.9 might still wonder about the legibility of the recovered writing. Experienced paleographers can derive a first attempt of a transcription of the text depicted in Figure 5.9. To the end of a transcription, we contacted a former paleographer at the University of Passau. Unfortunately, as of the time of writing, the transcription remains work under progress. Hence, we can not state a transcription here. Instead, let us analyze the outline of the recovered text. The text consists of a title followed by three paragraphs of running text. The title is depicted by Patches 1 and 6. The three paragraphs are depicted by Patches 1 and 6 (first paragraph), 1, 2, 5 and 6 (second paragraph) and 2 through 5 (third paragraph) respectively.

According to Dr. Drost, the title says "An Wismader". This title indicates that the following running text refers to fiefs consisting of meadows and pasture. According to Dr. Drost, the recovered text might establish a blue print of a text contained in the urbarium's inner book. Within the inner book, there are five more passages of similar wording.



(a) $\overline{\mathbf{w}}_0$



(b) $\overline{\mathbf{w}}_1$



(c) $\overline{\mathbf{w}}_2$

Figure 5.7: Means $\overline{\mathbf{w}}_k$ of orders $k = 0, \ldots, 2$ obtained by application of RPCA to the multispectral series depicting Patch 5 of the urbarium. The multispectral series is depicted in Figure 5.6.



(a) $F_{\mathbf{X}}^{\mathrm{Equal}}\left(1\right)$

(b) Post processed $F_{\mathbf{X}}^{\mathrm{Equal}}(1)$

Figure 5.8: $F_{\mathbf{X}}^{\text{Equal}}$ and the post processed $F_{\mathbf{X}}^{\text{Equal}}$ obtained by equally weighting pseudo image construction using the means depicted in Figure 5.7



Figure 5.9: The final results obtained by application of the analysis described throughout this thesis to the urbarium. The patches feature significant overlap especially in vertical direction. The subcaptions of the subfigures are omitted for the sake of the overall view. Beginning in the upper left corner and enumerating counterclockwise, the subfigures depict Patches 1 through 6.

Chapter 6

Future work

Throughout the previous chapters, we have mentioned possible enhancements of our work. This chapter is dedicated to summarizing all these enhancements. Furthermore, we give possible approaches to achieve the enhancements. We group the enhancements into enhancements regarding image acquisition, mid processing and post processing in Sections 6.1 through 6.3 respectively.

6.1 Image acquisition

The amount of information about covered writing varies across the images contained in a multispectral series. From Section 4.2 we restate the following rule of thumb: the greater the exposing radiation's wavelength, the greater the amount of information about covered writing. In case of the urbarium, the wavelength has to be greater than roughly 600 nm. Wavelengths greater than this threshold permit the depiction of information about the writing we aim at. Letting aside white-light exposed incident and transmitted light images, the experimental setup described throughout Section 2.4 exposes only three out of nine images with radiation of sufficient wavelength. Further images exposed with radiation of sufficient wavelength are expected to improve the results. Thus, to improve the obtained results, we could acquire further exposing capabilities.

As stated in the motivation of this thesis, the analysis of the urbarium is restricted to hardware already present at the Chair of Digital Humanities. Obviously, the further acquisition of lighting panels contradicts this initial setting. Nonetheless, the acquisition is stated for the sake of completeness.

6.2 Mid-processing

The mid-processing described in Section 4.4.3 is heavily user-depending. Thus, the quality of the analysis' results is user-depending as well. This reduces the comparability of results. Imagine we apply the very same analysis to two types of objects. Or we apply the analysis to the same object multiple times with different pre- or post-processing pipelines. On the one hand, differences in the results' qualities can be caused by differences in the objects' characteristics or the pipelines' suitability. On the other hand, differences can also be caused by disadvantageous choices by the user throughout mid-processing. To overcome this incomparability, the mid-processing pipeline has to be fully automated. A possible approach to automation is to implement a scoring function for principal components. Here, the scoring function rates how much information about the writing we recover a given component depicts. Based on this rating, we can then automate

6.3. POST-PROCESSING

manual pruning, manual negation and manual relevance weighting. To the end of pruning, we discard components of low scores. To the end of negating a given component, we rate both the original component and its negated counter part. We then keep the version of greater score. To the end of weighting, we weight components according to their scores. A suitable scoring function could be derived by machine learning approaches.

Throughout mid-processing, we negate principal components based on the way they encode information. Unfortunately, some components feature conflicting encoding of writing. In other words, some parts of the component have to be negated while others have to remain unchanged. Figure 4.8e gives an example of such a component. The center part of the component encodes yellowish writing on blueish background. Hence, the center part has to remain unchanged. In contrast, the lower right corner depicts blackish writing on yellowish background. Consequently, the lower right corner has to be negated. So far, our analysis appears robust against few occurrences of conflicting encoding. Nonetheless, resolving conflicting encoding is expected to improve the analysis' results. To cope with conflicting encoding upon negation, we can subdivide components into independent parts. The parts then get negated independently from each other. This subdivision is expected to improve the analysis' results. To integrate this subdivision with above automation approach, we can ignore the actual need for subdivisions and divide every component per default. A specific part's scores then determine whether or not to negate the part.

6.3 Post-processing

In Subsections 4.4.4 and 4.5.2 we introduced the pseudo image constructions $F_{\mathbf{X}}^{\text{Taylor}}$ and $F_{\mathbf{X}}^{\text{Equal}}$. While we motivated both constructions, we did not give any mathematically sound derivation. To gain final confidence in the constructed pseudo images, we have to provide a sound derivation. A first approach to this derivation might be to weight the means $\overline{\mathbf{w}}_k$ according to the principal components' variances instead of the Taylor polynomials' weights. The approach's derivation might be easier for linear mid-processing steps.

As stated in Section 4.6, the constructed pseudo images feature fine grained salt-and-pepper noise. To reduce the noise, we apply median filtering during post-processing. Median filtering is a rather simple noise reduction tool. More sophisticated noise reduction tools are expected to improve the analysis' results. Throughout the work on this thesis, we experimented with noise reduction as described by Rudin et al.. Rudin et al. define a scoring function on images. This scoring function encodes the amount of noise in an image. The function further encodes an image's similarity to a fixed reference image. Minimization of the scoring function then yields a denoised counterpart of the reference image. Encoding even more properties of images in the scoring function permits us to further enhance our analysis' results. E.g. we could encode an image's blurriness to cope with the slight blur caused by covering layers of paper. In the average case, Rudin et al.'s noise reduction tool outperforms median filtering. The reason for us not to implement this tool within our post processing pipeline is above minimization. Our images resolution renders the minimization problem non-trivial to tune. Furthermore, the resolution renders the solving of the minimization problem time consuming. Close examination of the minimization problem could eventually make our pseudo images accessible for a noise reduction tool based on the findings by Rudin et al..

Conclusion

To conclude this thesis, we briefly restate the major findings over the course of our work. We further restate the main open endings of our work.

Major findings

First, we covered a set of tool-set related topics. We gave a brief introduction to the bedrock of digital image processing. Our introduction focused on the modeling and processing of monochromatic images. Moreover, we gave a brief introduction to the multispectral imaging of objects. We gave insights on how multispectral imaging distinguishes from ordinary photography. The use of a discrete spectrum of radiation permits multispectral imaging to emphasize the characteristics of an imaged object. This emphasis is not easily possible with ordinary photography. To analyze the obtained images, we introduced the framework of principal component analysis. Principal component analysis permits us to focus on relevant variations of the obtained multispectral series. We encountered a limitation of principal component analysis where the desired information spreads across a number of principal components. To overcome this spread of information, we introduced recursive principal component analysis. Recursive principal component analysis not only permits to focus on the relevant variations of the input data. Instead, the recursive analysis permits to focus on relevant variations of arbitrary order. To condense the means of increasing order obtained by recursive principal component analysis, we introduced the Taylor inspired pseudo image construction. We later transmuted this construction into the equally weighting pseudo image construction. These pseudo image constructions condense all the desired information contained in a single multispectral series into one image.

Second, we applied our tool-set of multispectral imaging and (recursive) principal component analysis to a number of use- and test-cases. By examination of ink-writing on parchment, we re-legitimated the combination of multispectral imaging and principal component analysis as tool to emphasize faint text. We further legitimated the combination of multispectral imaging and recursive principal component analysis as tool to recover covered writing. To this end, we examined two test-case documents. Both documents featured a priori known ideal recovery outcomes.

The test-case documents revealed four main limitations of multispectral imaging. First, only radiation of comparably great wavelengths penetrates paper. Hence, to depict covered writing, and image has to be exposed using radiation of comparably great wavelengths. Consequently, the imaging hardware should feature a number of distinct exposure configurations especially "in the reddish and near infrared wavelengths". Second, the penetration of the covering layer of paper depends on its thickness. The threshold the radiation's wavelength has to exceed to penetrate paper increases with increasing thickness. Third, the physical distance of the covering and covered layer of paper can hinder the recovery of covered text. Air pockets contained beneath An early modern urbarium established our third and final use-case. With respect to the urbarium, we were able to recover the title and the general outline of the covered text.

Open endings

The analysis we developed over the course of this thesis is heavily user-depending. To provide final comparability of the obtained results, the analysis has to be fully automated. I.e. the classification of principal components into the classes "depicts relevant information" and "does not depict relevant information" has to be automated. The construction of the concluding pseudo images has been motivated by Taylor polynomials. Still, the motivation lacks a mathematically sound derivation. To provide eventual confidence in the construction, a mathematically sound derivation has to be given.

To fully recover the covered text found within the discussed urbarium, further examination is inevitable. This further examination can be conducted using extended hardware capabilities. Alternatively, the examination can rely on a more sophisticated analysis of the obtained multispectral series.

References

- Michael Attas, Edward Cloutis, Catherine Collins, Douglas Goltz, Claudine Majzels, James R. Mansfield, and Henry H. Mantsch. Near-infrared spectroscopic imaging in art conservation: investigation of drawing constituents. *Journal of Cultural Heritage*, 4 (2):127-136, April 2003. ISSN 1296-2074. doi: 10.1016/S1296-2074(03)00024-4. URL http://www.sciencedirect.com/science/article/pii/S1296207403000244.
- [2] John K. Delaney, Mathieu Thoury, Jason G. Zeibel, Paola Ricciardi, Kathryn M. Morales, and Kathryn A. Dooley. Visible and infrared imaging spectroscopy of paintings and improved reflectography. *Heritage Science*, 4(1):6, December 2016. ISSN 2050-7445. doi: 10.1186/s40494-016-0075-4. URL https://link.springer.com/article/10.1186/ s40494-016-0075-4.
- [3] R. L. Easton, K. T. Knox, and W. A. Christens-Barry. Multispectral imaging of the Archimedes palimpsest. In 32nd Applied Imagery Pattern Recognition Workshop, 2003. Proceedings., pages 111–116, October 2003. doi: 10.1109/AIPR.2003.1284258.
- [4] Alejandro Giacometti, Alberto Campagnolo, Lindsay MacDonald, Simon Mahony, Stuart Robson, Tim Weyrich, Melissa Terras, and Adam Gibson. The value of critical destruction: Evaluating multispectral image processing methods for the analysis of primary historical texts. *Digital Scholarship in the Humanities*, 32(1):101–122, April 2017. ISSN 2055-7671. doi: 10.1093/llc/fqv036. URL https://academic.oup.com/dsh/article/32/1/101/ 2957366.
- [5] D. M. Goltz, E. Cloutis, L. Norman, and M. Attas. Enhancement of Faint Text Using Visible (420-720 nm) Multispectral Imaging. *Restaurator*, 28(1):11-28, 2008. doi: 10.1515/REST.2007.11. URL https://www.degruyter.com/view/j/rest.2007.28. issue-1/rest.2007.11/rest.2007.11.xml.
- [6] Rafael C. Gonzalez. Digital image processing. Pearson, New York, NY, fourth edition edition, 2018. ISBN 978-0-13-335672-4. URL http://infoguide.ub.uni-passau.de/ InfoGuideClient.upasis/start.do?Login=igupa&Query=540="978-0-13-335672-4".
- [7] F. Hollaus, M. Gau, and R. Sablatnig. Enhancement of Multispectral Images of Degraded Documents by Employing Spatial Information. In 2013 12th International Conference on Document Analysis and Recognition, pages 145–149, August 2013. doi: 10.1109/ICDAR. 2013.36.
- [8] Fabian Hollaus, Melanie Gau, and Robert Sablatnig. Multispectral Image Acquisition of Ancient Manuscripts. In Progress in Cultural Heritage Preservation, Lecture Notes in Computer Science, pages 30–39. Springer, Berlin, Heidelberg, October 2012. ISBN

978-3-642-34233-2 978-3-642-34234-9. doi: 10.1007/978-3-642-34234-9_4. URL https://link.springer.com/chapter/10.1007/978-3-642-34234-9_4.

- [9] Ian T. Jolliffe. Principal component analysis. Springer series in statistics. Springer, New York [u.a.], 2. ed. edition, 2002. ISBN 978-0-387-95442-4. URL http://infoguide.ub.uni-passau.de/InfoGuideClient.upasis/start.do? Login=igupa&Query=540="0-387-95442-2".
- [10] Thomas Müller-Gronbach. Einführung in die Stochastik Introduction to stochastics. Lecture held during the winter term 2015/2016 at the University of Passau., 2016.
- Simone Pentzien, Ira Rabin, Oliver Hahn, Jörg Krüger, Florian Kleber, Fabian Hollaus, Markus Diem, and Robert Sablatnig. Can Modern Technologies Defeat Nazi Censorship? In Computer Vision - ACCV 2012 Workshops, Lecture Notes in Computer Science, pages 13-24. Springer, Berlin, Heidelberg, November 2012. ISBN 978-3-642-37483-8 978-3-642-37484-5. doi: 10.1007/978-3-642-37484-5_2. URL https://link.springer.com/chapter/ 10.1007/978-3-642-37484-5_2.
- [12] Leonid I. Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259-268, November 1992. ISSN 0167-2789. doi: 10.1016/0167-2789(92)90242-F. URL http://www.sciencedirect. com/science/article/pii/016727899290242F.
- [13] Jonathon Shlens. A Tutorial on Principal Component Analysis. arXiv:1404.1100 [cs, stat], April 2014. URL http://arxiv.org/abs/1404.1100. arXiv: 1404.1100.
- [14] A. H. (Albert Hugh) Smith. The photography of manuscripts. London, 1938. (Reprinted from London medieval studies, v.1, pt.2).

Appendix A RPCA mid-processing decisions

The mid-processing applied throughout RPCA highly depends on the user's decisions (cf. Subsection 4.4.3). To make our results presented throughout Chapters 4 and 5 reproducible, we thus state the decisions taken upon the application of RPCA.

Upon the k+1-th recursive application of PCA, we arrange the principal components obtained by the previous application of PCA as matrix $\mathbf{Q}_k \in \text{Mat}(\mathbb{R}, m, m)$. Moreover, we let $\mathbf{W}_k \in$ Mat (\mathbb{R}, m, n_k) denote the mid-processed components which are used as input to the k + 1-th recursive application of PCA (cf. Figure 4.9). We finally let $\mathbf{0}_{m,n} \in \text{Mat}(\mathbb{R}, m, n)$ denote the matrix of all zeros. We now state matrices $\mathbf{\Theta}_k \in \text{Mat}(\mathbb{R}, m_k, n_k)$ such that

$$\mathbf{W}_k = \mathbf{Q}_k \cdot \left(egin{array}{c} \mathbf{\Theta}_k \ \mathbf{0}_{m-m_k,n_k} \end{array}
ight)$$

In other words, the matrix $\begin{pmatrix} \Theta_k \\ \mathbf{0}_{m-m_k,n_k} \end{pmatrix}$ equals to the product $\mathbf{\Delta}_k \mathbf{\Lambda}_k \mathbf{\Xi}_k \Phi_k$ with $\mathbf{\Delta}_k, \mathbf{\Lambda}_k, \mathbf{\Xi}_k$ and $\mathbf{\Phi}_k$ as in Subsection 4.4.3.

A.1 Stack 2 (cf. Chapter 4)

A.2 Patch 1 (cf. Chapter 5)

A.3 Patch 2 (cf. Chapter 5)

A.4 Patch 3 (cf. Chapter 5)

A.5 Patch 4 (cf. Chapter 5)

A.6 Patch 5 (cf. Chapter 5)

A.7 Patch 6 (cf. Chapter 5)

$$\boldsymbol{\Theta}_{2} = \frac{1}{8} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$
$$\boldsymbol{\Theta}_{3} = \frac{1}{3} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -1 \end{pmatrix}$$

Appendix B Digital appendix

In order to enable the interested reader of this thesis to reproduce the reported results, we provide a digital appendix. The appendix consists of several zip-archives. The archives contain both the multispectral series and the analysis' implementation discussed throughout this thesis. The series are given as sets of dng-files. The analysis' implementation is given as a set of Matlabsource-files. For details, see the file README.md distributed with the archives.

Appendix C Eidesstattliche Erklärung

Hiermit erkläre ich an Eides statt, dass ich die vorliegende Arbeit selbstständig verfasst und nur die angegebenen Quellen und Hilfsmittel verwendet habe.

Die Arbeit wurde bisher weder einer anderen Prüfungsbehörde vorgelegt noch anderweitig veröffentlicht.

Passau, den 26. März 2018

David Spitzenberg