

This is the draft website for *Programming Historian* lessons under review. Do not link to these pages.

For the published site, go to <https://programminghistorian.org>
(<https://programminghistorian.org>)

A Table of Contents

Contents

- [A Table of Contents](#)
- [Overview](#)
 - [Motivation](#)
 - [Scope](#)
 - [Objectives](#)
- [Prerequisites](#)
 - [Skills](#)
 - [Software](#)
 - [Windows 10](#)
 - [MacOS](#)
 - [Linux](#)
 - [Topic Modelling](#)
 - [Data](#)
- [Assessing your PDF\(s\)](#)
 - [Text recognition in PDF files](#)
 - [Extract embedded text from PDFs](#)
 - [Extract embedded images from PDFs](#)
 - [Combine images and PDFs to a single PDF](#)
 - [Use Topic Modelling to analyze the corpus](#)
 - [Download the corpus](#)
 - [Prepare and clean up the corpus](#)
 - [Create the Topic Model](#)
 - [Evaluate the Topic Model](#)
- [Concluding remarks](#)
 - [Alternatives](#)

Overview

Motivation

Humanities scholars often work with text-based historical and contemporary sources. In most cases, the Portable Document Format (PDF) (<https://en.wikipedia.org/wiki/PDF>) format is used as an exchange format. This includes digital reproductions of physical sources such as books and photographs as well as digitally created documents. The digitisation of these objects increases their accessibility and availability, but mostly also the quantity. Archives have begun to digitise entire collections and make them accessible via the Internet. Even more dramatic is the increase in the amount of data in digitally created sources such as corporate and government reporting. As a result, humanities scholars are increasingly exploring larger collections by means of Distant Reading and other algorithmic tools. However, PDF documents are only suitable for digital processing to a limited extent and must first be converted into plain text files.

¶ 1

Scope

If you meet one or more criteria, this lesson will be instructive for you:

¶ 2

- You work with text-based sources and need to extract the content of the sources.
- Your files are in PDF file format or can be converted to this file format.
- You work with a large corpus and you do not want to touch each file individually (Batch processing).
- You want to examine your corpus by the means of Distant Reading (<https://programminghistorian.org/en/lessons/?topic=distant-reading>) and therefore need it to be in plain text format.
- You don't have access to commercial software, such as Adobe Acrobat Professional or Abbyy FineReader.

Objectives

In more technical terms, in this lesson you will learn the following:

¶ 3

- Recognize and extract texts in PDFs with Optical Character Recognition (OCR) ([https://en.wikipedia.org/wiki/Optical character recognition](https://en.wikipedia.org/wiki/Optical_character_recognition)).
- Extract embedded texts from PDFs.
- Extract embedded images from PDFs.
- Combine images and PDFs to a single PDF file.
- Do all of the above at once (batch processing) with a large corpus.
- Analyze a large corpus using Topic Modelling (https://en.wikipedia.org/wiki/Topic_model) to get a quick overview of the topics it contains.

Tesseract OCR software ([https://en.wikipedia.org/wiki/Tesseract_\(software\)](https://en.wikipedia.org/wiki/Tesseract_(software))) used in this lesson supports over 110 languages including non-western languages and writing systems.

Prerequisites



Skills

You feel comfortable using the command line of your computer. Windows users should take a look at Introduction to the Windows Command Line with PowerShell (<https://programminghistorian.org/en/lessons/intro-to-powershell>). MacOS and Linux users should take a look at Introduction to the Bash Command Line (<https://programminghistorian.org/en/lessons/intro-to-bash>).

¶ 4

Software

Windows 10

Some components of the unix-based open source software used in this lesson do not run on Windows systems natively. Fortunately, since Windows 10 Fall Creators Update there is a workaround. Open PowerShell (<https://docs.microsoft.com/en-us/powershell/scripting/getting-started/starting-windows-powershell>) as administrator and run `Enable-WindowsOptionalFeature -Online -FeatureName Microsoft-Windows-Subsystem-Linux`. Install Ubuntu 18.04 LTS (<https://www.microsoft.com/store/apps/9N9TNGVNDL3Q>) from the Microsoft Store. To initialize (<https://docs.microsoft.com/en-us/windows/wsl/initialize-distro>) the Windows Subsystem for Linux (WSL) click on the Ubuntu tile in the Start Menu and create a user account.¹

¶ 5

Follow these instructions carefully and do not lose your credentials. You will need them as soon as you run programs as administrator.

Once the WSL is up and running, navigate to your working directory (Downloads). Invoke `bash` through PowerShell and install all requirements via the built-in package manager Aptitude (<https://wiki.debian.org/Aptitude>).

¶ 6

```
bash
```

```
sudo apt install ocrmypdf tesseract-ocr-all poppler-utils imagemagick
```

MacOS

Installing all the requirements without a package manager is cumbersome. Therefore install the [Command Line Tools for Xcode \(https://itunes.apple.com/us/app/xcode/id497799835\)](https://itunes.apple.com/us/app/xcode/id497799835) and [Homebrew \(https://brew.sh\)](https://brew.sh) first. It offers an easy way to install all the tools and software needed for this lesson. ¶ 7

```
xcode-select --install
/usr/bin/ruby -e "$(curl -fsSL
https://raw.githubusercontent.com/Homebrew/install/master/install)"
brew install ocrmypdf tesseract-lang poppler imagemagick
```

Linux


On [Ubuntu 18.04 LTS \(https://ubuntu.com/download/desktop\)](https://ubuntu.com/download/desktop) and most Debian-based Linux distributions you can install all requirements via `aptitude`. ¶ 8

```
apt install ocrmypdf tesseract-ocr-all poppler-utils imagemagick
```


Even though all tools used in this lesson are shipped with Ubuntu, an update is recommended. ¶ 9

Topic Modelling

The Topic Modelling in the case study is performed with the [DARIAH Topics Explorer \(https://dariah-de.github.io/TopicsExplorer/\)](https://dariah-de.github.io/TopicsExplorer/). It is a very easy to use tool with a graphical user interface. You can download the open source program for Windows, Mac and Linux [here \(https://github.com/DARIAH-DE/TopicsExplorer/releases\)](https://github.com/DARIAH-DE/TopicsExplorer/releases). ¶ 10

If you are on a  Mac and receive an error message that the file is from an “unidentified developer,” you can overwrite it by holding control while double-clicking it. If that doesn't work, go to Systems Preferences, click on Security & Privacy, and then click Open Anyway.

Data

Throughout this lesson you will work with historical documents from the  [1st International Conference of Labour Statisticians \(https://ilostat.ilo.org/resources/methods/icls/icls-documents/\)](https://ilostat.ilo.org/resources/methods/icls/icls-documents/) from 1923. The data of all past conferences is provided by the [International Labour Organization \(ILO\) \(https://www.ilo.org/global/about-the-ilo/history/lang--en/index.htm\)](https://www.ilo.org/global/about-the-ilo/history/lang--en/index.htm) and is publicly available (<https://www.ilo.org/public/libdoc/ilo/ILO-SR/>). ¶ 11

To make it easier for you to navigate through the file system and create folders, here are some basic commands of the Bash Command Line: ¶ 12

- Display the path of the current folder with `pwd`
- Display the contents of the current folder with `ls`
- Display only PDF files in the current folder with `ls *.pdf`
- Create a folder named `proghist` with `mkdir proghist`
- Change to this folder with `cd proghist`
- Open your current folder with a file browser `open .` (On Windows use `explorer.exe .`)
- Change to the parent folder with `cd ..`
- Change to your users home directory with `cd`
- Paste the code snippets into explainshell.com (<https://explainshell.com/>) to see what the code actually does

Throughout the lesson I will assume that ``proghist`` is your working directory.

Save all files below to your working directory:

¶ 13

- [Classification of industries](https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N1_engl.pdf) (https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N1_engl.pdf)
- [Statistics of wages and hours of labour](https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N2_engl.pdf) (https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N2_engl.pdf)
- [Statistics of industrial accidents](https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N3_engl.pdf) (https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N3_engl.pdf)
- [Report of the Conference](https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N4_engl.pdf) (https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO-SR_N4_engl.pdf)
- [International labour review](https://www.ilo.org/public/libdoc/ilo/P/09602/09602(1924-9-1)3-30.pdf) ([https://www.ilo.org/public/libdoc/ilo/P/09602/09602\(1924-9-1\)3-30.pdf](https://www.ilo.org/public/libdoc/ilo/P/09602/09602(1924-9-1)3-30.pdf))

To illustrate image extraction and PDF merging you will include one more files to our corpus that is not directly related to the International Conference of Labour Statisticians from 1923.

¶ 14

- [Speeches made at the ceremony on 21st October 1923](https://www.ilo.org/public/libdoc/ilo/1923/23B09_5_engl.pdf) (https://www.ilo.org/public/libdoc/ilo/1923/23B09_5_engl.pdf)

For the Topic Modelling of the [case study](#) you will download more files later in the lesson. ¶ 15

Always make a backup copy of your data before using the commands in this course. Text recognition and combining PDFs can change the original files.

Assessing your PDF(s)

In order to make this lesson as realistic as possible, you will be guided by a concrete historical case study. The study draws on the extensive collection of the [International Labour Organization \(ILO\)](https://ilostat.ilo.org/resources/methods/icls/icls-documents/) (<https://ilostat.ilo.org/resources/methods/icls/icls-documents/>), in particular the sources of the [1st International Conference of Labour Statisticians](#). ¶ 16

You are interested in what topics were discussed by the labour statisticians. For this purpose you would like to analyze all available documents of this conference using Topic Modelling. This assumes that all documents are available in plain text. ¶ 17

First you will get an overview of our corpus. Large databases can create a false impression of evidence. Therefore, the documents must be subjected to qualitative analysis. For this you will use scientific methods such as [source criticism](https://en.wikipedia.org/wiki/Source_criticism) (https://en.wikipedia.org/wiki/Source_criticism). All documents are written in English and are set in the same font. `ILO-SR_N1_engl.pdf`, `ILO-SR_N2_engl.pdf`, `ILO-SR_N3_engl.pdf` and `ILO-SR_N4_engl.pdf` are part of the same series. In addition, you note that the `ILO-SR_N2_engl.pdf` file does not contain any embedded text. You also note that `23B09_5_engl.pdf` contains images. One of these images contains text.² ¶ 18

1. You will recognize the text of `ILO-SR_N2_engl.pdf`.
2. You will extract the text from all PDF files.
3. You will extract images from `23B09_5_engl.pdf`.
4. [Previously](#) for illustrative purposes, you will combine different images and documents into a single PDF document. This can be helpful if the scanning process involves individual image files that are to be combined into a single document.
5. You will analyze a lot of plain text files using Topic Modelling.

Text recognition in PDF files

For the text recognition, you will use [OCRmyPDF](https://ocrmypdf.readthedocs.io) (<https://ocrmypdf.readthedocs.io>). This software is based on the state-of-the-art open source text recognition software [Tesseract](https://github.com/tesseract-ocr/tesseract) (<https://github.com/tesseract-ocr/tesseract>), which is maintained and further developed by Google. The software automatically recognizes the page orientation, corrects skewed pages, cleans up image artifacts, and adds an OCR text layer to the PDF. Only the document language must be given as a parameter. ¶ 19

```
ocrmypdf --language eng --deskew --clean 'ILO-SR_N2_engl.pdf' 'ILO-SR_N2_engl.pdf'
```

```
→ pdftest ocrmypdf --language eng --deskew --clean 'ILO-SR_N2_engl.pdf' 'ILO-SR_N2_engl.pdf'
INFO - 5: [tesseract] Image too small to scale!! (2x36 vs min width of 3)
INFO - 5: [tesseract] Line cannot be recognized!!
WARNING - 19: [tesseract] lots of diacritics - possibly poor OCR
WARNING - 23: [tesseract] lots of diacritics - possibly poor OCR
WARNING - 39: [tesseract] lots of diacritics - possibly poor OCR
WARNING - 51: [tesseract] lots of diacritics - possibly poor OCR
WARNING - Some input metadata could not be copied because it is not permitted in PDF/A. You may
wish to examine the output PDF's XMP metadata.
INFO - Optimize ratio: 1.00 savings: 0.0%
INFO - Output file is a PDF/A-2B (as expected)
```

[\(/ph-submissions/images/working-with-batches-of-pdf-files/working-with-batches-of-pdf-files1.png\)](#)

Figure 1: The status messages of the software indicate recognition errors in the OCR process. ¶ 20

The status messages of the software indicate recognition errors during the OCR process (see Figure 1). If certain errors occur systematically, it may be worthwhile to write a correction script. See [Cleaning OCR'd text with Regular Expressions](#)

(<https://programminghistorian.org/en/lessons/cleaning-ocrd-text-with-regular-expressions>) ¶ 21

OCRmyPDF has many useful parameters to optimize your results. See the [documentation \(https://ocrmypdf.readthedocs.io/en/latest/cookbook.html\)](https://ocrmypdf.readthedocs.io/en/latest/cookbook.html). The output might look slightly different depending on the version used.

To process all PDF files in your working directory at once. OCRmyPDF automatically skips PDFs that already contain embedded text. ¶ 22

```
find . -name '*.pdf' -exec ocrmypdf --language eng --deskew --clean '{}' '{} ' \;
```

Extract embedded text from PDFs

To extract the embedded texts from the PDF files, use [Poppler](#) ([https://en.wikipedia.org/wiki/Poppler_\(software\)](https://en.wikipedia.org/wiki/Poppler_(software))). It is a very powerful command line tool for processing PDF files that is used by many other programs. ¶ 23

```
pdftotext 'ILO-SR_N1_engl.pdf' 'ILO-SR_N1_engl.txt'
```

To process all PDF files in your working directory at once. The status message `Syntax Warning: Invalid Font Weight` means that the file contains formatting that does not meet the standard specifications of PDF. You can safely ignore this message. ¶ 24

```
find . -name '*.pdf' -exec pdftotext '{}' '{}.txt' \;
```

Once you have extracted all the embedded text from the PDFs, you can easily browse the text files. You can use the Windows Explorer, [macOS Finder](#), or a command line program like `grep`. You can display all the mentions of the term "statistics". ¶ 25

```
grep 'statistic' . -R
```

grep is also able to process complicated search queries (so-called regular expressions (<https://manpages.ubuntu.com/manpages/bionic/en/man1/grep.1.html#regular%20expressions>)). For example, you can also search for all files containing either "labour statistics" or "wage statistics". ¶ 26

```
grep -E 'labour statistics|wage statistics' . -R
```

Regular expressions also include numbers. This is particularly interesting for historians. This command displays all years in the 2nd century. ¶ 27

```
grep -E '19[0-9][0-9]' . -R
```

Once you have successfully extracted all text from the PDF files, they can be further analyzed using methods of Distant Reading (<https://programminghistorian.org/en/lessons/?topic=distant-reading>) such as Topic Modelling (<https://programminghistorian.org/en/lessons/topic-modeling-and-mallet>). You will apply such methods to the case study later in this lesson. ¶ 28

Extract embedded images from PDFs

PDF is a container file format and can contain multiple embedded images per page. You can also use Poppler to extract those images. The program allows us to select a target format for the extracted images. It is recommended use a lossless image format like PNG when working with the images. ¶ 29

```
pdftimages -png '23B09_5_engl.pdf' '23B09_5_engl'
```

For digitally created documents, Poppler extracts all embedded image files. This often includes image files that are outside the visible area or overlaid by other objects.

To process all PDF files in your working directory at once. ¶ 30

```
find . -name '*.pdf' -exec pdftimages -png '{}' '{} ' \;
```

Poppler can only extract illustrations if they are available as individual images in the PDF file. If you want to extract illustrations from a scanned page take a look at this lesson: Extracting Illustrated Pages from Digital Libraries with Python (<https://programminghistorian.org/en/lessons/extracting-illustrated-pages>). ¶ 31

Combine images and PDFs to a single PDF

Although OCRmyPDF can process image files directly, there are cases where you first want to combine the images into a PDF document. Because most image formats do not support multiple pages, each page of a document has to be saved as a single file. With the widespread command line image editing software [ImageMagick \(https://imagemagick.org/\)](https://imagemagick.org/) you can resize this very easily. ¶ 32

```
convert '23B09_5_engl-002.png' '23B09_5_engl-004.png' '23B09_5_engl-006.png'
'23B09_5_engl-007.png' 'some-images-combined.pdf'
```

To combine all images to a PDF file at once use the wildcard operator `*.png`. This step could take a few minutes. ¶ 33

```
convert '*.png' 'all-images-combined.pdf'
```

If you want to combine different PDF files, you can fall back on Poppler. The tool does this job much faster than ImageMagick and preserves attributes of the original documents. ¶ 34

```
pdfunite 'ILO-SR_N1_engl.pdf' 'ILO-SR_N2_engl.pdf' 'ILO-SR_N3_engl.pdf' 'ILO-
SR_N4_engl.pdf' 'some-pdfs-combined.pdf'
```

Use Topic Modelling to analyze the corpus

Now that you have performed all the steps of the PDF processing on some examples, you can return to the historical question of the case study. Which topics were discussed by the labour statisticians at the international conferences of the ILO? In order to answer this question using Topic Modelling, the following steps are necessary. ¶ 35

1. Download the corpus.
2. Prepare and clean up the corpus.
3. Create the Topic Model.
4. Evaluate the Topic Model.

Both the download and the processing of the corpus very time and resource consuming. At doi.org/10.5281/zenodo.3582736 (<https://zenodo.org/record/3582818/files/20191218-ilo-dataset.zip?download=1>) you can download the collection as a ZIP file and go directly to step 3.

Download the corpus

To avoid confusion create a new folder with `mkdir` and open it with `cd`. ¶ 36

```
mkdir case_study
cd case_study
```

You can download the corpus from the [ILO website](https://www.ilo.org/public/libdoc/ilo/ILO-SR/) (https://www.ilo.org/public/libdoc/ilo/ILO-SR/). All English documents contain 'engl' in the title. It's over a gigabyte of data. Depending on your internet speed this may take a while. ¶ 37

To automate this step you can alternatively use the following command line commands. This will download all English documents (340 files) at once. ¶ 38

```
curl https://www.ilo.org/public/libdoc/ilo/ILO-SR/ |
grep -o 'ILO[^"]*engl[^"><\/]*' |
uniq |
sed 's,ILO,https://www.ilo.org/public/libdoc/ilo/ILO-SR/ILO,g' >
list_of_files.txt
xargs -n 1 curl -O < list_of_files.txt
rm list_of_files.txt
```

Prepare and clean up the corpus

Now you can batch process all downloaded PDF files. First, perform text recognition on all files that don't have embedded text. Then extract all embedded text from the files. Depending on the performance of your computer, this step will take several hours. ¶ 39


```
find . -name '*.pdf' -exec ocrmypdf --language eng --deskew --clean '{}' '{} ' \;
&&
find . -name '*.pdf' -exec pdftotext '{}' '{}.txt' \;
```



Create the Topic Model

In order to create a Topic Model with the DARIAH Topics Explorer, you don't need to have any deeper mathematical knowledge about the used method Latent Dirichlet Allocation (LDA).³ Nevertheless, it is worth clarifying some implicit assumptions of the model before you begin.


- A corpus consists of documents. Each document consists of words. Words are carriers of meaning. The order (sentences, sections, etc.) of the words is completely irrelevant. Only the frequency of words in a document or corpus (or more precisely the co-occurrence of words) is measured.
- You determine how many topics are present in the corpus.
- Each word has a specific probability to belong to a topic. The algorithm finds the corresponding probabilities of the individual words.
- Words that occur very frequently do little to discriminate between the individual topics. They are often function words such as "and" and so forth. Therefore, they should not be included in the analysis.
- Topic modeling using LDA is non-deterministic. This means that a different result can be obtained for each run. Fortunately, the result usually converges towards a stable state. Run it several times and compare the results. You will quickly see if the topics remain

stable.

Now open the [DARIAH Topics Explorer](https://dariah-de.github.io/TopicsExplorer/) (https://dariah-de.github.io/TopicsExplorer/) and follow the steps given in the software  ¶ 41

1. Select all 340 text files for the analysis. Remove the 150 most common words. (Alternatively, you can also load the file with the English stop words contained in the [example Corpus](https://github.com/DARIAH-DE/TopicsExplorer/tree/master/data) (https://github.com/DARIAH-DE/TopicsExplorer/tree/master/data) of the [DARIAH Topics Explorer](#).) 
2. Choose 30 for the number of topics and 200 for the number of iterations. (You should play with the number of topics and choose a value between 10 and 100. With the number of iterations you increase the accuracy to the price of the calculation duration.)
3. Click on ["Train Model"](#) . Depending on the speed of your computer, this process may take several minutes.

Evaluate the Topic Model

The [DARIAH Topics Explorer](https://dariah-de.github.io/TopicsExplorer/) (https://dariah-de.github.io/TopicsExplorer/) has a graphical user interface that makes it very easy to explore and evaluate the Topic Model and its thirty topics. In this run the second topic  looks like this (see Figure 2). ¶ 42

pension, benefits, benefit, ...

On this page you can find the 15 most relevant words for this topic, as well as the 10 most relevant documents, whose bar width indicates the respective weight, and the three most similar topics, where the [cosine similarity](#) between all *topic vectors* was calculated and ranked.

TOP 15: RELATED WORDS

pension

benefits

benefit

fund

sickness

funds

contributions

pensions

TOP 10: RELATED DOCUMENTS

ILO-SR_NS23_ENGL.PDF

ILO-SR_M13_ENGL_V.2.PDF

ILO-SR_M11_ENGL.PDF

ILO-SR_M13_ENGL_V.1.PDF

ILO-SR_M10_ENGL.PDF



</ph-submissions/images/working-with-batches-of-pdf-files/working-with-batches-of-pdf-files2.png>

Figure 2: DARIAH Topics Explorer showing related words, related documents and similar topics of a single topic. ¶ 43

This topic deals with various social insurance schemes. Both old-age provision and unemployment benefits are included. The words are sorted in descending order of relevance and give a good overview of the topic. You can also see which documents have the highest correspondence with this topic. As you can see from a look at [related topics](#), this topic differs from the topics on accident insurance and legislation. ¶ 44

To further process or visualize the results with a spreadsheet program, click on the [Export Data](#) button. The paper ["Parliament's Debates about Infrastructure"](#) by Jo Guldi illustrates how Topic Modelling can be put to use for historical research.⁴ ¶ 45

Concluding remarks

Over the past decades, PDF has become the de facto standard for archiving and exchanging digital text documents.⁵ This development has also not stopped at the humanities and will accelerate even further in view of the countless digitisation and editing projects. ¶ 46

However, this is not only the case for projects that focus primarily on digitized historical sources. For many digitally generated content such as websites and interactive documents, no generally accepted archiving formats have yet been established. Therefore, PDF is often used in these cases as well. Sometimes contemporary source documents present us with the same challenges as inferior scans of historical documents. As a technical analysis by Duff Johnson shows, the Mull Report was digitally created, printed, scanned at least once, and in an inferior version sent for text recognition. Text and metadata were lost, which would have made working with the document much easier.⁶ ¶ 47

Alternatives

This lesson focused on tools that are easy to use and are available as open source software free of charge. There are a lot of open source and commercial alternatives (https://en.wikipedia.org/wiki/List_of_PDF_software) to process PDF files.⁷ Getting Started with Topic Modeling and MALLET (<https://programminghistorian.org/en/lessons/topic-modeling-and-mallet>) covers one of many alternatives for Topic Modelling. ¶ 48

1. If you run into troubles activating (<https://docs.microsoft.com/en-us/windows/wsl/install-win10>) the WSL check out the troubleshooting (<https://docs.microsoft.com/en-us/windows/wsl/troubleshooting>), documentation (<https://aka.ms/wsl/docs>), or the learning (<https://aka.ms/learnwsl>) resources. ↩ ¶ 49
2. In the case of a larger corpus, it is advisable to carry out random sampling instead of a detailed analysis. If no text is embedded in certain files, text recognition can be run over the entire corpus. Text recognition recognizes embedded text and performs text recognition only when one is missing. ↩ ¶ 50
3. If you still want to learn more, see Ganegedara, Thushan. "Intuitive Guide to Latent Dirichlet Allocation." Medium August 23, 2018. <https://towardsdatascience.com/light-on-math-machine-learning-intuitive-guide-to-latent-dirichlet-allocation-437c81220158>. ↩ ¶ 51
4. Guldi, Jo. "Parliament's Debates about Infrastructure: An Exercise in Using Dynamic Topic Models to Synthesize Historical Change." *Technology and Culture* 60, no. 1 (2019): 1–33. <https://doi.org/10.1353/tech.2019.0000>. ↩ ¶ 52

5. The fourth chapter describes in detail the history of the PDF file format and the associated socio-technological upheaval, see Gitelman, Lisa. [Paper Knowledge: Toward a Media History of Documents](#). Durham: Duke University Press, 2014. ↩ ¶ 53
6. Johnson, Duff. "A Technical and Cultural Assessment of the Mueller Report PDF." *PDF Association* (blog), April 19, 2019. <https://www.pdfa.org/a-technical-and-cultural-assessment-of-the-mueller-report-pdf/>. ↩ ¶ 54
7. Especially worth mentioning are the German wiki pages of the Ubuntu community about [PDF](https://wiki.ubuntuusers.de/PDF/) (<https://wiki.ubuntuusers.de/PDF/>) and [OCR](https://wiki.ubuntuusers.de/Texterkennung/) (<https://wiki.ubuntuusers.de/Texterkennung/>). These pages contain references to free software for working with PDF files and improving text recognition. Unfortunately, there are no translations into other languages for these pages, so a translation service should be used. ↩ ¶ 55

license

[Hosted on GitHub](https://github.com/programminghistorian/jekyll) (<https://github.com/programminghistorian/jekyll>)

[Site last updated 21 January 2020](#)

<https://github.com/programminghistorian/jekyll/commits/gh-pages>

[See page history](https://github.com/programminghistorian/jekyll/commits/gh-pages/lessons/working-with-batches-of-pdf-files.md) (<https://github.com/programminghistorian/jekyll/commits/gh-pages/lessons/working-with-batches-of-pdf-files.md>)

[Make a suggestion](http://programminghistorian.org/feedback) ([/ph-submissionshttp://programminghistorian.org/feedback](http://programminghistorian.org/feedback))