# MET CS688
# *Web Analytics and Mining*

## Zlatko Vasilkoski

TEXT MINING EXERCISES

# *Exercise*: Extracting the text content from a pdf file with Xpdf

Implement the following steps:

1. Install Xpdf Tools.

2. Test the installation on your OS by running "pdftotext.exe" in a terminal mode.

3. Run the terminal mode command "pdftotext.exe" from R using R's *system()* function.

4. Use R to extract content from several pdf files in a folder by running the terminal mode command "pdftotext.exe" in a loop using *lapply()*.

# 1. Extracting the content from a pdf file

- Installing Xpdf - Extracts the content from a pdf file

- You might find it helpful to use it in Text mining projects on your own.

- Install pdftotext.exe (open-source PDF viewer) part of the Xpdf software suite.

- It can be downloaded from: http://www.xpdfreader.com/download.html

- Choose the download for your operating system.

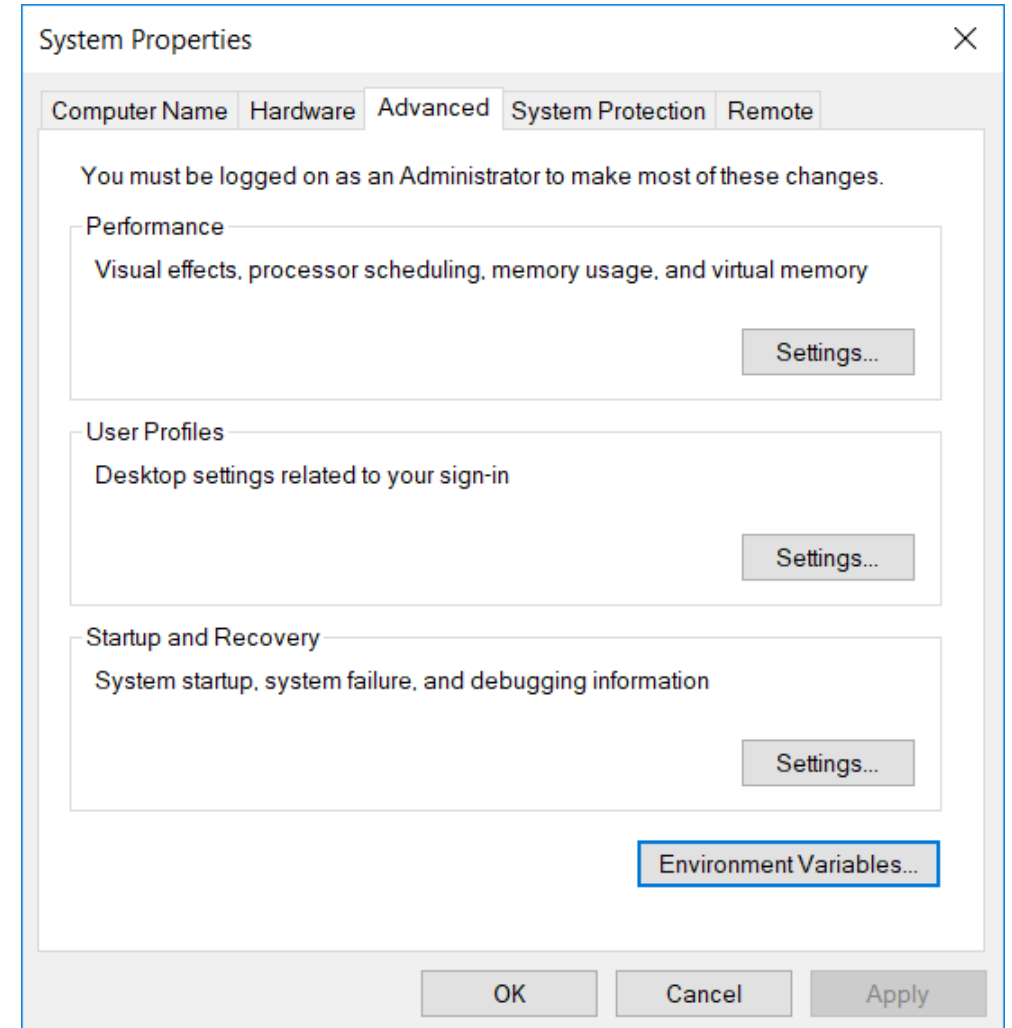- The Windows installation is more involved so it is illustrated on the next few slides.

# Additional Xpdf Utilities

These are all of the Xpdf utilities:

1. pdftotext -- generates a text file from a pdf file
2. pdftops -- generates a PostScript file from a pdf file
3. pdfinfo -- dumps a PDF file's Info dictionary (plus some other useful information)
4. pdffonts -- lists the fonts used in a PDF file along with various information for each font
5. pdfdetach -- lists or extracts embedded files (attachments) from a PDF (archived) file
6. pdftoppm -- converts a PDF file to a series of PPM/PGM/PBM-format bitmaps
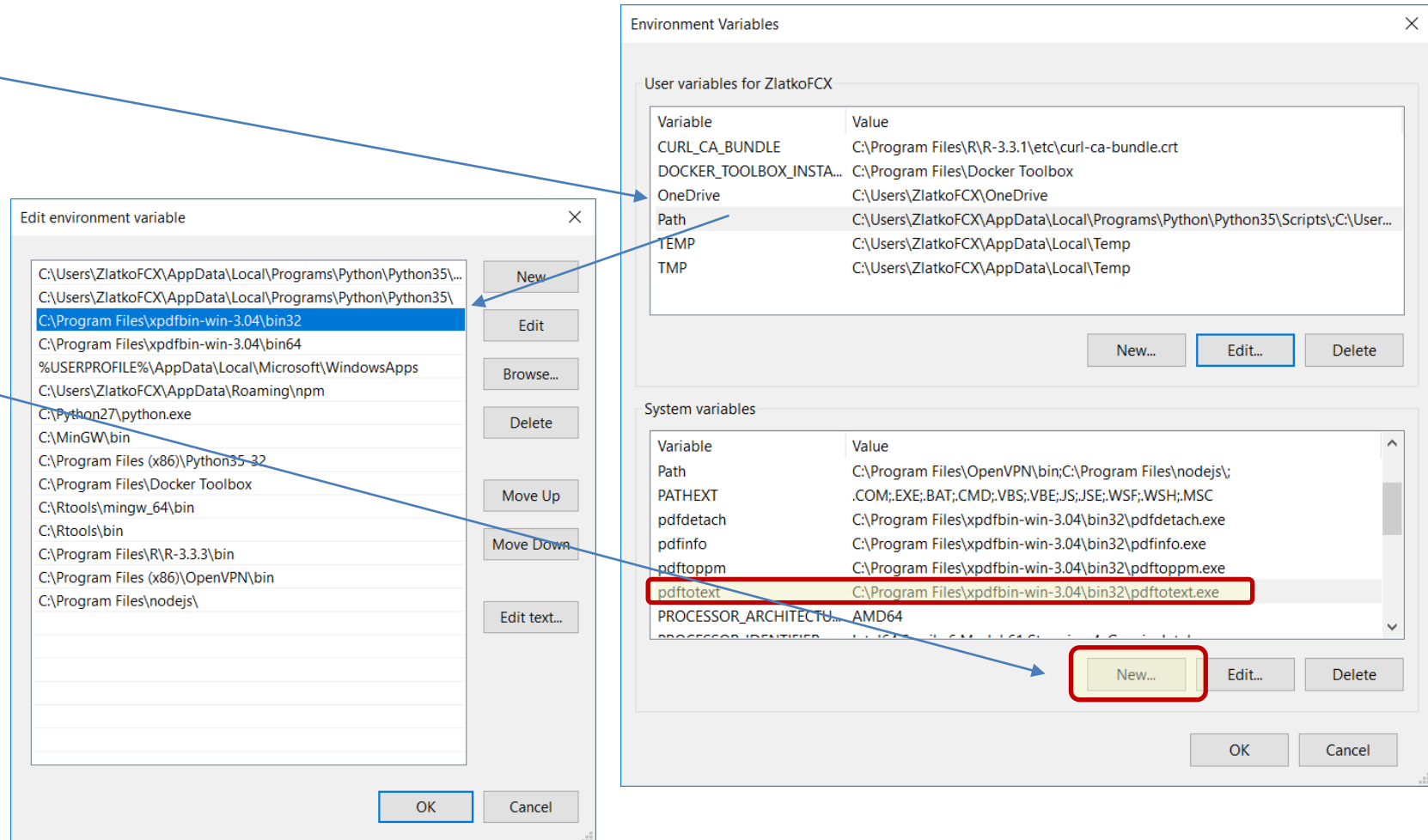7. pdfimages -- extracts the images from a PDF file

# Installing Xpdf

- For mac users follow:
  [http://macappstore.org/pdftotext/](http://macappstore.org/pdftotext/)

- Un-compress the downloaded files in a folder and set the environmental variables and the path to point to the folder where downloaded files are so R can find them.

- On the Windows operating system you would do that by going to **"System"**, choosing **"Advanced system settings"**, and **"Environment Variables"**.

# Installing Xpdf

- Add to the "Path" user-variable the **folder** where the downloaded Xpdf files are.

- By clicking "New" add to the System the 2 new environmental variables with "Variable name" "pdfinfo" and "pdftotext".

- For the "Variable value" enter the **path** to the executables "pdfinfo.exe" and "pdftotext.exe".

# Install Xpdf on Mac

- Download the Mac source code to your computer (Desktop, etc.).
- Uncompressed source code will expand into a folder.
- Open up a Terminal and  follow the directions in the "INSTALL" text file (shown below):
  - To install this binary package:
  - 1. Copy the executables (xpdf, pdftotext, etc.) to to /usr/local/bin.

  - 2. Copy the man pages (*.1 and *.5) to /usr/local/man/man1 and /usr/local/man/man5.

  - 3. Copy the sample-xpdfrc file to /usr/local/etc/xpdfrc.  You'll probably want to edit its contents (as distributed, everything is commented out) -- see xpdfrc(5) for details.

- To test the installation:
- In the Terminal using "cd /location of file", navigate to the directory where the PDF file is and then type:
  - pdftotext -layout pdfname.pdf
- Depending on the size of the PDF file, your output text file (with the same name as the original) will be in the same folder in a matter of seconds.

# 2. Test the "pdftotext.exe" installation

- To test the installation, convert a pdf file in a specific folder to a text file.

- The "pdftotext.exe" conversion from pdf to text in a terminal mode is implemented at as:

  > pdftotext "file.pdf"

- Note:
  - This usage produces a text file with the same name as the input file
  - The text file is created in the same directory as the PDFs.

# Test the "pdftotext.exe" installation

- Open a terminal window (cmd on Windows).

- Either navigate to the directory where the PDF file is or include the path to it in the filename.

- List all the files in the folder using "dir"

- Type: pdftotext "Class 1.pdf"

- This produces a text file, with the same name as the pdf file, created in the same directory as the PDFs.

- List all the files in the folder again using "dir" to see it.

- Note that the file name (and the path) needs to be enclosed in quotation.

- In R we can use the function *system()* to implement a cmd command as you would implement in a terminal window.

# 3. Try implementing these R code examples

Download several pdf class notes from Blackboard and place them in a folder "PDF Files", then

**Task1:** Use R to get the text content by implementing terminal mode command such as

> > pdftotext "file.pdf"

- To accomplish this from an R script you can use R's function *system()*
- To insert the **""** around the pdf filename you need to escape them with **''** . To merge, use *paste0().* Here is the R code example:

  system(paste(Sys.which("pdftotext"), paste0('"', myPDFfiles[1], '"')), wait=FALSE)

**Task2:** How to extend this to multiple files in a folder?

- A wildcards (*), for example *pdftotext "*pdf",* for converting multiple files, cannot be used because pdftotext expects only one file name.

- Using R's *lapply()* several pdf files that are contained in a single folder (specified by the R object "myPDFfiles") can be converted with an in line function such as this,

  > lapply(myPDFfiles, function(i) system(paste(Sys.which("pdftotext") , paste0('"', i, '"')), wait = FALSE))

- Note how each PDF file converted into a text file is indexed by *"i"*
- Note: Quotes ("") in R are tricky. Make sure you properly "escape" them with single quotes such as in paste0('"', i, '"')) . **Copy/pasting the above code in R may not work. You need to type it!**

# 4. Getting several pdf files from a folder

**Task:** How to get path to **several** pdf files that are contained in a **single** folder.

```
# Example 1: Convert to text single pdf files that is contained in a single folder.
exe.loc <- Sys.which("pdftotext")  # location of "pdftotext.exe"
pdf.loc=file.path(getwd(),"PDF Files")  # folder "PDF Files" with PDFs

# Get the path (character vector) of PDF file names
myPDFfiles <- normalizePath(list.files(path = pdf.loc, pattern = "pdf",  full.names = TRUE))

# Convert single pdf file to text by placing "" around the character vector of PDF file name
system(paste(exe.loc, paste0('"', myPDFfiles[1], '"')), wait=FALSE)
```

Note:

- *Sys.which* gives the path to *pdftotext.exe* (the one you set it up during installation).
- *file.path(getwd(),"PDF Files")* Gets the current folder (*getwd()*) and forms a path.
- *normalizePath()* Converts the file paths to a canonical form for the operating system.
- *list.files()* Lists the files in a Directory/Folder.
- *system()* Invoke a system command.
- myPDFfiles[1] Access the first element of vector "myPDFfiles" (contains path to the PDF file).

# Using the "pdftools" package

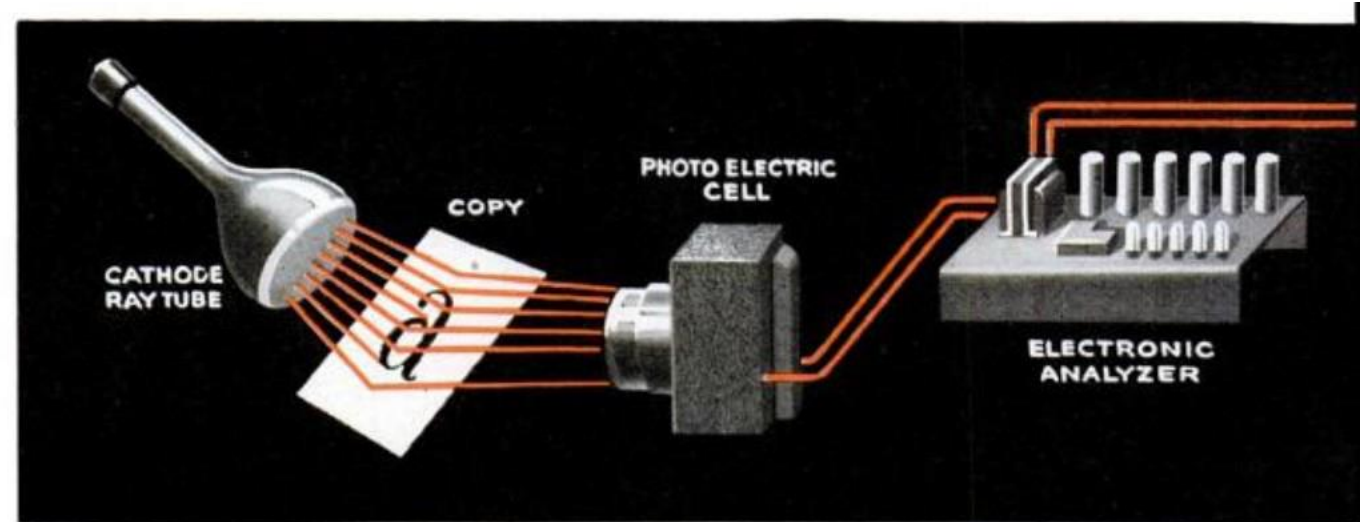- Another way is to use the package "pdftools".

- To get the text content from a PDF file use

    > my.text <- pdf_text(myPDFfiles[1])

- To save text content as txt file use

    > write.table(my.text, file=paste0(pdf.loc,"/text.txt"), quote = FALSE, row.names = FALSE, col.names = FALSE, eol = " " )

- To convert and save several pdf files you can create a function as illustrated below.

```
1   # Extracting the content from a pdf file
2   rm(list=ls()); cat("\014") # Clear Workspace and Console
3   library(pdftools)
4
5   pdf.loc <- file.path(getwd(),"PDF Files") # folder "PDF Files" with PDFs
6   myPDFfiles <- normalizePath(list.files(path = pdf.loc, pattern = "pdf",  full.names = TRUE)) # Get the path (chr-vector) of PDF file names
7
8   my.text <- pdf_text(myPDFfiles[1]) # Get the text content from the PDF file
9   write.table(my.text, file=paste0(pdf.loc,"/text.txt"), quote = FALSE, row.names = FALSE, col.names = FALSE, eol = " " ) # Save as txt file
10
11
12  # Convert to text several pdf files that are contained in a single folder.
13  convert.PDF <- function(myPDFfiles) {
14    for (ff in 1:length(myPDFfiles)) {
15      pdf.file <- myPDFfiles[ff]
16      my.text <- pdf_text(pdf.file) # Get the text content from the PDF file
17      File.Name <- sub(".pdf",".txt",pdf.file)
18      write.table(my.text, file=File.Name, quote = FALSE, row.names = FALSE, col.names = FALSE, eol = " " ) # Save as txt file
19    }
20  }
21
22  convert.PDF(myPDFfiles)
23
24  # Use lapply with in line function to convert each PDF file indexed by "i" into a text file
25  lapply(1:length(myPDFfiles),
26         function(ff, myPDFfiles)
27           {my.text = pdf_text(myPDFfiles[ff]); write.table(my.text, file=sub(".pdf",".txt",myPDFfiles[ff]),
28                                  quote = FALSE, row.names = FALSE, col.names = FALSE, eol = " " )},
29         myPDFfiles)
30
```

# History of OCR

**Google Book Search initiative** – a servicer that searches the full text of books and magazines that Google has scanned, converted to text using optical character recognition (OCR), and stored in its digital database.

- It has opened up many avenues for future research in document understanding and recognition.
- Resulted in developing Google's Tesseract software.



How machine reads: Electric eye looks at letter and reports its shape to electronic "brain."

An image of RCA's 1949 OCR system: M. Martin, "Reading Machine Speaks Out Loud," Popular Science, vol. 154, no. 2, Feb 1949, pp. 125-127. Used under fair use, 2014. The system was discontinued prior to completion due to its high costs .

# Types of problems OCR encounters

This removes from the property list stored in *place* the property with an indicator eq to *indicator*. The property indicator and the corresponding value are removed by destructively splicing the property list. `remf` returns `nil` if no such property was found, or some non-`nil` value if a property was found. The form *place* may be any generalized variable acceptable to `setf`. See `remprop`.

y-critical system istinction can be tes. The mission haviour while the y controller when more, the aims of mission controller ed – this will also er into an unsafe ied with avoiding unsafe states that

y-critical system istinction can be tes. The mission haviour while the y controller when more, the aims of mission controller ed – this will also er into an unsafe ied with avoiding unsafe states that

$p_i(\mathbf{x}) = \mathbf{P}(\theta = \omega_i \mid \mathbf{X} = \mathbf{x}), \; i = 1, ..., c$ are the posteriori probabilities. Let $R^*$ denote the Bayes risk, i.e., the risk of the Bayes rule. In practice we rarely have any information about the distribution of the pair $(\theta, \mathbf{X})$, instead there is in our disposal a training set $\eta_n = \{(\theta_1, \mathbf{X}_1), ..., (\theta_n, \mathbf{X}_n)\}$, i.e., a sequence of pairs $(\theta, \mathbf{X})$ distributed like $(\theta, \mathbf{X})$, where $\mathbf{X}$ is the feature vector and $\theta$ is its class assignment. An empirical classification rule $\psi_n$ is a measurable function of $\mathbf{X}$ and $\eta_n$. It is natural to construct a rule which resembles the Bayes rule, i.e., by replacing $p_i(\mathbf{x})$ by its estimate $p_{in}(\mathbf{x})$. A popular nonparametric classification technique is the kernel classifier being defined as follows

$$\psi_n(\mathbf{x}) = \arg\max_{1 \le i \le c} \sum_{j=1}^{n} \mathbf{1}(\theta_j = \omega_i) W\left(\frac{\mathbf{x} - \mathbf{X}_j}{b}\right), \tag{1.1}$$

THE not so hidden COST$

MEALS $350

FUEL $600

LODGING $300

FUN $150

# Amazon Textract

Textract is a machine learning service that automatically extracts text, handwriting and data from scanned documents that goes beyond simple optical character recognition (OCR) to identify, understand, and extract data from forms and tables.

# Exercise: OCR with Tesseract

OCR (pattern recognition in general) is a very difficult problem for computers.

The R tesseract package provides R bindings to Google's OCR library Tesseract.

- It is a powerful optical character recognition (OCR) engine that supports over 100 languages. The engine is highly configurable in order to tune the detection algorithms and obtain the best possible results.

- Results are rarely perfect and the accuracy rapidly decreases with the quality of the input image. But if you can enhance your input images to reasonable quality, Tesseract can often help to extract most of the text from the image.

    - One such image enhancements used with together with tesseract is the "magick" package.

# Tesseract

**Using Tesseract, the following file types can be created:**

- **alto** — Output in ALTO format (*OUTPUTBASE*.xml).

- **hocr** — Output in hOCR format (*OUTPUTBASE*.hocr).

- **pdf** — Output PDF (*OUTPUTBASE*.pdf).

- **tsv** — Output TSV (*OUTPUTBASE*.tsv).

- **txt** — Output plain text (*OUTPUTBASE*.txt).

- **makebox** — Write box file (*OUTPUTBASE*.box).

- **get.images** — Write processed input images to file (tessinput.tif).

- **lstm.train** — Output files used by LSTM training (*OUTPUTBASE*.lstmf).

The general syntax for converting a *eng -* English language *png.file* into a text file (*txt*) called *output.file* is:

*system( paste("tesseract", png.file, output.file, ' --oem 1 -l eng txt'), wait=TRUE)*

# Install Tesseract

Search online where to find a link to install Tesseract for your OS

- One such site is: https://digi.bib.uni-mannheim.de/tesseract/

Once you download and run tesseract executable, on Windows OS you will need to add Tesseract path (most likely "C:\Program Files\Tesseract-OCR") to your path in the system's environment variable.

- Add "C:\Program Files\Tesseract-OCR"

To verify if the installation worked open command prompt in "C:\Program Files\Tesseract-OCR" and type

    tesseract --version

If you see any error like

    tesseract command not found

most probably you have made some mistake in the installation.

# Example 1: OCR Image to Extract Text

OCR a PNG File to extract the text and save it as text file.

1. Start with a pdf that we are going to turn into a png file (image)

2. OCR the PNG File to extract the text and save it into a file.

```r
1   rm(list=ls()); cat("\014") # Clear Workspace and Console
2   library("pdftools")
3   library("tesseract"); library("magick"); library('tabulizer')
4
5   # File Locations
6   pdf.file <- "Data/myPDFfile_pg1.pdf"
7   png.file <- "Data/myPDFfile_pg1.png"
8
9   # 1) Create PNG Image from PDF files (using  pdftools package)
10  pngfile <- pdf_convert(pdf.file, pages = 1, filenames = png.file, dpi = 100)
11
12  # 2) OCR a PNG => TXT
13  txt.file <- sub('.png','', png.file)
14  system( paste("tesseract", png.file, txt.file), wait=TRUE)
15
```

*PDF File*

*PNG File*

*TXT File*

# Example 2: Tesseract and Magick

The Life and Work of
Fredson Bowers

*by*

G. THOMAS TANSELLE

IN EVERY FIELD OF ENDEAVOR THERE ARE A FEW FIGURES WHOSE ACCOM-plishment and influence cause them to be the symbols of their age; their careers and oeuvres become the touchstones by which the field is measured and its history told. In the related pursuits of analytical and descriptive bibliography, textual criticism, and scholarly editing, Fredson Bowers was such a figure, dominating the four decades after 1949, when his *Principles of Bibliographical Description* was published. By 1973 the period was already being called "the age of Bowers": in that year Norman Sanders, writing the chapter on textual scholarship for Stanley Wells's *Shakespeare: Select Bibliographies*, gave this title to a section of his essay. For most people, it would be achievement enough to rise to such a position in a field as complex as Shakespearean textual studies; but Bowers played an equally important role in other areas. Editors of nineteenth-century American authors, for example, would also have to call the recent past "the age of Bowers," as would the writers of descriptive bibliographies of authors and presses. His ubiquity in the broad field of bibliographical and textual study, his seemingly complete possession of it, distinguished him from his illustrious predecessors and made him the personification of bibliographical scholarship in his time.

When in 1969 Bowers was awarded the Gold Medal of the Bibliographical Society in London, John Carter's citation referred to the *Principles* as "majestic," called Bowers's current projects "formidable," said that he had "imposed critical discipline" on the texts of several authors, described *Studies in Bibliography* as a "great and continuing achievement," and included among his characteristics "uncompromising seriousness of purpose" and "professional intensity." Bowers was not unaccustomed to such encomia, but he had also experienced his share of attacks: his scholarly positions were not universally popular, and he expressed them with an aggressiveness that almost seemed calculated to

- Image cleaning and pre-processing before OCR is always advised
  - It improves the quality of the OCR output into text
- Cleaning and pre-processing steps typically involve cropping out the text area, rescaling, increasing contrast, etc.
- The "magick" package is an excellent tool for this task, as illustrated by the following code.

```
1   # OCR Image to Extract Text and save it as searchable PDF
2   rm(list=ls()); cat("\014") # Clear Workspace and Console
3   library("tesseract"); library("magick")
4
5   input <- image_read("Data/Image_1.jpg")
6
7   # 1) JPG => PNG => OCR => TXT
8   text <- input %>%
9       image_resize("2000x") %>%
10      image_convert(type = 'Grayscale') %>%
11      image_trim(fuzz = 40) %>%
12      image_write(format = 'png', density = '300x300') %>%
13      tesseract::ocr()
14
15  # 2) OCR a PNG => TXT
16  cat(text) # Display in Console
17  cat(as(text, "character"), sep = "\n", file = 'Data/Image_1.txt', append = FALSE) # Save as txt file
```

```
> cat(text)
The Life and Work of
Fredson Bowers
by
G. THOMAS TANSELLE

N EVERY FIELD OF ENDEAVOR THERE ARE A FEW FIGURES WHOSE ACCOM-
plishment and influence cause them to be the symbols of their age;
their careers and oeuvres become the touchstones by which the
field is measured and its history told. In the related pursuits of
analytical and descriptive bibliography, textual criticism, and scholarly
editing, Fredson Bowers was such a figure, dominating the four decades
after 1949, when his Principles of Bibliographical Description was pub-
lished. By 1973 the period was already being called "the age of Bowers":
in that year Norman Sanders, writing the chapter on textual scholarship
for Stanley Wells's Shakespeare: Select Bibliographies, gave this title to
a section of his essay. For most people, it would be achievement enough
to rise to such a position in a field as complex as Shakespearean textual
studies; but Bowers played an equally important role in other areas.
Editors of nineteenth-century American authors, for example, would
also have to call the recent past "the age of Bowers," as would the writers
of descriptive bibliographies of authors and presses. His ubiquity in
the broad field of bibliographical and textual study, his seemingly com-
plete possession of it, distinguished him from his illustrious predeces-
sors and made him the personification of bibliographical scholarship in

his time.
```

# Example 3: Searchable PDF from JPG

The Life and Work of
Fredson Bowers
*by*
G. THOMAS TANSELLE

IN EVERY FIELD OF ENDEAVOR THERE ARE A FEW FIGURES WHOSE ACCOM-
plishment and influence cause them to be the symbols of their age; their careers and oeuvres become the touchstones by which the field is measured and its history told. In the related pursuits of analytical and descriptive bibliography, textual criticism, and scholarly editing, Fredson Bowers was such a figure, dominating the four decades after 1949, when his *Principles of Bibliographical Description* was pub-lished. By 1973 the period was already being called "the age of Bowers": in that year Norman Sanders, writing the chapter on textual scholarship for Stanley Wells's *Shakespeare: Select Bibliographies*, gave this title to a section of his essay. For most people, it would be achievement enough to rise to such a position in a field as complex as Shakespearean textual studies; but Bowers played an equally important role in other areas. Editors of nineteenth-century American authors, for example, would also have to call the recent past "the age of Bowers," as would the writers of descriptive bibliographies of authors and presses. His ubiquity in the broad field of bibliographical and textual study, his seemingly com-plete possession of it, distinguished him from his illustrious predeces-sors and made him the personification of bibliographical scholarship in his time.

When in 1969 Bowers was awarded the Gold Medal of the Biblio-graphical Society in London, John Carter's citation referred to the *Principles* as "majestic," called Bowers's current projects "formidable," said that he had "imposed critical discipline" on the texts of several authors, described *Studies in Bibliography* as a "great and continuing achievement," and included among his characteristics "uncompromising seriousness of purpose" and "professional intensity." Bowers was not unaccustomed to such encomia, but he had also experienced his share of attacks: his scholarly positions were not universally popular, and he expressed them with an aggressiveness that almost seemed calculated to

- Image cleaning and pre-processing before OCR is always advised
  - It improves the quality of the OCR output into text
- Cleaning and pre-processing steps typically involve cropping out the text area, rescaling, increasing contrast, etc.
- The "magick" package is an excellent tool for this task, as illustrated by the following code.

```r
# OCR Image to Extract Text and save it as searchable PDF
rm(list=ls()); cat("\014") # Clear Workspace and Console
library("tesseract"); library("magick")

file_name <- "Data/Image_1.jpg"
input <- image_read(file_name)

# 1) JPG => PNG
input %>%
  image_resize("2000x") %>%
  image_convert(type = 'Grayscale') %>%
  image_trim(fuzz = 40) %>%
  image_write('Data/Image_1.png', format = 'png', density = '300x300')

# 2) OCR a PNG => PDF
png.file <- sub('.jpg','.png', file_name)
pdf.file <- sub('.jpg','', file_name)
system( paste("tesseract", png.file, pdf.file, ' -l eng pdf'), wait=TRUE)
```

Find (1/1)
Life
Previous   Next

The Life and Work of
Fredson Bowers
*by*
G. THOMAS TANSELLE

IN EVERY FIELD OF ENDEAVOR THERE ARE A FEW FIGURES WHOSE ACCOM-
plishment and influence cause them to be the symbols of their age; their careers and oeuvres become the touchstones by which the field is measured and its history told. In the related pursuits of analytical and descriptive bibliography, textual criticism, and scholarly editing, Fredson Bowers was such a figure, dominating the four decades after 1949, when his *Principles of Bibliographical Description* was pub-lished. By 1973 the period was already being called "the age of Bowers":

# Tesseract and Other Languages

- To better identify OCR-ed words, Tesseract has capability for installing and using vocabularies for additional languages. Use the tesseract_download() function to install additional languages:

  – tesseract_download("deu")

  To OCR German text use:

  – (german <- tesseract("deu"))

  – text <- ocr("Data/127193473.png", engine = german)

*Tesseract Output*

```
> cat(text)
Der Streuwald. 309
viele Waldungen mit ähnlicher Beitodung an und lejen wir in den Forft-
einrichtungsmwerken der abgelaufenen Beitabjehnitte nad, To zeigt fih, daß
diefe Waldungen eigentlich jchon immer in diefer BVerfaffung gewejen
find — die Holgnugung war gering, der Ausjhluß von der Stra
nußung hat feinen oder geringen Einfluß auf die Holzerzeugung aus=
zuüben vermocdt, Ummandlungen jcheiterten an der Schwierigfeit der
Aufforftung befonders bei mangelnden Geldmitteln.

Sn folden und ähnliden Waldungen haben wir zuerft
den Hebel anzujegen.

Füllen wir die Lüden mit viel Streuwerf abwerfenden und bildenden
Laube, Nadel- und Strauhhölzern — wie fie den Böden, dem Stlime,
der Lage entjprehen — aus und betradten wir die Holnußung als
Nebennugung, jo haben wir einen Wald, der, wenn aud nicht augen=
blidlich, jo doch bald jahraus, jahrein oder in kurzem Wechjel zur Streu:
nusung herangezogen werben Fann, ohne daß wir befürchten müfjen, daß
er dabei zu Grunde geht.
                              ift: Die guten Teile der Waldungen können
                              rweiterten Streunußungsmechjel unterftellt werden.
                              e Weije Iaffen fih Streuwaldungen jhaffen.

                              lder, Wiejen, Ödungen, Abhänge, Gruben,
                              einer richtigen Ausnügung. Dur Aufforftung mit
                              und verjchiedenen Verbefferungen fönnen Ddieje
                              ur Streugewinnung eingerichtet werden. Dbmwohl die
                              ftung folder Orte im Wald jhon vielfach erfolgt
                              t viel gefchehen, vielleicht zieht Die Sadhe unter dem
                              befier, bejonders wenn fie durch unentgeltliche Ab-
                              eförbert und durch die Staatsforfiverwaltung ge-

                              ine Forderung der Zeit, fuhen wir ihr
                              e gerecht zu werden.
```

```r
1   # Tesseract and Other Languages
2   rm(list=ls()); cat("\014") # Clear Workspace and Console
3   library("tesseract")
4
5   # Use the tesseract_download() function to install additional languages:
6   # tesseract_download("deu")
7
8   (german <- tesseract("deu"))
9   text <- ocr("Data/127193473.png", engine = german)
10
11  cat(text) # Display in Console
12  cat(as(text, "character"), sep = "\n", file = 'Data/Image_1.txt',
13      append = FALSE) # Save as txt file
14
```

Der Streuwald. 309

viele Waldungen mit ähnlicher Beſtockung an und leſen wir in den Forſt-
einrichtungswerken der abgelaufenen Zeitabſchnitte nach, ſo zeigt ſich, daß
dieſe Waldungen eigentlich ſchon immer in dieſer Verfaſſung geweſen
ſind — die Holznutzung war gering, der Ausſchluß von der Streu-
nutzung hat keinen oder geringen Einfluß auf die Holzerzeugung aus-
zuüben vermocht, Umwandlungen ſcheiterten an der Schwierigkeit der
Aufforſtung beſonders bei mangelnden Geldmitteln.

In ſolchen und ähnlichen Waldungen haben wir zuerſt
den Hebel anzuſetzen.

Füllen wir die Lücken mit viel Streuwerk abwerfenden und bildenden
Laub-, Nadel- und Strauchhölzern — wie ſie den Böden, dem Klima,
der Lage entſprechen — aus und betrachten wir die Holznutzung als
Nebennutzung, ſo haben wir einen Wald, der, wenn auch nicht augen-
blicklich, ſo doch bald jahraus, jahrein oder in kurzem Wechſel zur Streu-
nutzung herangezogen werden kann, ohne daß wir befürchten müſſen, daß
er dabei zu Grunde geht.

Die Hauptſache aber iſt: Die guten Teile der Waldungen können
geſchont und einem erweiterten Streunutzungswechſel unterſtellt werden.

Aber auch auf andere Weiſe laſſen ſich Streuwaldungen ſchaffen.

Viele unrentable Felder, Wieſen, Ödungen, Abhänge, Gruben,
Sümpfe uff. harren einer richtigen Ausnützung. Durch Aufforſtung mit
paſſenden Holzarten und verſchiedenen Verbeſſerungen können dieſe
Ländereien leicht zur Streugewinnung eingerichtet werden. Obwohl die
Anregung zur Aufforſtung ſolcher Orte im Wald ſchon vielfach erfolgt
iſt, iſt bisher nicht viel geſchehen, vielleicht zieht die Sache unter dem
Schlagwort „Streu" beſſer, beſonders wenn ſie durch unentgeltliche Ab-
gabe von Pflanzen gefördert und durch die Staatsforſtverwaltung ge-
leitet wird.

Der Streuwald iſt eine Forderung der Zeit, ſuchen wir ihr
auf mannigfache Weiſe gerecht zu werden.

_____

Die Linde im Pfälzerwald und in den übrigen Waldgebieten
der Pfalz.

Von Johann Keiper.

(Fortſetzung.)

Zu Fragen 3 und 4.

Von einigen Ausnahmen der Zwiſchenſtändigkeit abgeſehen, tritt die
Linde beider Arten in den pfälziſchen Hochwaldungen hauptſtändig auf

22*

# Example 4: OCR a PNG and create HTML file that looks the same

*Original PNG Image*

*Generated HTML*

SCIENCE ADVANCES | RESEARCH ARTICLE

NETWORK SCIENCE

## Optimal network topology for responsive collective behavior

David Mateo[1]*, Nikolaj Horsevad[1], Vahid Hassani[1], Mohammadreza Chamanbaz[1], Roland Bouffanais[1]

Animals, humans, and multi-robot systems operate in dynamic environments, where the ability to respond to changing circumstances is paramount. An effective collective response requires suitable information transfer among agents and thus critically depends on the interaction network. To investigate the influence of the network topology on collective response, we consider an archetypal model of distributed decision-making and study the capacity of the system to follow a driving signal for varying topologies and system sizes. Experiments with a swarm of robots reveal a nontrivial relationship between frequency of the driving signal and optimal network topology. The emergent collective response to slow-changing perturbations increases with the degree of the interaction network, but the opposite is true for the response to fast-changing ones. These results have far-reaching implications for the design and understanding of distributed systems: a dynamic rewiring of the inter-action network is essential to effective collective operations at different time scales.

Copyright © 2019

The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S.Government Works. Distributed under a Creative Commons Attribution NonCommercial License 4.0 (CC BY-NC).

### INTRODUCTION

A wide range of complex systems are characterized by relatively simple dynamical rules while still producing excessively complex emergent col-lective behaviors. Examples abound in the natural world [e.g., a flock of birds, a school of fish, a swarm of insects (J-9)], in social systems [e.g., social networks (10-12)], and in engineered multi-agent systems [e.g., self-organized networks of mobile sensors, multi-vehicle coordination, and swarm robotics systems (13-16)].

Historically, particular attention has been directed toward in-vestigating varieties of collective behaviors obtained by testing a wide range of local agent-to-agent interaction rules (6, 9). Collective beha-viors have also been investigated from the network-theoretic perspec-tive (4, 8, 17-21). It is now clear that such rich collective behaviors are the outcome of a complex interplay between network topology— characteristic of the group-level organization—and the dynamical laws at the agent's level (4, 8, 20-22).

Many collective behaviors can be studied through the lens of distrib-uted consensus problems, including collective motion in animal groups and multi-robot systems. Over the past decade, the number of studies on

# Example 4: OCR a PNG to create HTML file

Convert the PNG image to a HTML file using the following syntax:

- "tesseract Data/www/myPDFfile_pg1.png Data/www/myPDFfile_pg1.png --oem 1 -l eng hocr"
- This will produce a file with "hocr" extension that you need to rename it with "html" extension before you open it in a browser.

```r
1   # Convert image of scanned PDF file into HTML
2   rm(list=ls()); cat("\014") # Clear Workspace and Console
3   library("pdftools")
4   library("tesseract"); library("magick")
5
6   pdf.file <- "Data/www/myPDFfile_pg1.pdf"
7   png.file <- "Data/www/myPDFfile_pg1.png"
8   hocr.file <- "Data/www/myPDFfile_pg1.png"
9
10  # 1) Convert PDF => PNG
11  file.remove(png.file) # Remove existing PNG file
12  # system( paste("convert -density 200", pdf.file, '-alpha remove -quality 200 -scale 125%', png.file), wait=TRUE)
13  pdf.file <- pdf_convert(pdf.file, pages = 1, filenames = png.file, dpi = 200)
14
15  # 2) OCR the PNG file, Extract text & Create searchable HTML
16  system( paste("tesseract", png.file, hocr.file, ' --oem 1 -l eng hocr'), wait=TRUE) # OCR PNG & Convert to HTML
17
18  # At the end just rename ".hocr" file into ".html" and open it in a browswer.
19
```

# Find Coordinates of the OCR words

Another Tesseract feature is to

- OCR a PDF page,
- Find Coordinates and
- recreate page with OCR'ed text



```r
# OCR PDF page, Find Coordinates and recreate page with OCR'ed text
rm(list=ls()); cat("\014") # Clear Workspace and Console
library(tesseract)
library(grid)
eng <- tesseract("eng")
pdf.file <- normalizePath(list.files(path = "Data/", pattern = "pdf",  full.names = TRUE))[1]
pdf.file <- "Data/135737664.pdf"

image.file <- pdftools::pdf_convert(pdf.file, format = 'tiff', pages = 1, dpi = 400)
results <- tesseract::ocr_data(image.file, engine = eng)
results.XML <- tesseract::ocr(image.file, engine = eng, HOCR=TRUE)
results

# Get Words & their coordiates
words <- unlist(lapply(results$word, function(x)  x))
wcoord <- do.call('rbind', lapply(results$bbox, function(x) as.numeric( unlist(strsplit(x, ",")))))

# Re-Scale coordinates
z <- data.frame(words=words, coord=wcoord, stringsAsFactors = FALSE)
co.x <- z$coord.1/max(z$coord.1); co.y <- (max(z$coord.2) - z$coord.2)/max(z$coord.2)
zz <- data.frame(words=words, x=co.x, y=co.y, stringsAsFactors = FALSE)

# ==== Plot extracted text into a grid
grid.newpage()
draw.text <- function(txt, x, y, just) {
  grid.text(txt, x, y, just=just, gp=gpar(col="grey", fontsize=8))
  # grid.text(txt, x=x[j], y=y[i], just=just)
  # grid.text(deparse(substitute(just)), x=x[j], y=y[i] + unit(2, "lines"),
  #           gp=gpar(col="grey", fontsize=8))
}

draw.text(zz$words, zz$x, zz$y, "left")
```