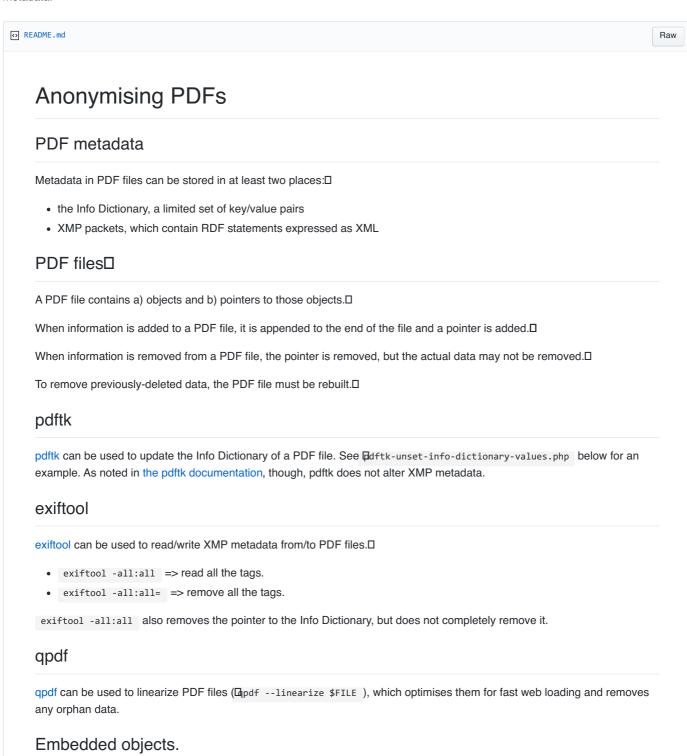


Remove metadata from a PDF file, using exiftool and qpdf. Note that embedded objects may still contain metadata.



After running qpdf, there may be new XMP metadata, as it extracts metadata from any embedded objects. To read the XMP tags of embedded objects, use exiftool -extractEmbedded -all:all \$FILE .

```
    pdftk-unset-info-dictionary-values.php

                                                                                                                                          Raw
       <?php
      $file = 'example.pdf';
   5
       // get the current metadata
   6
       $command = sprintf('pdftk %s dump data', escapeshellarg($file));
       $output = array(); $return = null; exec($command, $output, $return);
   9
       //print_r($output);
  10
       if ($return) {
           throw new Exception('There was an error reading metadata from the PDF file');
  13 }
  14
       // set any metadata values to null
  16
       foreach ($output as $index => $line) {
          if (strpos($line, 'InfoValue:') === 0) {
              $output[$index] = 'InfoValue:';
  18
  19
  20 }
       // write the updated metadata to a file
       $metadataFile = tempnam(sys get temp dir(), 'pdf-meta-');
  24
       file_put_contents($metadataFile, implode("\n", $output));
  26
       // create a new PDF using the updated metadata
       $tmpFile = tempnam(sys_get_temp_dir(), 'pdf-tmp-');
  $$ $command = sprintf('pdftk %s update_info %s output %s',
  29
         escapeshellarg($file), escapeshellarg($metadataFile), escapeshellarg($tmpFile));
  $ $output = array(); $return = null; exec($command, $output, $return);
       if ($return) {
           throw new Exception('There was an error writing metadata to the PDF file');
  3.4
  36 // clean up the temporary files
  37 rename($tmpFile, $file);
  38 unlink($metadataFile);
```

```
remove-pdf-metadata.sh
                                                                                                                                       Raw
       #!/bin/bash
      FILE=example.pdf
   5 # read tags from the original PDF
   6 #exiftool -all:all $FILE
   8
      # remove tags (XMP + metadata) from the PDF
       exiftool -all:all= $FILE
  10
  # linearize the file to remove orphan data
  12 qpdf --linearize $FILE
  14 \, # read XMP from the modified PDF
      #exiftool -all:all $FILE
       # read all strings from the modified PDF
  18
       #cat $FILE | strings > $FILE.txt
  20
      # read XMP from embedded objects in the modified PDF
       #exiftool -extractEmbedded -all:all $FILE
```



Nichtraucher commented Dec 29, 2013

Could you possibly add functionality that makes it possible to a) remove metadata for files in a directory (and its subdirectories), and b) make it a Nautilus script (in order to edit metadata in selected files/directories)? That would make it a lot easier to use! cheers!



This is short enough to make it a shell function.

```
clean_pdf() {
  pdftk $1 dump_data | \
    sed -e 's\(\infoValue:\)\s.*/\l\ /g' | \
    pdftk $1 update_info - output clean-$1

exiftool -all:all= clean-$1
  exiftool -all:all clean-$1
  exiftool -extractEmbedded -all:all clean-$1
  qpdf --linearize clean-$1 clean2-$1

pdftk clean2-$1 dump_data
  exiftool clean2-$1

pdfinfo -meta clean2-$1
}
```

via http://blog.snapdragon.cc/2015/08/28/shell-function-to-remove-all-metadata-from-pdf/



naught101 commented Aug 31, 2015

These methods don't seem to remove EXIF data from images embedded within a PDF. For example, the adobe photoshop editing history in a JPEG.



Telekor commented Sep 29, 2015

I have scripyt bymanuelRiel: now the word "clean" is added at the end of the file name (without extension). The Dicky line is this:

```
FILE="${FILE%%.*}"
```

And this is the fool script:

```
clean_pdf() {
   FILE=$1
   FILE="${FILE%%.*}"
   echo "##########"
   echo $1
   echo "##########"
   if [ -e $1 ]
       then
       pdftk $1 dump_data | \
       sed -e 's/\(InfoValue:\)\s.*/\1\ /g' | \
       pdftk $1 update_info - output ${FILE}.clean.pdf
       exiftool -all:all= ${FILE}.clean.pdf
       exiftool -all:all ${FILE}.clean.pdf
       exiftool -extractEmbedded -all:all ${FILE}.clean.pdf
       qpdf --linearize ${FILE}.clean.pdf ${FILE}.clean2.pdf
       pdftk ${FILE}.clean2.pdf1 dump_data
       exiftool ${FILE}.clean2.pdf
       echo "###########"
       echo "Metadata of file "${FILE}.clean2.pdf
       pdfinfo -meta ${FILE}.clean2.pdf
       echo "##########"
   else
       echo "File not found!"
}
```



rbreton74 commented Sep 30, 2015

Sorry, there is a small mistake. This one works fully:

```
clean_pdf() {
   FILE=$1
   FILE="${FILE%%.*}"
    echo "#########""
   echo $1
   echo "##########"
   if [ -e $1 ]
       then
       pdftk $1 dump_data | \
       sed -e 's/\(InfoValue:\)\s.*/\1\ /g' | \
       pdftk $1 update_info - output ${FILE}.clean.pdf
       exiftool -all:all= ${FILE}.clean.pdf
       exiftool -all:all \{FILE\}.clean.pdf
        exiftool -extractEmbedded -all:all ${FILE}.clean.pdf
       qpdf --linearize ${FILE}.clean.pdf ${FILE}.clean2.pdf
       pdftk ${FILE}.clean2.pdf dump data
       exiftool ${FILE}.clean2.pdf
       echo "#########"
       echo "Metadatos de fichero "${FILE}.clean2.pdf
       pdfinfo -meta ${FILE}.clean2.pdf
       echo "##########"
    else
       echo "File not found!"
        fi
}
```



ande2101 commented Nov 22, 2015

The previous script doesn't work with files with spaces in the filename.  $\square$ 



Changaco commented Dec 20, 2015

Other option: install pdf-redact-tools and run pdf-redact-tools -s \$FILE



danielneis commented Mar 7, 2016

This pdf-redact-tools uses exiftool to remove some tags as you can see in <a href="https://github.com/firstlookmedia/pdf-redact-tools/blob/master/pdf-D">https://github.com/firstlookmedia/pdf-redact-tools/blob/master/pdf-D</a> redact-tools#L115



bluesceada commented Sep 24, 2016 • edited

Hi, I wonder if exiftool is still a valid (or ever was) approach.

If I run exiftool on my file it warns me that tags are not really removed:  $\!\square\!$ 

```
$ exiftool -all:all= myfile.pdf
Warning: [minor] ExifTool PDF edits are reversible. Deleted tags may be recovered! - myfile.pdf
1 image files updated
```

My files also grow from 542.9 kB to 543.2 kB by exiftool and then from 543.2 kB to 544.6 kB by qpdf. So it seems there is actually more information added?

Let's see if these pdf-redact-tools do anything more. However I do for sure not want to follow one of their approaches, stacking PNG files and call it a PDF (that won't be searchable, has no vector graphic figures, and is probably larger in file size...)

//edit: OK that is actually the only approach they support, that's not applicable for me (and shouldn't be for most people that don't want to give away very larger or bad quality PDFs)



bertalanimre commented Feb 2, 2017

I'm wondering if there is any enhancement for this script to remove the embeded meetadata as well. For example, I have metadata about embeded Word documents. Unfortunatelly after the cleaning some sensitive data remains like filename, DocumentID and Instance ID. ☐

How can I delete these embeded metadata in the fist place? ☐



## bluesceada:

Unfortunately, exiftool was never a really sanitizing approach due to its limitation: http://www.sno.phy.queensu.ca/%7Ephil/exiftool/ - "Writer Limitations: PDF - The original metadata is never actually removed."

## Rut

qpdf --pages myfile.pdf 1-z -- --empty clean-myfile.pdf□

/\* creates a new (empty) PDF document from scratch and add (all: 1-z) the pages from the original PDF file into it \*/□ does the trick as the top-level (=file itself) metadata are concerned. It does not clean metadata of embedded objects.□

(Remark 1.: It is possible to use

pdftk myfile.pdf cat 1-end output clean-myfile.pdf  $\square$ 

instead abovementioned as well.

Remark 2.: On MS Windows, you can use BeCyPDFMetaEdit to obtain the same result, too; but for PDF version >1.6 the result is not guaranteed.)

## bertalanimre:

It may perhaps be done by filtering the PDF file through an editor (sed, tr?) capable of deleting characters between (and including) "
<x:xmpmeta" and "</x:xmpmeta>" strings. But I have never needed it so never tried it.



## RootLUG commented Mar 23, 2017

Did you guys find saome way how to remove metadata from file like this: Inttps://publications.usa.gov/USAFileDnld.php? PubType=P&PubID=6099&httpGetPubID=0?

I also tried the approach suggested by @9991212 but there is still a lot of metadata left like Creator, For, Create Date etc... which should be on the top level PDF

Sign up for free

to join this conversation on GitHub. Already have an account? Sign in to comment

© 2017 GitHub,

Inc. Terms Privacy Security Status Help

(7)

Contact GitHub

API Training Shop Blog About