

pdfSweep: what it does and why you need it

pdfissep) is in filest add on that immove (inducti) sensitive information from a FDF Portable Countered Common document. Confidentiality is assured, because the inducted information cannot be recovered in a count in the original process, pdf-issep, deletes test and in says as sea endead coordinate, or as defined by a regular expression. After having pamel the medientry information in the original FDF document, a new document is created without the induction context.

Why do we need redaction?

What is pdfSweep?

Reduction can be useful whenever the publisher or author of a document wishes to take out certain information. Common use cases include:

- Freedom of Information Act (USA) and similar legi
 Commented declaración
 Proprietary information
 Proprietary information
 Trade secreta
 General Data Protection Regulation (Dursye)
 All data that would impact the privacy of people
 Social security numbers (USA)
 National register identification numbers
 Phone numbers
 Bank account details
 Bank account details
 Bank account details

A short history of redaction

In the part, 'reduction' simply involved printing a document, blacking out the necessary information and making a photocopy of the document. That way, all information covered by dark ink simply does not get copied. This worked because paper in a simple WYSWYG format. There is no hidden data, no metadata that needs to be eased.

The President set that there is every referee that our position in health is a trengly supported by the special three med we are committed to that trees. So, Durandoner says that we are not a test of very. This is incorrect, it would be well if it the development of Un-053 Training the set of the

Challenges in redacting PDF documents

- Text rendering instructions do not need to appear in logical (reading) order.
 Text rendering instructions do not always constitute complete worth/duraks of text.
 Text rendering instructions as the applied (bith, read, lost transformation, etc.).
 Text instructions can be applied (bith, read, lost transformation, etc.).
 PDI documents can contain metadata, which should also be checked.
 PDI documents can contain norigin. Thus adding the possibility of ginarric content.

Images can be defined as a series of drawing operations.

All instructions of a page have to be processed to determine where e coder to find a neat way to avoid drawing in that area.

Images can be added to a PDF document under many formats.

Solect those parts of the document that must be reducted: either by specifying the coordinates, or by inputiting a regular
expression that its your needs. We have already provided a substantial list of cermon regular expressions to do some
of the heapy filing fory, such as social exactly numbers, places marriers, dates, etc.
 Pass the locations to pfilineep, or invoke pdiRutoSureep with the pattern(s) of your choice.

This is a pdfAutoSweep example that reducts the words Wilce' and White Rabbit' and Yabbit' (regardless of casing), it marks all occurences of Alabbit' with a pink rectangle, and all occurences of Rabbit' with a gay rectangle.

// define a composite strategy
CompositeCleansptErategy strategy = new CompositeCleansptErategy();
strategy addition spendiamedCleansptErategy(*Noo*) satReductionColor(close PRINT);
strategy addition should be repetitive should be repetitive strategy addition.

How does pdfSweep work?

CHIPTERS Some for includes.

All the training of the training state of the training stat

CHAPTON L Burn for The second is in the control of the the first house the same of the process containing here is for which the same of the same does, to reduce the limit of the state of th the bill past is a boost. Also und a fall as fas, I deal from entang of building about cloth from how four if it fast on others Was I would find an include about a control of the deal of the fast of the fall of the second of the control of the fall of the control of the control of the fall of the control of the control of the fall of the control of the control of the fall of the control of the

US_ZIP_CODE

Alternative example

- Select those parts of the document that must be reducted.
 Call pollutoSeeps with the method fentativeChars/b, This will produce a document where the content was not actually advantable, but where all content that supposedly must be reducted in marked with an amoutation. Adobe Acrobact can manage these amoutations for you, allowing you to remove reduction amoutations, add extra amoutations and add comments.

'Redacting' Image elements







How does pdfAutoSweep work?



Pre-configured regular expressions

NAME	MATCHES	
ROMAN_NUMERALS_STRICT	a Roman number, matching only those numbers that are valid.	
ROMAN_NUMERALS_FLEXIBLE	a Roman number, also matching III (4), whereas STRCT only matches $\ensuremath{\mathcal{W}}$	
US_SOCIAL_SECURITY_NUMBER	a US social security number (ddd-dd-dddd)	

US_CURRENCY DATE_MM_DD_YYYY_HH_MM_SS a date, specified in MM/DD/YYYY HHMMM:SS a date, specified in DD/MM/YYYC separator can be 'space' 7' or ' DATE_DD_MM_YYYY IPV6_ADDRESS MAC_ADDRESS a Media Access Control address

NAME MATCHES

To give a brief idea of the performance of pdffierep, we've used the iPhone user manual as a reference. We've reducted the regular expression 'III(Phone'. You can imagine this occurs quite a lot in that document.

COPIES	#PAGES	RINPUT FILE SIZE (MB)	#MILLISECONDS
1	130	3.23	6407
2	360	5.55	10253
4	520	ma	17946
	1040	22.5	23445

In this whilespaper, we're briefly presented our add-on-pdf/sweep, pdf/sweep allows you to seamlessly integrate data reduction is your existing worldfow. Reduction rectangels have to be defined, either by using pdf/kets/sweep which uses angular appression to world for marking rot and consoftante, not programmatically enterting the coordinates. pdf/sweep will reduct both to and images or existing complete confiderability.