



- [DOC Data Sets](#)
- [Description](#)
- [Creation](#)
- [Information](#)
- [Related Data Sets](#)
- [Provenance](#)
- [License](#)
 - [Edinburgh Napier University License Agreement](#)
 - [Open Government License Agreement](#)
- [Dataset Contents](#)

Description

This dataset contains 5,000 examples of documents created in the Microsoft© DOC format. This data set forms part of a much larger mixed file data set, called NapierOne, created by the School of Computing at Edinburgh Napier University in 2021. When using this dataset attribution should be given to Edinburgh Napier University as well as following the Open Government attribution guidelines. Details of which appear in the License section below. All files in the dataset have been validated to be of the correct format, having the correct extension and magic numbers where appropriate.

The files have also been scanned for viruses using Windows Defender© at the following level:

Scan Component	Component Version
Antimalware Client Version	4.18.2101.9
Engine Version	1.1.17800.5
Antivirus Version	1.331.1067.0
Anti-spyware Version	1.331.1067.0

Creation

To populate this dataset, the gov.uk websites and domains specified in the 'Provenance' section were queried. A combination of the following techniques were used to identify and retrieve the files:

- Searching the website using [Google Dorks](#)
- Browsing the website and identifying the data files.
- Using the Firefox plugin [Downloadthemall](#)

Examples of dorks used on these website are:

- site:<website name>
- filetype:csv
- filetype:doc
- filetype:docx
- filetype:jpg
- filetype:pdf
- filetype:ppt
- filetype:pptx
- filetype:xls
- filetype:xlsx

To provide file name consistency across the entire dataset, the retrieved files were renamed using the following syntax:

<sequence-number>-<original-extension>-<extra-information>.<extension>

Where:

File name fragment	Description
<sequence-number>	is unique number within this dataset. Files derived from this DOC dataset, for example the DOC-NOMAGIC and DOC-PASSWORD datasets, will retain the same sequence number across all datasets. This facilitates the easy identification of related files between related datasets.
<original-extension>	is the original extension that the file had, indicating the format of the file. Some datasets contain files, where the actual extension doesnot reflect the actual format of the original file. This especially relates to files that have been encrypted by malware.
<-extra-information>	is a placeholder in the file name that can be used to provided optional additional information relating to the file. For example if it is password protected.
<extension>	is the extension of the file name. Normally the characters after the '.' in a file name. The extension is normally used to identify the content and format of the file.

So for this sub dataset an example of a file name could be: 0001-doc.doc

A mapping of the original file name to the new dataset file name is held by the researchers and is available on request.

Information

If you require further information, then please feel free to contact the developers below:

Researcher	Contact
Simon Davies	s.davies@napier.ac.uk
Rich Macfarlane	r.macfarlane@napier.ac.uk
Bill Buchanan	b.buchanan@napier.ac.uk

Related Data Sets

Datasets that share the same sources as the dataset described here.

Dataset	Dataset description
CSS	A dataset containing examples of files found in the Cascading Style Sheet format found on the investigated websites.
CSV	A dataset containing examples of files found in the Comma Separated file format found on the investigated websites.
DOC-NOMAGIC	A copy of the DOC dataset, with the magic numbers removed from the file headers.
DOC-PASSWORD	A copy of the DOC dataset, the files being encrypted with a password. The password used is 'napierone'
DOCX	A dataset containing renamed modern versions of Microsoft Word document files found on the investigated websites.
DOCX-NOMAGIC	A copy of the DOCX dataset, with the magic numbers removed from the file headers.
DOCX-PASSWORD	A copy of the DOCX dataset, the files being encrypted with a password. The password used is 'napierone'
HTML	A dataset containing examples of files in the Hypertext Markup Language (HTML) format found on the investigated websites.
JAVASCRIPT	A dataset containing examples of files found in the Javascript(JS) file format found on the investigated websites.
JPG	A dataset containing examples of images in the JPG format found on the investigated websites.
JSON	A dataset containing examples of files in the Java Script Object Notation (JSON) format found on the investigated websites.
PDF	A dataset containing files found on the investigated websites that were in the Portable Document Format (PDF)
PDF-NOMAGIC	A copy of the PDF dataset, with the magic numbers removed from the file headers.
PDF-PASSWORD	A copy of the PDF dataset, with the magic numbers removed from the file headers.
PPT	A dataset containing renamed legacy versions of Microsoft Powerpoint document files found on the investigated websites.
PPT-NOMAGIC	A copy of the PPT dataset, with the magic numbers removed from the file headers.

Dataset	Dataset description
PPT-PASSWORD	A copy of the PPT dataset, the files being encrypted with a password. The password used is 'napierone'
PPTX	A dataset containing renamed modern versions of Microsoft Powerpoint document files found on the investigated websites.
PPTX-NOMAGIC	A copy of the PPTX dataset, with the magic numbers removed from the file headers.
PPTX-PASSWORD	A copy of the PPTX dataset, the files being encrypted with a password. The password used is 'napierone'
TXT	A dataset containing examples of flat files found on the investigated websites.
WEBP	A dataset containing examples of images in the WEBP format found on the investigated websites.
XLS	A dataset containing renamed legacy versions of Microsoft Excel document files found on the investigated websites.
XLS-NOMAGIC	A copy of the XLS dataset, with the magic numbers removed from the file headers.
XLS-PASSWORD	A copy of the XLS dataset, the files being encrypted with a password. The password used is 'napierone'
XLSX	A dataset containing renamed modern versions of Microsoft Excel document files found on the investigated websites.
XLSX-NOMAGIC	A copy of the XLSX dataset, with the magic numbers removed from the file headers.
XLSX-PASSWORD	A copy of the XLSX dataset, the files being encrypted with a password. The password used is 'napierone'
XML	A dataset containing examples of files in the Extensible Markup Language format found on investigating websites.

Provenance

The original files contained in this dataset were gathered from the following U.K. government domains and websites between January and March 2021.

Websites Accessed

acaf.food.gov.uk
acmsf.food.gov.uk
acnfp.food.gov.uk
acss.food.gov.uk
ahvla.defra.gov.uk
armedforcescomplaints.independent.gov.uk
www.armedforcescovenant.gov.uk
astro.ukho.gov.uk/
www.avonfire.gov.uk
www.bcomm-scotland.independent.gov.uk
www.belfastcity.gov.uk
www.birmingham.gov.uk
www.blogs.fco.gov.uk
boundarycommissionforengland.independent.gov.uk
www.brighton-hove.gov.uk
www.bristol.gov.uk
www.broads-authority.gov.uk
www.bucksfire.gov.uk
www.cabinetoffice.gov.uk
www.cafcass.gov.uk
www.cardiff.gov.uk
careers.dwp.gov.uk
www.castlepoint.gov.uk
www.ccrcc.gov.uk
census.gov.uk
www.childcarechoices.gov.uk
www.childrenscommissioner.gov.uk
www.cityoflondon.gov.uk
www.civilservicecommission.independent.gov.uk
www.complaints.judicialconduct.gov.uk
consult.defra.gov.uk

Websites Accessed

www.kirklees.gov.uk
www.lakedistrict.gov.uk
laqm.defra.gov.uk
www.lawcom.gov.uk
www.lbhf.gov.uk
www.leeds.gov.uk
www.legislation.gov.uk
www.leicester.gov.uk
www.liverpool.gov.uk
www.local.gov.uk
www.london.gov.uk
londonbridgeinquests.independent.gov.uk
www.londoncouncils.gov.uk
www.london-fire.gov.uk
lordsappointments.independent.gov.uk
www.luton.gov.uk
www.manchester.gov.uk
www.manchesterfire.gov.uk
www.merseysidewda.gov.uk
www.metoffice.gov.uk
www.mi5.gov.uk
www.midsussex.gov.uk
naei.beis.gov.uk
www.nationalarchives.gov.uk
nationalcareersservice.direct.gov.uk
www.nationalcrimeagency.gov.uk
www.natlotcomm.gov.uk
www.naturalresourceswales.gov.uk
www.ncsc.gov.uk
www.newforestnpa.gov.uk
www.newport.gov.uk

consult.education.gov.uk
consult.justice.gov.uk
cot.food.gov.uk
www.cps.gov.uk
www.cyberaware.gov.uk
www.danielmorganpanel.independent.gov.uk
www.dartmoor.gov.uk
www.data.gov.uk
data.justice.gov.uk
www.ddfire.gov.uk
www.derby.gov.uk
devtracker.dfid.gov.uk
www.dft.gov.uk/vca
digitalblog.ons.gov.uk
www.dsfire.gov.uk
dvlaregistrations.dvla.gov.uk
www.dwi.gov.uk
www.edinburgh.gov.uk
employerview.ofsted.gov.uk
www.essex-fire.gov.uk
europeanmemoranda.cabinetoffice.gov.uk
www.exmoor-nationalpark.gov.uk
www.faststream.gov.uk
www.ffc-environment-agency.metoffice.gov.uk
www.firescotland.gov.uk
www.food.gov.uk
www.forestry.gov.uk
www.fraudinvestigationjobs.co.uk
www.fsa.gov.uk
www.gamblingcommission.gov.uk
www.gchq.gov.uk
gdorb.decc.gov.uk
www.getingofar.gov.uk
getintoteaching.education.gov.uk
gigabitvoucher.culture.gov.uk
www.gla.gov.uk
www.glasgow.gov.uk
www.gloucestershire.gov.uk
www.gosportpanel.independent.gov.uk
www.gov.uk/performance/central-government-websites/website-domains
www.gov.wales
www.great.gov.uk
www.greatbusiness.gov.uk
www.hants.gov.uk
www.hantsfire.gov.uk
www.harrow.gov.uk
www.healthwatch.co.uk
www.helptobuy.gov.uk
helpwithchildarrangements.service.justice.gov.uk
www.hfea.gov.uk
panel.hillsborough.independent.gov.uk
hillsboroughinquests.independent.gov.uk
hireanapprentice.campaign.gov.uk
www.hmgcc.gov.uk
www.hse.gov.uk
www.hsl.gov.uk
www.nidirect.gov.uk
www.northampton.gov.uk
www.nottingham.gov.uk
www.notts-fire.gov.uk
www.nrscotland.gov.uk
www.ofgem.gov.uk
www.ofwat.gov.uk
old.food.gov.uk
www.online.hmrc.gov.uk/webchatprod/community/forums/list.page
www.ons.gov.uk
www.open.justice.gov.uk
www.orr.gov.uk
www.ownyourhome.gov.uk
www.oxford.gov.uk/downloadx
www.peakdistrict.gov.uk
www.pensionwise.gov.uk
www.plymouth.gov.uk
www.pointsoflight.gov.uk
www.policeconduct.gov.uk
www.portsmouth.gov.uk
www.ppo.gov.uk
www.princeofwales.gov.uk
www.prisonandprobationjobs.gov.uk
www.privycouncil.independent.gov.uk
publicappointments.cabinetoffice.gov.uk
publicappointmentscommissioner.independent.gov.uk
qavs.direct.gov.uk
www.reading.gov.uk
www.righttobuy.gov.uk
www.royal.gov.uk
www.rugby.gov.uk
science-council.food.gov.uk
service.gov.uk
www.sfo.gov.uk
sharp.dft.gov.uk
www.sheffield.gov.uk
www.sia.homeoffice.gov.uk
www.sis.gov.uk
www.smallbusinesscommissioner.gov.uk
www.southlanarkshire.gov.uk
www.southwales-fire.gov.uk
www.statisticsauthority.gov.uk
www.stoke.gov.uk
www.swansea.gov.uk
www.terrorismlegislationreviewer.independent.gov.uk
www.thepensionsregulator.gov.uk
think.direct.gov.uk
trade.great.gov.uk
www.tunisiainquests.independent.gov.uk
www.ukho.gov.uk
ukinventory.nda.gov.uk
www.uksport.gov.uk
www.understandinguniversalcredit.gov.uk
www.valuationtribunal.gov.uk
vehicleenquiry.service.gov.uk
www.wigan.gov.uk

www.hta.gov.uk
www.humbersidefire.gov.uk
iapdeathsincustody.independent.gov.uk
icai.independent.gov.uk
www.independent.gov.uk
www.invest.great.gov.uk/int
isc.independent.gov.uk
www.jncc.defra.gov.uk
www.judicialappointments.gov.uk
www.justice.gov.uk
www.justiceinspectores.gov.uk
www.keytosuccess.education.gov.uk

www.workplacepensions.gov.uk
www.ybtj.justice.gov.uk
www.york.gov.uk
www.yourpension.gov.uk

License

This dataset is covered by two separate license agreements. The Edinburgh Napier University license and the Open Government License. Both of which need to be respected and attribution given, when using this dataset.

Edinburgh Napier University License Agreement

ENU License Copyright (c) 2021 Edinburgh Napier University

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

Open Government License Agreement

The documents contained in this data set were recovered from gov.uk websites, these are covered by the Open Government License for public sector information.






Open Government Licence
for public sector information


<http://www.nationalarchives.gov.uk/doc/open-government-licence/version/3/>

The main points of which are:

You are free to:

-  Copy, publish, distribute and transmit the Information;
-  Adapt the Information;
-  Exploit the Information commercially and non-commercially for example, by combining it with other Information, or by including it in your own product or application.

You must (where you do any of the above):

 Acknowledge the source of the Information in your product or application by including or linking to any attribution statement specified by the Information Provider(s) and, where possible, provide a link to this licence;

If the Information Provider does not provide a specific attribution statement, you must use the following text:

Contains public sector information licensed under the Open Government Licence v3.0.

If elements of this DOC dataset are used then the correct attribution must also be given.

Dataset Contents

Documentation created at 13:41 on April 29, 2021