# Brave New World:
# Artificial Intelligence and Archives

**International Council on Archives**
Conseil International des Archives

Dr Anthea Seles, Secretary General of ICA

26 November 2019

14th General Conderence and Seminar

Tokyo, Japan

1

---

## Overview

**International Council on Archives**
Conseil International des Archives

- What is 'Artificial Intelligence'?
- Use of AI in government:
  - Acknowledging AI as evidence and archival record of the future
  - Ethical challenges and the role of archivists
- Impact of information management practices and the implications this has for using AI technologies
- Automating archival practice: appraisal, selection and sensitivity review
- Access and re-use of born-digital records in research and the use of automation in research

2

## Definitions

- DATA:
  - Structured data: Information, more often numerical information, put in tabular form to enable quantitative analysis.
  - Unstructured data: Information consisting of word processing documents, power point presentations, videos, sound records, photographs etc.

- ENVIRONNEMENT
  - Structured record-keeping environments: Environments where documents and data are placed in an ordered fashion to allow for retrieval. Ex: Information management system or shared drives with a unified classification scheme.
  - Non structured record-keeping environments: An environment where documents and information are not organised and can be comprised of a running sequence of document or a shared drive with no unified classification scheme.

3

## What is Artificial Intelligence(AI)?

- Artificial intelligence can be defined in many different ways; there is no standard definition
- There are really two categories
  - Supervised
  - Unsupervised
- Supervised: Requires a human to mark up or compile a homogeneous dataset to train an algorithm to recognise patterns or terms in the data. This process requires a lot of up front work and also requires you to have some level of understanding of the dataset.
- Unsupervised: Data is loaded into the system and without any upfront human intervention, analyses the data and provides result.
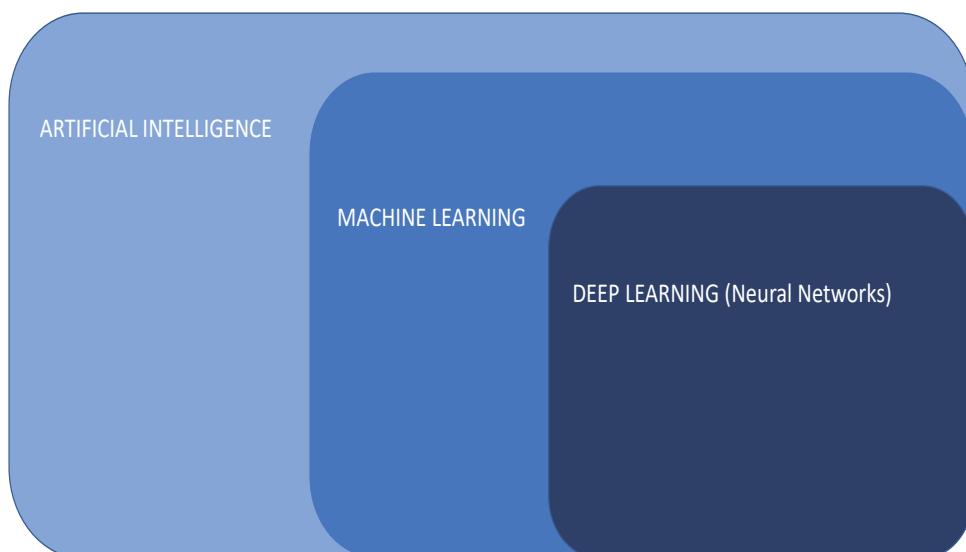
4

# Artificial Intelligence, Machine Learning and Neural Networks

- *Artificial Intelligence:* It's an all-encompassing definition for any activity where a machine/system takes information (structured and unstructured) to predict an outcome

- *Machine Learning*: Process of training a system to 'learn' how to make a decision using a pre-tagged dataset.

- *Neural Networks*: Just like we use our brains to identify patterns and classify information, neural networks can be trained to accomplish similar tasks.
  - Deep learning: Layering multiple neural networks

5

---

# Artificial Intelligence, Machine Learning and Deep Learning

ARTIFICIAL INTELLIGENCE

MACHINE LEARNING

DEEP LEARNING (Neural Networks)

Artificial Intelligence vs Machine Learning vs Deep Learning (7 May 2018) https://www.datasciencecentral.com/profiles/blogs/artificial-intelligence-vs-machine-learning-vs-deep-learning

## Archival Considerations

Impact of artificial intelligence, machine learning and data mining in government

Use of artificial intelligence in archival processes

Making records accessible and readable for research

7

## Government Use of Artificial Intelligence and *machine learning*

- Decisions are being made now using machine-learning and artificial intelligence
- These techniques are being used by data science or statistical analysis units in government departments and private corporations
  - Data science and the ability to mine data is seen as a competitive advantage.
  - Platforms that is common usage use these techniques: Netflix, Google, Facebook etc
  - For government it is seen as a way to parse through large volumes of data (structured and unstructured) to make a decision
  - Visualisations for policy decisions

8

# Government Use of Artificial Intelligence and *machine learning*

**International Council on Archives**
Conseil International des Archives

- There are challenges with the data science approach and the use of machine-learning and AI algorithms in government decision-making:
  - Is the data we are combining meant to be combined? Are we simply comparing apples and oranges?
  - Is the data biased and how does that affect the output of the algorithm? How does that affect what we see and how we interpret it?
- Archivists have often played a role in advising organisations on the creation and preservation of records and data to ensure their evidentiary value:
  - What advice would we give in the creation and preservation of 'algorithmic/computational records'?
  - Does the archivist have a role to play in advising how algorithms and code are created for decisions-making? How do we know what to preserve and how?
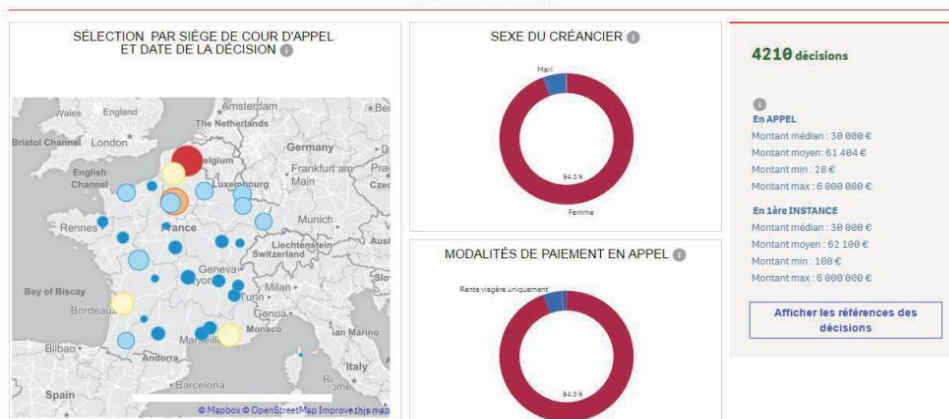
9

---

# Government Use of Artifical Intelligence and *machine learning*

**International Council on Archives**
Conseil International des Archives



10

## Slide 11

**ICA** International Council on Archives — Conseil International des Archives

# Government Use of Artificial Intelligence and *machine learning*

**Considerations:**
- If this becomes standard practice in government and passes into policy how do we begin to advise on what documentation needs to exist to document the training data and subsequent information that is input or not into the system? What does integrity and accountability look like in this context? By extension, what do we preserve?
- Does the archivist have a role as an ethical advisor in this context?
- To read the article: https://news.sky.com/story/handwriting-to-help-govt-catch-gangs-behind-mass-scale-benefit-fraud-11190448
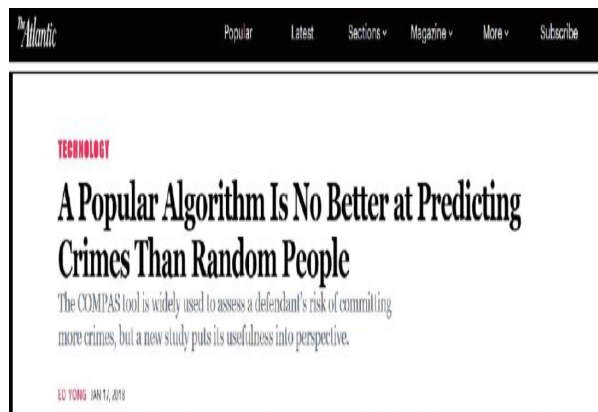
**sky news**

Home  UK  World  Politics  US  Ocean Rescue  Science & Tech  Business  Ents & Arts  Offbeat  Weather

### Handwriting to help Govt catch gangs behind mass-scale benefit fraud

Artificial intelligence is going to be used to clamp down on cheats claiming bogus benefit payments worth millions of pounds.

19:31, UK, Sunday 31 December 2017

CONCEPTION GRATUITE ET SANS ENGAGEMENT !
Les Bons Plans sont chez Darty*
-15%
sur les meubles Sorbonne, Concorde, Rio et Pérou

A record £1.1bn in overpaid benefits was recovered from fraudsters last year

11

## Slide 12

**ICA** International Council on Archives — Conseil International des Archives

# Government Use of Artificial Intelligence and *machine learning*

**Example:**
- Cathy O'Neil *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*
- In some US states they use algorithms to help determine recidivism rates (COMPAS- Correctional Offender Management Profiling for Alternative Sanctions)
- Some context of the data that was used to train COMPAS the algorithm created by Northpoint
  - Sentences given to African-American prisoners in the federal system is 20% longer than those given to white convicts for similar crimes
  - African-American represent 13% of the population of the United States, but account for 40% of the prison population
- Base training data set is biased and then the algorithm is created by a private company, which makes it a black box

*The Atlantic*  Popular  Latest  Sections ˅  Magazine ˅  More ˅  Subscribe

**TECHNOLOGY**

### A Popular Algorithm Is No Better at Predicting Crimes Than Random People

The COMPAS tool is widely used to assess a defendant's risk of committing more crimes, but a new study puts its usefulness into perspective.

ED YONG  JAN 17, 2018

12

# Why should this matter to you?

- Algorithms are the historical documents of tomorrow or NOW!
- Governments need to be held accountable if they use these technologies to make a decision that has an impact on the lives of their citizens, and we are responsible for identifying and preserving that information:
  - But what should we preserve? All the components that contributed to training the algorithm? (e.g. documents, data, social media information, algorithm and the results)? Only algorithm and the results?
- All this requires us to have the capacity and the skills to advise decision-makers in departments and ministries that are seeking to implement these technologies
  - Are we invited to the table?

13

# Why should this matter to you?

- Challenges and issues:
  - We will be responsible for preserving these algorithms in intermediate and historical archives
  - We are not currently considered stakeholders when it comes to discussions connected to the development and the implementation of AI technologies
  - We do not currently have the capacity or the skills to play our role as trusted adviser on information management questions related to AI records to ensure their preservation and durability.
  - We will need not only to advise decision-makers on the preservation of algorithms but we need to understand how to manage significant ethical challenges that will be posed by AI technologies
  - It is sometimes difficult to understand how an algorithm arrives at a result or decision, even if we preserve everything related to that decision.

14

## Impact of Information Management Practices

- Information management systems are not always easy to use, and they can be quite rigid, meaning that users try and find other, easier ways to file their information.
  - They use shared drives in parallel with information management systems, resulting in incomplete folders and duplication
- In the UK, we carried out a study to assess the state of record-keeping in government departments and understand the amount of 'legacy data' they held.
  - See: The Digital Landscape in Government 2014-2015
    http://www.nationalarchives.gov.uk/documents/digital-landscape-in-government-2014-15.pdf

15

## Impact of Information Management Practices

- Study results:
  - **1 TB: ~25 TB**
    - For each terabyte in an information management system there was about 25 TB in shared drives and this does not include data or information held in email servers.
  - **1.5 PB = approximately 1.5 billion Word documents**
    - Once we accounted for the totality of the information holdings which includes email servers and data sets it added up to over 1.5 petabytes of data that needed to be appraised and selected
  - Information management teams did not know what was contained in legacy data holdings and did not know what documents or data needed to be preserved
  - This information could also have differing levels of contextual information and limited metadata. Metadata could also be compromised because of previous migrations.

16

# Impact of Information Management Practices on Appraisal and Selection

- Volume can greatly complicate the appraisal and selection process, along with the ability of archivists to carry out large scale evaluations of unstructured data

- Due to the amount of information that required evaluation, we decided to start a second study to examine off the shelf systems that had *machine learning* capabilities for the purposes of assessing their viability to carry out appraisal and selection
  - See: *The Application of Technology Assisted Review to Born-Digital Records Transfers, Inquiries and Beyond*. http://www.nationalarchives.gov.uk/documents/technology-assisted-review-to-born-digital-records-transfer.pdf
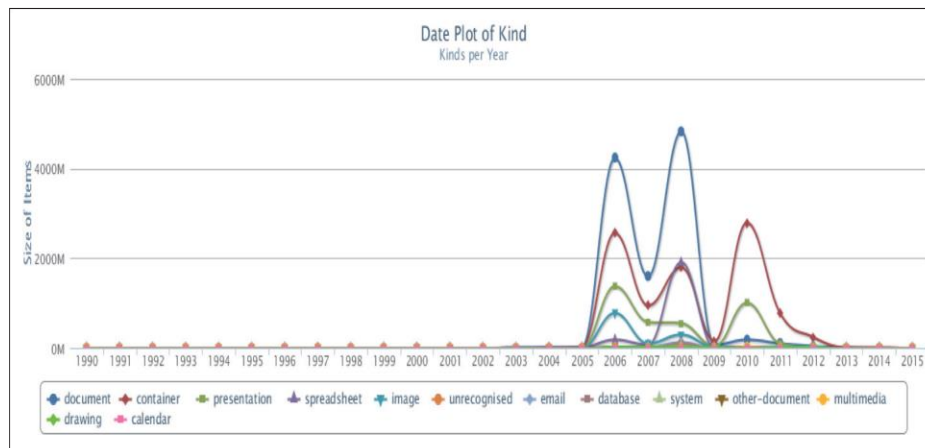
17

---

# Artificial Intelligence and *machine learning* in Records Management and Archives

- What can machines do well?
  - Boolean and keyword searches ✔
  - Regular expressions ✔
  - Process at scale ✔
  - Understand context and inference ✘
  - Handwriting analysis ✘
- What can humans do well?
  - Process at scale ✘
  - Understand and infer context ✔
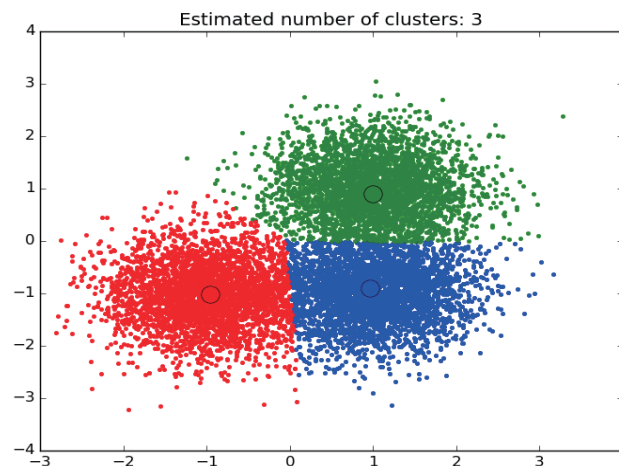  - Handwriting analysis ✔

18

# Artificial Intelligence and *machine learning* in Records Management and Archives

**Date Plot of Kind**
Kinds per Year

*Representation of a digital records collection by date and format*

# Artificial Intelligence and *machine learning* in Records Management and Archives

Estimated number of clusters: 3

*Concept clustering*

## Artificial Intelligence and *machine learning* in Records Management and Archives

- Problems and limits encountered during testing
  - Lack of understanding regarding the content and the context of creation
  - Corruption or alteration of metadata
  - Difficulty understanding the visualisations generated by the machines
  - Understanding the reliability (precision and recall) of the results and the acceptable level of risk
  - Distrust in technology and the results generated by the systems
    - However in other instances the results are accepted without question with an imprecise understanding of how the results were arrived at.
  - Significant time required to 'train' the system, departments wanted something much more automated (i.e. unsupervised)

21

## Artificial Intelligence and *machine learning* in Records Management and Archives

- Automation is no longer a choice, but a necessity. However, that does not mean that humans/archivists are irrelevant in this process
- The challenge with automating appraisal and selection, along with the sensitivity review process:
  - How do you measure accuracy? What does 'good enough' look like? What are the risks? What is acceptable risk appetite?
  - How can we determine what might be missing?
  - How can be accountable for the decisions we make based on machine outputs? How do we equally hold the machines to account?
  - How do we compensate for the change in the digital record over time? Re-tune the algorithm?
- We are dealing with 'Black Boxes'
- RISK: Biasing the historical record and by proxy history and our collective memory

22

# Artificial Intelligence and *machine learning* in Records Management and Archives

- Archival codes of ethics need to be studied and revised
- We are lacking the necessary competencies and skills to properly work with these types of technologies
- Algorithmic accountability and transparency
  - Corporations and businesses need to be accountable for how their machines arrive at a result or they must disclose the workings of their algorithms
  - Declaration of algorithmic transparency from the Association of Computing Machinery (ACM)
    - https://www.acm.org/binaries/content/assets/public-policy/2017_usacm_statement_algorithms.pdf
    - Seven (7) principles: Awareness; access and redress; accountability; explanation; data provenance; auditability; validation and testing
  - Partnership on AI – Partnership between Google, Microsoft, IBM et Facebook to promote AI for social good  https://www.partnershiponai.org/
  - Montreal Declaration: https://www.declarationmontreal-iaresponsable.com/
  - EU Regulations and principles around AI: https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai
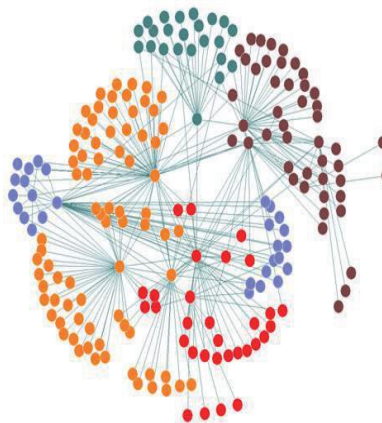
23

# Artificial Intelligence and Machine Learning in Research

- Two issues for the archival community to consider:
  - Impact of researchers trying to mine archival data
  - Digitisation of historical data and information

- Researchers are starting to use data mining techniques to parse through large volumes of digital data.
  - Ex: Researchers are using tools like Google NGRAM to mine literature to trace things like stereotypes in literature
    - Susan Mason. 'Analysing Stereotypes Across Time Using Google Ngram Viewer' *SAGE Research Methods Cases Part 2* (2018) doi:10.4135/9781526436245

- There are also many other tools, sometimes bespoke, that researchers are or will begin using.

24

## Artificial Intelligence and Machine Learning in Research

- There is a question for archivists about how much access we may wish to allow researchers access to public records and data
  - Data mining and machine learning tools breakdown siloes created by archival description (i.e. fonds, series, files)
  - Can reveal unknown connection that become sensitive or problematic by virtue of making that connection
  - Can surface sensitive information that was missed during sensitivity review
  - Also once the data is mined and put into a system outside the archives, what else can it can be combined to?
- Let's not get tunnel vision with AI. There is a danger of focusing too much on the impact on our individual collections, but what about linked data? And the semantic web? What will this mean for archives and opening our collections?

25

---

## Artificial Intelligence and Machine Learning in Research

- We also need to consider the impact of future digitisation.
  - The re-purposing and re-use of archival records and data has enormous value and I think we sacrificed much of digitisation and allowing companies to digitize archival records and data, in order that we can get a 'free' copy'. We must be savvier.
  - Companies are beginning to realise the value of data held in historical records. Digitising them and applying OCR is a method for gaining access to large volumes of data to train algorithms.
  - We need to start asking ourselves:
    - Why is the digitisation free?
    - Will this data be used to train an algorithm?
    - What is the company's ethical stance?
    - What happens to the data once the digitisation is done?
    - Will there be an impact on people's lives?
  - Scenario: Paper death registrations

26

# Conclusion

- Government Use of Artificial Intelligence:
  - What role does the archives and information communities have to play in this space? Do we have a role?
  - What skills do we have, or do we need if we have a role to play?
  - What is the 'record'? How do we capture and preserve that record?
  - Who are our partners? How do we begin to work with them?
- Machine Learning and Artificial Intelligence in Archival Processes
  - What is accuracy? What risks are we willing to accept?
  - How can we ensure the accountability of the decision we make based on machine-learning and AI processes?
- Artificial Intelligence and Machine Learning in Research
  - How much access is too much when machines are involved?
  - What are the right questions to ask when private companies offer us free digitisation?
  - How do researchers want to use our records to carry out digital research?

27

# A parting tought...

*Whether you are using an algorithm, artificial intelligence, or machine learning, one thing is certain: If the data being used if flawed, then the insights and information will be flawed.*

-Venkatesan M *Artificial Intelligence vs Machine Learning vs Deep Learning*

28

# References and Further Reading

- *The Application of Technology Assisted Review to Born-Digital Records Transfers, Inquiries and Beyond*. (2016) London:The National Archives UK http://www.nationalarchives.gov.uk/documents/technology-assisted-review-to-born-digital-records-transfer.pdf

- Bhaskar, Michael. *Curation: The Power of Selection in a World of Excess.* (2017) London:Piatkus

- Caplan,Robyn, Joan Donovan, Lauren Hanson and Jeanna Matthews. 'Algorithmic Accountability: A Primer' *Data and Society* (2018) https://datasociety.net/wp-content/uploads/2018/04/Data_Society_Algorithmic_Accountability_Primer_FINAL-4.pdf

- Chumtong, Jason and David Kaldewey. 'Beyond the Google NGRAM Viewer: Bibliographic Databases and Journal Archives As Tools for Quantitative Analysis of Scientific and Meta-Scientific Concepts. *FIW Working Paper No 8* (2017) https://www.fiw.uni-bonn.de/publikationen/FIWWorkingPaper/fiw-working-paper-no.-8

- Delort, Pierre. *Le Big Data* (2015) Paris : PUF

- Domingos, Pedro. *The Master Algorithm: How the Quest for the Ultimate Learning Machine Will Remake Our World* (2015) New York:Basic Books

- Engin, Zeynep and Philip Treleaven. 'Algorithmic Government: Automating Public Services and Supporting Civl Sevants in using Data Science Technologies' *The British Computer Society* (August 2018)https://academic.oup.com/comjnl/advance-article/doi/10.1093/comjnl/bxy082/5070384

- Ertzscheid, Oliver. *L'appétit des géants: pouvoir des algorithmes, ambitions des plateformes* (2017) Paris : C&F

- Information Privacy Commissioner. *Big Data, Artificial Intelligence, Machine Learning and Data Protection.* (2017) London:ICO https://ico.org.uk/media/for-organisations/documents/2013559/big-data-ai-and-ml-and-data-protection.pdf

- Jerven, Morten. *Poor Numbers: How We Are Misled by African Development Statistics and What to Do about It.* (2013) Ithica:Cornell University Press

- LeSueur, Andrew. 'Robot Government: Automated Decision-Making and its Implications for Parliament' [Draft chapter for publication in *Parliament: Legislation and Accountability* (Oxford:Hart Publishing) 2016] https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2668201

- Lorenzi, Jean-Hervé et Mickaël Berrebi. *L'avenir de notre liberté* (2017) Paris : Eyrolle

# References and Further Reading

- Lynch, Clifford. Stewardship in the 'Age of Algorithms' *First Monday* Vol 22 (12) December 2017 http://firstmonday.org/article/view/8097/6583

- Mason, Susan . 'Analysing Stereotypes Across Time Using Google Ngram Viewer' *SAGE Research Methods Cases Part 2* (2018) doi:10.4135/9781526436245

- Mason, S. E., C.V. Kuntz, & , C. M. McGill**.** 'Oldsters and ngrams: Age stereotypes across time'. *Psychological Reports: Sociocultural Issues in Psychology*, (2015),116, 324–329. doi:http://dx.doi.org/10.2466/17.10.PRO.116k17w6

- Musser, George. 'Artificial Intelligence: How Machines could learn creativity and common sense, among other human qualities'. *Scientific American* Vol 320, No 5 (May 2019) 47-51

- O'Neill, Cathy. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy* (2016) New York: Crown Publishing

- Padilla, Thomas, Laurie Allen, Sarah Potvin, Elizabeth Roke Russey, and Stewart Varner. 'Collections as Data', 7 March 2017. https://doi.org/10.17605/OSF.IO/MX6UK.

- Rolan, Gregory, Glen Humphries, Lisa Jeffrey, Evanthia Samaras, Tatiana Antsoupova and Katharine Stuart. 'More Human than Human? Artificial intelligence in the archive' *Archives and Manuscripts* Vol 47, No 2 (November 2018) 179-203

- Venkatesan M *Artificial Intelligence vs. Machine Learning vs. Deep Learning* (7 May 2018) https://www.datasciencecentral.com/profiles/blogs/artificial-intelligence-vs-machine-learning-vs-deep-learning

- Villani, Cédrique. *Donné un sens à l'intelligence artificielle: Pour une stratégie nationale et européenne* (8 mars 2018) https://www.aiforhumanity.fr/pdfs/9782111457089_Rapport_Villani_accessible.pdf

- World Wide Web Foundation. 'Algorithmic Accountability: Applying the Concept to Different Country Contexts'. *A Smart Web for a More Equal Future* (2017) https://webfoundation.org/docs/2017/07/Algorithms_Report_WF.pdf

- Zambonelli, Franco, Flora Salim, Seng W. Loke, Wolfgang De Meuter and Salil Kanhere. 'The Algorithmic Governance in Smart Cities: The Conundrum and the Potential of Pervasive Computing Solutions' *IEEE Technology and Society Magazine* (June 2018) pp 80-87

Thank you.

**Dr Anthea Seles**
Secretary General
International Council on Archives
**seles@ica.org**