RESPONSIBLE ARTIFICIAL INTELLIGENCE

Prof. Dr. Virginia Dignum

Chair of Social and Ethical Artificial Intelligence - Department of Computer Science

Email: virginia@cs.umu.se - Twitter: @vdignum



DECISION-MAKING AND ETHICS



"All my decisions are well thought out."





AI AND ETHICS - SOME CASES

- Self-driving cars
 - Who is responsible for the accident by self-driving car?
 - How can car decide in face of a moral dilemma?
- Automated manufacturing
 - How can technical advances combined with education programs (human resource development) help workers practice new sophisticated skills so as not to lose their jobs?
- Chatbots
 - Mistaken identity (is it a person or a bot?)
 - Manipulation of emotions / nudging / behaviour change support



WHAT WE TALK ABOUT WHEN WE TALK ABOUT AI

- Autonomy
- Decision-making
- Algorithms
- Robots
- Data
- Learning
- End of the world!?
- A better world for all?



WHAT IS AI? – NOT JUST THE ALGORITHM







WHAT IS AI? – NOT JUST MACHINE LEARNING



- Not all problems can be solved by correlation
 - Causality, abstraction...
- Rubbish in rubbish out
- Brittleness



ARTIFICIAL INTELLIGENCE





RESPONSIBLE AI: WHY CARE?

- AI systems act autonomously in our world
- Eventually, AI systems will make *better* decisions than humans

AI is designed, is an artefact

• We need to sure that the purpose put into the machine is the purpose which we really want

Norbert Wiener, 1960 (Stuart Russell) King Midas, c540 BCE



RESPONSIBLE AI

- AI can potentially do a lot. Should it?
- Who should decide?
- Which values should be considered? Whose values?
- How do we deal with dilemmas?
- How should values be prioritized?

•

How can we develop AI to benefit humanity?

TAKING RESPONSIBILITY

- Responsiblity / Ethics in Design
 - Ensuring that development <u>processes</u> take into account ethical and societal implications of AI as it integrates and replaces traditional systems and social structures
- Responsibility / Ethics **by** Design
 - Integration of ethical <u>reasoning</u> abilities as part of the behaviour of artificial autonomous systems
- Responsibility /Ethics <u>for</u> Design(ers)
 - Research integrity of <u>researchers</u> and manufacturers, and certification mechanisms



ETHICS <u>IN</u> DESIGN

- Design for values
- Design for all
- Doing the right thing
- Doing it right

Do things right, and do the right things."



PETER DRUCKER

ETHICS IN DESIGN - DOING THE RIGHT THING

- Taking an ethical perspective
 - Ethics is the new green
 - Business differentiation
 - Certification to ensure public acceptance



- Regulation is drive for transformation
 - Better solutions
 - Return on Investment





ETHICS IN DESIGN- DOING IT RIGHT

- Principles for Responsible AI = ART
 - <u>A</u>ccountability
 - Explanation and justification
 - Design for values
 - **<u>R</u>esponsibility**
 - Autonomy
 - Chain of responsible actors
 - Human-like AI
 - **<u>T</u>**ransparency
 - Data and processes
 - Algorithms





- AI systems (will) take decisions that have ethical grounds and consequences
- Many options, not one 'right' choice
- Need for design methods that ensure

ART METHODOLOGY

https://medium.com/@virginiadignum/on-biasblack-boxes-and-the-quest-for-transparency-inartificial-intelligence-bcde64f59f5b

- Socially accepted
 - Participatory
- Ethically acceptable
 - Ethical theories and human values
- Legally allowed
 - Laws and regulations
- Engineering principles
 - $\circ \ \ Cycle: Analyse synthetize evaluate repeat$
 - Report: Identify, Motivate, Document





ACCOUNTABILITY CHALLENGES

https://xkcd.com/1838/



- Optimal AI is explainable AI
- Optimal is not that which optimizes the result ignoring the context but the one that gives the best result for the context
- Design for values
 - include values of ethical importance in design
 - Explicit, systematic
 - Verifiable







RESPONSIBILITY CHALLENGES

- Chain of responsibility
 - researchers, developerers, manufacturers, users, owners, governments, ...
- Levels of autonomy
 - $\circ~$ Operational autonomy: Actions / plans
 - Decisional autonomy: Goas/ motives
 - Attainable autonomy: dependent on context and task complexity





Human-like AI

- Mistaken identity / expectations
- Vulnerable users: children / elderly

UMEÅ UNIVERSITY

TRANSPARENCY CHALLENGES

- Manage expectations
 - Training wheels / L-plates
- Openness
 - Data, processes, stakeholders
- Data
 - Bias is inherent in human behavior
 - Provenance: Where does it come from? Who is involved?
 - Training data: the cheapest/easiest or the best?
 - Data responsibility
 - Governance, storage, updated
 - Environment: energy costs
 - Market mechanisms



UMEÅ UNIVERSITY





ACCOUNTABILITY – DEALING WITH BIAS

• Bias

- Expectations derived from experienced regularities
- Heuristics used to deal with uncertainty produce bias
 - Portugal has the best footballers
 - Most programmers are male

• Stereotype

- those bias that we don't want to have persisting
- Most programmers are male

• **Prejudice**: acting on stereotypes

- *Hiring only male programmers*
- Bias are inherent on human data;
- We dont want AI to be prejudiced!
 - How to evaluate/clean existing data?
 - Historical, culturally dependent, contextual
 - Are we creating new bias ?





TRANSPARENCY - DEALING WITH BIAS

• COMPAS – recidivism risk identification







DEALING WITH BIAS

- COMPAS recidivism risk identification
- Is the algorithm fair to all groups?







DEALING WITH BIAS

- COMPAS recidivism risk identification
- Is the algorithm fair to all groups?







CONFIRMATION BIAS

 tendency to search for, interpret, favour, and recall information in a way that confirms one's pre-existing beliefs or hypotheses.





Predictive policing

GUIDELINES TO DEVELOP ALGORITHMS RESPONSIBLY

- Who will be affected?
- What are the decision criteria we are optimising for?
- How are these criteria justified?
- Are these justifications acceptable in the context we are designing for?
- How are we training our algorithm?
 - Does training data resemble the context of use?



IEEE standard Algorithmic bias https://standards.ieee.org/project/7003.html

ART IS ABOUT BEING EXPLICIT

- Question your options and choices
- Motivate your choices
- Document your choices and options
- Regulation
 - External monitoring and control
 - Norms and institutions
- Engineering principles for policy
 - $\circ \quad Analyze-synthetize-evaluate\ -\ repeat$





https://medium.com/@virginiadignum/on-bias-black-boxesand-the-quest-for-transparency-in-artificial-intelligencebcde64f59f5b

ETHICS <u>BY</u> DESIGN – ETHICAL ARTIFICIAL AGENTS

• Can AI artefacts be build to be ethical?

- What does that mean?
- What is needed?
- Understanding ethics
- Using ethics
- Being ethical





NOT SO BRAVE NEW WORLD?



ETHICS BY DESIGN

1. Value alignment

- Identify relevant human values
- Are there universal human values?
- Who gets a say? Why these?



2. How to behave?

- Ethical theories: How to behave according to these values?
- How to prioritize those values?

3. How to implement?

- Role of user
- Role of society
- Role of AI system



PARTICIPATION – AI VALUES

UMEÅ UNIVERSITY

- Sources
 - Stakeholders: Designer, User, Owner, Manufacturer
 - Society: codes of ethics, codes & standards, law
- Identify what is
 - $\circ \quad \text{Socially accepted} \quad$
 - Morally acceptable
 - Legally possible



- Who decides who has a say?
- How to make choices and tradeoffs between conflicting values?
- How to verify whether the designed system embodies the intended values?





SOCIAL ACCEPTANCE – DEMOCRACY











- Binary choice
 - Brexit or Remain ?
- Information
 - "Are you for or against the European Union's Approval Act of the Association Agreement between the European Union and Ukraine?"
- Involvement
 - Colombia: city dwellers outvoted country side, where people had suffered by far the most from the FARC guerilla
- Legitimacy
 - Colombia: 50.2% No to 49.8% Yes, a difference of fewer than 54,000 votes out of almost 13 million cast
- Counting
 - Simple majority? Districts? Rankings?

IMPLEMENTATION: FROM VALUES TO FUNCTIONALITIES





ETHICAL REASONING? - AN EXAMPLE

- Design a self-driving can that makes ethical decisions
- Value: "human life"
- Implementation?
- Utilitarian car
 - The best for most; results matter
 - maximize lives
- Kantian car
 - $\circ \quad \text{Do no harm} \quad$
 - $_{\rm o}$ $\,$ do not take explicit action if that action causes harm
- Aristotelian car
 - Pure motives; motives matter
 - Harm the least; spare the least advantaged (pedestrians?)



Ethical theories

- Many different theories, each emphasizing different points
 - Utilitarian, Kantian, Virtues....
- Highly abstract
- None provide ways to resolve conflicts
- Deontology and Virtue Ethics focus on the individual decision makers while Teleology considers on all affected parties.



IMPLEMENTATION CHOICES



COMPUTATIONAL REQUIREMENTS

| Shared awarenessExplanationReal-time decision | Formal ethical rulesInstitutionsOffline reasoning |
|---------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------|
| collaboration | regulation |
| algorithmic | random |
| Formal ethical rules Ethical reasoning Real-time reasoning Learning ethics | • Trust ! |

IMPLEMENTATION ISSUES - DILEMMAS

- Moral dilemma
 - \circ You cannot have all
 - There is not one right solution!



- So the issue should not be what is the answer the AI system will give to a moral dilemma
- The issue is one of responsibility and openness
 - Make assumptions clear
 - Make options clear
 - Open data
 - Inspection
 - \circ Explanation
 - o ...





GOVERN AND VERIFY - GLASS BOXES

- Verify limits to action and decision
- Define the ethical borders
 - \circ Formal
 - Monitoring input output
- Governance
 - \circ input and output
 - Monitor
 - "block" undisarable





extending (TC King et al, AAMAS 2015)

TRANSPARENCY – GLASS DOMES





Terms and Conditions

1. Gauge, "Colors only the first well to one as as "as well-by" hand out in the colors and as a small of only in the Colorson of these the model on the first based of the decay of the only "based on the model on the first based of the decay of the only "based of the set of generic with a markey of galaxy of only."

1. Janes Conception Material and extension and extension protocologic and and Adalian any toward in Profession and Keng and an interfaent Assort Conception is producing a second and approximate for protocol or any interface and an any conception for the protocol or any interface and any conception for the the protocol or any interface and protocol prior for Conception and any comparison of protocol prior and protocols and by Conception and Protocols and any comparison of the Conception and protocols and protocols and protocols and by Conception

Astronomic, Property of Institute of editive systems: International International Contents on the International Contents and Property.

Used by Dispecting Provides over on the light in the providence of the providence of the second secon

 Elizzon Balle, Politika sector de rijk to construité der le obte au acturg ministrativitation and Construité la titul, dat sector qui es de acceleration de Palado, de forgeneral el dell'acturate.

* Supplicity and a second s

Description and the second seco

 Example. Designed due of charge shall be put in Mimetric strategy of models. J Publication provide: A second design of 1-20 per minist shall be assessed on all put do second.

 <u>Compared Manager</u>: Science (spin-tell art cleaner has be provided by plantas of parameter provided by the rest of the provided parameter of the plantas of the parameter of the parameter provided or the rest of the parameter of the parameter.

processory and an antience property space for a first property The Thiggs C, character can be a Character to a star of the property and the transmission of the star of the transmission of the transmissi

 Lans Gammanus Greener we creat the second loss in the second rest contents of the first means and the in the second rest. Contents and its and the internal second rest. Contents and its and the internal second rest. (cprop) (per second an er and a link in the indeparture, an left or the control of the second and the second of the two of the address the second of the second of the two of the second of the second second of the second of the second of the second second of the second of the second of the second second of the second of the second of the second second of the second of the second of the second second of the second of the second of the second second of the second of the second of the second second second second of the second of the second of the second second second second of the second of the second seco

 $\Phi_{\rm eff}$ and the second of the Ref. Physican stars in the observed second s

10. Management and the series of the state of the series and the series of the seri

into **Exception Comparation**. Follows not style, while stops in Contains Advanced of processing and the approxim in the local statement of the Contains Statement and interpreter statement and statement. Advanced and advanced to obtain the particular.

H. Gannal and Jan. The scenario did by generating the soft control and the based on hear d'Arque. Applications for the incident of an damas did by hole these d'Brages rate scenario Milliograd Comp.

Bookil Scoper of the prevention for the second secon

as appropriately provide the second strategies of the second second second second strategies and second second

 Associated. This happeness may not be associated within the map by and associate arring this manufacture of particulation.

in Automatical Content party integral to safe, a might an analysis of the same set of the same

 Security Assessed To a second relation over the strain production of the second second relation of the second secon















APPROACH



ETHICS FOR DESIGN(ERS)

- Regulation
- Certification
- Standards
- Conduct





ETHICS FOR DESIGN(ERS) – REGULATION, CONDUCT

- A code of conduct clarifies mission, values and principles, linking them with standards and regulations
 - \circ Compliance
 - Risk mitigation
 - \circ Marketing
- Many professional groups have regulations
 - Architects
 - Medicine / Pharmacy
 - Accountants
 - Military
- Is what happens when society relies on you!





ETHICALLY ALIGNED DESIGN

- identify and find broad consensus on pressing ethical and social issues and define recommendations regarding development and implementations of these technologies
- Standards
 - System design
 - Dealing with transparency
 - Dealing with privacy
 - Dealing with algorithmic bias
 - Data protection
 - Robotics
 - o ...
- Auditing
 - Certified agency





https://ethicsinaction.ieee.org/

AI ETHICS, AI FOR GOOD, AI FOR PEOPLE,...

- Harness the positive potential outcomes of AI in society, the economy
- Ensure inclusion, diversity, universal benefits
- Prioritize UN2020 Sustainable Development Goals
- The objective of the AI system is to maximize the realization of human values









http://www.ai4people.eu

AI AND EUROPE

- High level expert group on AI
 - Ethical guidelines
 - Policy and investment strategy
- AI Alliance: <u>https://ec.europa.eu/digital-single-market/en/european-ai-alliance</u>
 - forum engaged in a broad and open discussion of all aspects of Artificial Intelligence development and its impacts.
- AI4People Global Forum: <u>http://www.eismd.eu/ai4people/about/</u>
- CLAIRE (ELLIS): <u>https://claire-ai.org/</u>
 - Confederation of Laboratories for Artificial Intelligence Research in Europe
 - Research hubs ("CERN for AI")
- H2020
 - EU "AI-on-Demand Platform"
 - EU Flagship proposal "Humane AI" ...



UMEÅ UNIVERSITY

ETHICAL GUIDELINES

- IEEE Ethically Aligned Design
- EGE Statement on AI
- Asilomar principles
- EESC report on AI
-



Tim Dutton: <u>https://medium.com/politics-ai/an-overview</u> -of-national-ai-strategies-2a70ec6edfd Human dignity Privacy Security Responsibility Justice Solidarity Democracy

Rumman Chowdhury's list: <u>https://goo.gl/ca9YQV</u>

....

TAKE AWAY MESSAGE

- AI influences and is influenced by our social systems
- Design in never value-neutral
- Society shapes and is shaped by design
 - The AI systems we develop
 - The processes we follow
 - The institutions we establish
- Knowing ethics is not being ethical
 - $\circ\quad$ Not for us and not for machines
 - Different ethics different decisions
- Artificial Intelligence needs ART
 - Accountability, Responsibility, Transparency
 - Be explicit!
- AI systems are artefacts built by us for our own purposes
- We set the limits





RESPONSIBLE ARTIFICIAL INTELLIGENCE

