

Class Notes

Table of contents

- [Class Notes](#)
 - [Table of contents](#)
- [Notes 1](#)
 - [1. Course Title and Main Theme](#)
 - [Why "Ethics in AI"?](#)
 - [2. Course Structure – Topics to Be Covered](#)
 - [2.1. The Right Thing to Do](#)
 - [2.2. Why Ethics of AI?](#)
 - [2.3. Is Big Data Value Neutral? Ethics of Big Data](#)
 - [2.4. The Opacity of Algorithms: Fairness and Transparency](#)
 - [2.5. Responsibility and Explainability](#)
 - [2.6. Privacy and the Question of Data Ownership](#)
 - [2.7. Ethics and the Design of Social Media](#)
 - [2.8. Ethics of AI in Healthcare](#)
 - [2.9. Ethics of Robots](#)
 - [2.10. Ethics of Autonomous Systems \(Self-driving cars and Warfare\)](#)
 - [2.11. Embedding Ethics in AI](#)
 - [2.12. Designing Moral Machines](#)
 - [2.13. AI for Social Good](#)
 - [3. Evaluation System – Undergraduate \(UG\)](#)
 - [4. Evaluation System – Postgraduate \(PG\) and PhD](#)
 - [Notable Differences](#)
 - [5. General Rules](#)
 - [6. Synthesis and Reflection](#)
 - [7. Quotes and References Recap](#)
 - [8. Concluding Remarks](#)
- [Notes 2](#)
 - [1. Framing the Fundamental Questions](#)
 - [Analysis](#)
 - [2. The Nature of Moral Knowledge](#)
 - [Analysis](#)
 - [3. Universality vs. Relativism](#)
 - [Analysis](#)
 - [4. Enumerating the Sources of Moral Obligation](#)
 - [\(i\) Conforming to Social Norms and Behavior](#)
 - [Analysis](#)
 - [\(ii\) Conforming to Religious or Sect Norms](#)
 - [Analysis](#)
 - [\(iii\) Producing the Best Consequences](#)
 - [Analysis](#)
 - [\(iv\) Conforming to Norms of Reason](#)

- Analysis
 - (v) Actions That "Good People" Do
 - Analysis
 - (vi) Mutual Agreement, Promises, or Contracts
 - Analysis
 - (vii) Caring for Someone
 - Analysis
 - (viii) Sympathy/Empathy
 - Analysis
 - (ix) Acting in Self-Interest Without Harming Others
 - Analysis
 - 5. Bringing It All Together
 - 6. Conclusion: Reflective Moral Practice
 - Key Takeaways
- Notes 3
 - 1. Why Think of Ethics in AI?
 - 1.1 Challenging Common Assumptions
 - Analysis & Example
 - 2. Fundamental Questions in AI Ethics
 - 2.1 Intrinsic Moral Properties vs. Interactional Morality
 - Analysis & Example
 - 2.2 Agency, Autonomy, and Intelligence
 - Analysis
 - 3. Where Does the Question of Ethics Arise in AI?
 - 3.1 The Impact Question
 - Analysis & Example
 - 3.2 The Question of Knowing
 - Analysis & Example
 - 3.3 Is It the Machine or the Human?
 - Analysis
 - 3.4 Speed of Development
 - Example
 - 3.5 Superintelligence & the Problem of Control
 - Analysis
 - 3.6 Epistemic Reasons
 - Example
 - 3.7 Time
 - Analysis
 - 3.8 Nature of Ethics: Universal vs. Contextual
 - Example
 - 4. Ethical Challenges and Open Questions
 - 4.1 Universal Frameworks vs. Cultural Differences
 - 4.2 Mathematical Modeling
 - 4.3 Conceptual Discrepancies in Intelligence, Autonomy, Agency
 - 5. Bringing It All Together
 - Concluding Thoughts

- Notes 4
 - 1. Defining Big Data
 - Analysis
 - Example
 - 2. Assumptions Underlying Big Data
 - Analysis
 - 3. Big Data and the Limits of Knowing
 - 3.1 Probability vs. Explanation
 - Example
 - 3.2 The Role of Theory
 - Analysis
 - 3.3 Objectivity vs. Human Involvement
 - Example
 - 4. The Problem of Context
 - Analysis
 - 5. The Problem with Correlation
 - Analysis
 - Example
 - 6. Big Data and the Digital Divide
 - 6.1 Impact on Decision-Making
 - Real-World Example
 - 6.2 Salient Examples from the Text
 - 7. Broader Ethical Implications for AI
 - 8. Concluding Reflections
- Notes 5
 - 1. Why Think About Algorithmic Accountability?
 - 1.1 The Opaque Nature of Algorithmic Decisions
 - 1.2 Biased Data and Embedded Values
 - 1.3 The Need for Explicit Values
 - 2. The Rationale Behind Transparency
 - 2.1 Observability and Knowledge
 - 2.2 Transparency as Performative
 - 3. Three Forms of Opacity
 - 4. Defining Algorithmic Accountability
 - 4.1 Justification, Sanction, and Transparency
 - 5. What Kind of Justifications "Count"?
 - 5.1 Reasonable vs. Acceptable Justifications
 - 5.2 Universally Accepted Norms?
 - 6. Strategies for Algorithmic Accountability
 - 7. Is Opacity Always a Problem?
 - 7.1 The "Neural Net Produced It" Defense
 - 7.2 Contextual Goals vs. Outputs
 - 7.3 The Case Against Building the System
 - 8. Concluding Reflections
- Notes 6
 - 1. Why Think of Transparency?

- 1.1 Logic of Accumulation
 - 1.2 Performative Aspect of Transparency
- 2. Three Forms of Opacity
 - Example
- 3. Accountability in Social Contexts
 - 3.1 Social Embedding of AI
 - 3.2 Purpose and Impact
- 4. Kinds of Responsibility
 - 4.1 The Problem of Many Hands
 - 4.2 Oversight Mechanisms
- 5. Oversight as Operationalizing Accountability
 - 5.1 Evidence and Record-Keeping
 - 5.2 Contextual Norms
- 6. Algorithmic Accountability Defined (Binns)
 - 6.1 Justifications and Sanctions
 - 6.2 Answerability and Outcome Responsibility
- 7. What Kind of Justifications Count?
 - 7.1 Criteria for Valid Justifications
 - 7.2 Universally Accepted Norms
- 8. Conclusion: Tying It All Together

Notes 1

1. Course Title and Main Theme

The document is titled “Notes 1” and centers on **Ethics in AI**. It lists the broad themes, the evaluation system, and the general rules for the course. Each bullet point references a core area of study within AI ethics.

Why “Ethics in AI”?

Ethics in AI looks at how artificial intelligence technologies, algorithms, and data collection impact society in terms of fairness, justice, privacy, accountability, and a host of other moral and ethical considerations. The course covers **both theoretical frameworks** (such as understanding what constitutes “the right thing to do”) and **practical implications** (like designing fair, transparent, and beneficial AI systems). □cite□turn0file0□

2. Course Structure – Topics to Be Covered

The document explicitly lists the following bullet points as the core topics. Each one captures an essential aspect of AI ethics:

1. **The Right thing to do**
2. **Why Ethics of AI?**
3. **Is Big Data Value Neutral? Ethics of Big Data**
4. **The Opacity of Algorithms. Fairness and Transparency**
5. **Responsibility and Explainability**

6. **Privacy and the Question of Data Ownership**
7. **Ethics and the Design of Social Media**
8. **Ethics of AI in Healthcare**
9. **Ethics of Robots**
10. **Ethics of Autonomous Systems (Self-driving cars and Warfare)**
11. **Embedding Ethics in AI**
12. **Designing Moral Machines**
13. **AI for Social Good**

Below is a deep-dive into each of these focal areas.

2.1. The Right Thing to Do

- **Core Idea:** This introduces the fundamental philosophical question behind ethics: how do we determine “the right thing to do”? This question frames the rest of the course, as students need to learn not just technical details of AI but also the moral frameworks (e.g., utilitarianism, deontology, virtue ethics) that help us decide how AI should behave.
 - **Example:** A self-driving car faces a sudden dilemma: should it protect the occupant at all costs or minimize overall harm (e.g., potentially hitting fewer pedestrians)? The “right thing” might differ depending on the underlying ethical theory. □cite□turn0file0□
-

2.2. Why Ethics of AI?

- **Core Idea:** AI can greatly enhance human capabilities, but it also carries risks: bias, invasion of privacy, manipulation, and unintended societal impacts. The question “Why Ethics of AI?” addresses why we must go beyond mere programming and technical performance to analyze how AI aligns with ethical values.
 - **Quote from the Notes:** While not an explicit quote in the document, it repeatedly emphasizes the notion of ensuring “we do not plagiarise” or cause harm—this is part of a broader ethical approach that highlights responsibility.
 - **Reference:** “The Right thing to do,” from the bullet above, ties in seamlessly here. If we understand *why* ethics is essential, we can better pursue *what* is ethically correct in an AI context. □cite□turn0file0□
-

2.3. Is Big Data Value Neutral? Ethics of Big Data

- **Core Idea:** At first glance, “big data” might seem like an objective, neutral resource. But whenever data is collected, processed, or used, it can contain hidden biases and value judgments. The phrase “Is Big Data Value Neutral?” challenges the assumption that large-scale datasets are purely factual. Instead, it raises questions about **who collects the data**, **why** they collect it, and **how** it is being interpreted.
 - **Example:** A company that aggregates social media data to determine credit risk might inadvertently discriminate against certain demographics if the data (and the algorithms) reflect historical prejudices.
 - **Important Quote:** The document asks: “Is Big Data Value Neutral?” and references the “Ethics of Big Data,” pointing to the moral obligations in data use. □cite□turn0file0□
-

2.4. The Opacity of Algorithms: Fairness and Transparency

- **Core Idea:** Many AI algorithms, especially deep learning models, are “black boxes” whose decision-making processes can be difficult to interpret. Such opacity raises concerns about fairness—are the models discriminating based on race, gender, or other protected attributes?
 - **Transparency:** The course will discuss if and how to make these algorithms explainable and transparent. Transparency includes letting people know they are interacting with an AI system, clarifying why certain decisions are made, and showing the underlying logic or data used in the decision process.
 - **Real-World Example:** Credit-scoring algorithms that do not reveal why a certain user is denied credit, or hiring algorithms that rank candidates but never explain their rationale. Lack of transparency leads to challenges in detecting bias. □cite□turn0file0□
-

2.5. Responsibility and Explainability

- **Core Idea:** Closely related to fairness and transparency is the question of **who is responsible when AI goes wrong**. Is it the developer, the company that deploys it, or the AI itself (through some notion of artificial agency)?
 - **Explainability** is a step toward responsibility. If a system can explain its outputs in a human-understandable way, it becomes easier to hold the right parties accountable.
 - **Quote from the Document:** While not a direct quote, the topics clearly list “Responsibility and Explainability” as a dedicated bullet, suggesting a major component of the course. □cite□turn0file0□
-

2.6. Privacy and the Question of Data Ownership

- **Core Idea:** Modern AI systems rely heavily on user data. This section raises concerns about consent, surveillance, and data property rights. Who truly “owns” data once collected? Does a user have a right to have their data deleted?
 - **Examples:**
 - **Social Media:** Users uploading personal photos inadvertently granting usage rights to the platform.
 - **Healthcare:** Patients sharing medical records—should these be used for research without their explicit knowledge or only under strict anonymization protocols?
 - **Quote/Reference:** The course aims to unpack “Privacy and the Question of Data Ownership” because it is integral to building ethical AI that respects user autonomy. □cite□turn0file0□
-

2.7. Ethics and the Design of Social Media

- **Core Idea:** Social media platforms leverage AI to recommend content, moderate posts, and personalize user experiences. Ethical dilemmas arise around **filter bubbles**, **echo chambers**, **mental health implications**, and **manipulative design** (e.g., addictive features).
 - **Real-World Example:** Recommendation algorithms that only show users content they already agree with, leading to polarization and misinformation.
 - **Why It Matters:** Understanding how design choices in social media can propagate harmful social consequences is vital to designing more responsible systems. □cite□turn0file0□
-

2.8. Ethics of AI in Healthcare

- **Core Idea:** AI in healthcare can diagnose diseases, propose treatments, and manage patient data. With these benefits come ethical questions: do algorithms inadvertently discriminate? Are diagnoses transparent to doctors/patients? How are patient data and consent handled?
 - **Example:** An AI system recommending a certain cancer treatment, but not being transparent about the studies or the data behind its decision. This can impact patient trust and legal liability. □cite□turn0file0□
-

2.9. Ethics of Robots

- **Core Idea:** Robotic systems (e.g., humanoid robots, home assistants) raise questions of autonomy and moral standing. If a robot learns from its environment, to what extent can it be considered morally responsible for its actions?
 - **Discussion:** Topics might include the emotional bond humans form with robots, ethical constraints on how robots interact with vulnerable populations, and robots in hazardous industries.
 - **Quote/Reference:** The phrase “Ethics of Robots” signals that the course will tackle these fundamental concerns about the nature and rights (or non-rights) of machines. □cite□turn0file0□
-

2.10. Ethics of Autonomous Systems (Self-driving cars and Warfare)

- **Core Idea:** Autonomous systems operate with minimal human oversight. This raises extremely high-stakes ethical concerns:
 1. **Self-driving Cars:** Trolley-problem-style dilemmas in real traffic, liability questions, and the standards for safety.
 2. **Warfare:** Autonomous weapons deciding who to target. Is it ethically permissible to deploy lethal autonomous weapons without direct human control?
 - **Example:** Debates around the use of drones that can independently select targets. International bodies discuss whether to ban such weapons.
 - **Quote/Reference:** The bullet specifically mentions “Ethics of Autonomous Systems (Self-driving cars and Warfare).” □cite□turn0file0□
-

2.11. Embedding Ethics in AI

- **Core Idea:** How do we instill moral principles or constraints directly into AI systems? This might include value alignment techniques, rule-based restrictions, or robust auditing.
 - **Practical Angle:** Designing frameworks so that an AI’s objectives and behaviors match human ethical considerations—sometimes known as the “alignment problem.”
 - **Quote:** “Embedding Ethics in AI” is a recognized challenge: engineers often ask *how* to incorporate moral guidelines into code. □cite□turn0file0□
-

2.12. Designing Moral Machines

- **Core Idea:** A more direct extension of “embedding ethics.” If machines can act independently, how do we ensure they “choose” moral outcomes?
- **Contrast:** This goes beyond merely analyzing data ethically; it moves toward engineering machines that follow ethical imperatives even in unforeseen circumstances.

- **Example:** A “moral machine” might be a nursing robot that prioritizes patient well-being over cost-saving, or a self-driving car that respects all traffic laws and moral constraints.
 - **Quote/Reference:** The course bullet “Designing Moral Machines” is broad but central to AI ethics research. □cite□turn0file0□
-

2.13. AI for Social Good

- **Core Idea:** While many points address the pitfalls of AI, “AI for Social Good” highlights how AI can *positively* impact society: disaster response, medical breakthroughs, educational tools, climate change modeling, poverty alleviation, etc.
 - **Quote/Reference:** The bullet states “AI for Social Good,” suggesting an optimistic focus on leveraging AI ethically to bring tangible societal benefits. □cite□turn0file0□
-

3. Evaluation System – Undergraduate (UG)

The document specifies how UG students will be evaluated in this course. The breakdown is as follows:

1. **Individual Assignment:** 20 points
2. **Group Project:** 30 points
3. **End Sem (Final Exam):** 25 points
4. **Class Participation:** 10 points
5. **Response Paper (3):** 15 points

Here, “Response Paper (3) – 15” likely means students have to produce three separate response papers; the total of these is worth 15 points. It is not explicitly stated whether each paper is worth 5 points or if it is aggregated differently, but presumably each paper might carry equal weight.

Key Insight: The varied nature of evaluation—individual assignments, group projects, final exam, participation, and response papers—indicates that the course aims to engage students both in collective, collaborative thinking (group project) and personal reflection (individual and response papers).

□cite□turn0file0□

4. Evaluation System – Postgraduate (PG) and PhD

For PG and PhD students, the evaluation structure is slightly different:

1. **Individual Assignment (2):** 20 points
2. **Individual Presentation (2):** 20 points (best of 2 out of 3 presentations)
3. **End Sem (Final Exam):** 20 points
4. **Response Paper (3):** 15 points (best of 3 out of 4 papers)
5. **Individual Project:** 25 points

Notable Differences

- PG/PhD students have to do **two individual assignments** instead of just one.
- They deliver multiple presentations; only two best out of three are counted for the final grade.
- They have more response papers (four possible, best three counted).

- They have an **Individual Project** worth 25 points, distinct from the UG Group Project.

□cite□turn0file0□

Why the Different Structure?

Graduate-level courses often demand more in-depth individual research and presentation skills, reflecting a higher level of specialization and academic rigor. □cite□turn0file0□

5. General Rules

The document also lists general rules for the course:

1. No Plagiarism

- Direct statement: "Please do not plagiarise in the course as it will get you into trouble." □cite□turn0file0□
- *Analysis:* Academic integrity is crucial, especially in an ethics class.

2. Teaching Fellows (TF) and TAs

- "We will have a TF and TAs for the course whose help you can seek at any time." □cite□turn0file0□
- They have office hours where students can discuss problems or clarify doubts. If students need to speak to the main instructor, they should seek an appointment.

3. Deadlines and Extensions

- "All deadlines and assignments will be discussed and announced in advance. Please do not negotiate for an extension." □cite□turn0file0□
- *Interpretation:* The instructor sets firm deadlines to teach responsibility and time-management—key ethical values in academic work.

4. Mental Health and Stress

- "If at any time you are feeling stressed out feel free to reach out..." □cite□turn0file0□
- This underscores the importance of well-being and the open-door policy for students who may need emotional or academic support.

6. Synthesis and Reflection

Bringing it all together:

- **Broad Scope:** The range of topics—from "The Right thing to do" to "AI for Social Good"—reveals that the course aims to address both *conceptual/philosophical questions* and *practical/technical concerns* in AI ethics.
- **Hands-On Evaluation:** The evaluation structure (assignments, projects, papers) ensures students engage with real-world examples and develop an in-depth understanding, rather than only learning abstract concepts.
- **Rules Emphasize Integrity and Well-being:** The explicit mention of plagiarism, seeking help from TAs/TF, and addressing stress points to a supportive environment where ethical conduct is expected not

just in AI but also in the students' academic work.

The note also reminds students that everything from data collection to algorithm design has ethical implications and that these are not optional considerations—rather, they define the trustworthiness and societal impact of AI systems. [cite](#) [turn0file0](#)

7. Quotes and References Recap

While the PDF itself is brief and mostly in bullet points, there are key references we can highlight as quotes or direct paraphrases:

1. "The Right thing to do" – frames the philosophical question at the heart of ethics.
2. "Is Big Data Value Neutral? Ethics of Big Data" – challenges assumptions of neutrality.
3. "The opacity of Algorithms. Fairness and Transparency" – underscores the black-box problem.
4. "Privacy and the Question of Data Ownership" – addresses who truly has rights over data.
5. "Please do not plagiarise in the course as it will get you into trouble." – emphasizes academic integrity.
6. "We will have a TF and TAs for the course whose help you can seek at any time." – highlights available support.
7. "If at any time you are feeling stressed out feel free to reach out..." – fosters an environment of open communication and support.

Each of these points is crucial to understanding the broader mission of the course: to ensure that students become ethically aware AI practitioners (or researchers) who can identify and mitigate potential harms.

8. Concluding Remarks

"Notes 1" sets the stage for a comprehensive journey through AI ethics. It demonstrates that:

- Ethical deliberation is not a side topic but central to responsible AI design and deployment.
- Students will be evaluated through a variety of assignments aimed at ensuring deep engagement with ethical, technical, and societal dimensions of AI.
- The course environment prioritizes **integrity, transparency, responsibility, and student well-being**—values that mirror the ethical principles the curriculum aims to teach.

No part of this document is superfluous: each bullet in the "Topics to be covered" is an essential puzzle piece in the broader conversation about how AI can and should serve the greater good while minimizing unintended negative consequences. [cite](#) [turn0file0](#)

Notes 2

1. Framing the Fundamental Questions

The notes begin by asking:

"How do we decide what is the right thing to do? What are the sources of our obligations? How do we know what is the right thing to do?" [cite](#) [turn0file0](#)

This cluster of questions underscores the complexity of moral decision-making. The text immediately situates ethics as a domain of inquiry into “obligations” and the methods by which we identify, recognize, or justify them. Instead of assuming a universal, one-size-fits-all answer, the notes highlight that moral action can be influenced by instinct, reason, training, or social context.

Analysis

- **Obligations vs. Preferences:** The word “obligation” typically refers to moral duties that bind us, as opposed to personal preferences (e.g., “I like chocolate ice cream” is a preference, whereas “I am obligated to be honest” is a moral imperative). This distinction is crucial to clarify what kind of “right thing” we are talking about—something that we owe to ourselves, to others, or to society.
- **Role of Education or Training:** The text raises the question: “How do we train our moral intuitions to act in the right direction?” [cite](#) [turn0file0](#). This implies that morality is not always an inborn capacity but may require cultivation—through education, reflective practice, or critical thinking.

A concrete example to illustrate this might be a child learning about telling the truth. At first, they might not fully understand why lying is wrong, but through consistent moral training (parental guidance, religious instruction, cultural norms), they develop an intuitive aversion to lying. Over time, this becomes “the right thing” to do in their worldview.

2. The Nature of Moral Knowledge

Next, the notes pose a series of questions:

“Do you instinctively know the right from the wrong? Or do you reason the right from the wrong? Are all moral actions about instincts or reasons?” [cite](#) [turn0file0](#)

And further:

“Is moral thinking and action a matter of training? How then one should be trained? What sort of reflection is required for moral training? In what direction should one be thinking?” [cite](#) [turn0file0](#)

Analysis

- **Instinct vs. Reason:** This dilemma mirrors an age-old philosophical debate: Are we naturally inclined toward certain moral truths (perhaps guided by empathy or innate moral feelings), or must we rely on rational argumentation to distinguish right from wrong? Thinkers like David Hume emphasized sentiment (instinct/feeling), while Immanuel Kant emphasized reason. The notes invite us to see that either extreme—pure emotion or pure rationality—may be incomplete.
- **Practical Training:** The questions about moral education or “training” address how we *develop* these intuitions and reasoning capacities. One might train moral judgment through:
 - **Case studies** (examining moral dilemmas),
 - **Role-modeling** (observing the behavior of admired individuals),
 - **Reflection** (journaling, meditation, philosophical study).

An illustrative example can be found in professional ethics training (e.g., medical ethics). Students in medical school do not rely solely on instincts. They also learn frameworks like the Hippocratic oath, principles of non-maleficence (“do no harm”), beneficence (promoting the patient’s best interests), and autonomy (respecting patient choice). This mixture of reason, tradition, and empathy is a form of moral training.

3. Universality vs. Relativism

“What are the ways that we can know what is the right? Is it same for everyone? Is it different from one individual to the other? Is it dependent on the context? Is it dependent on the culture? Is it different for the west and the east?” □cite□turn0file0□

Here, the text introduces the debate over universal moral principles versus moral relativism. It questions whether moral truths or obligations might vary across cultures, societies, or even individuals.

Analysis

- **Universalism:** Some moral theories (e.g., Kantian ethics, various religious traditions) hold that moral truths apply universally—regardless of culture, context, or personal preference. According to this view, certain actions are always right or always wrong.
- **Relativism:** On the other side, moral relativism suggests that standards of right and wrong depend on cultural norms, societal pressures, or personal contexts. An action might be morally acceptable in one culture but not in another.
- **Contextual Nuances:** The notes’ question “Is it dependent on the circumstances one is exposed to?” □cite□turn0file0□ acknowledges that even if some moral principles seem universal, their *application* can vary widely due to cultural conditions or different life circumstances.

An everyday example: Attitudes toward social norms around dress codes or dietary restrictions might differ. A specific act—like eating pork—could be morally neutral in one culture while being strongly frowned upon in another, for either religious or cultural reasons. Thus, the “rightness” of that action is influenced by context.

4. Enumerating the Sources of Moral Obligation

The text offers nine specific “sources of moral obligation” □cite□turn0file0□. These are not necessarily exhaustive, but each one captures a significant moral motivation that could drive our actions.

(i) Conforming to Social Norms and Behavior

“(i) Conforming to social norms and behaviour (Deviance may be a costly affair/ we are trained to act in certain ways so do act out of habit etc.)” □cite□turn0file0□

Analysis

1. **Social Pressure and Habits:** Often, people act morally (or at least in line with certain norms) because violating these norms leads to punishment, ostracism, or disapproval. We might hold a door open for someone because society teaches us this is polite.
2. **Cost of Deviance:** If you choose not to follow norms—say you lie repeatedly or engage in theft—you risk legal repercussions or social stigma. Over time, many of these norms become ingrained habits, making conformity feel like the “natural” choice.

An example is the practice of queuing in public spaces. People wait their turn largely because society frowns upon cutting in line. Over time, this social norm is internalized to the point that it feels morally wrong to skip ahead.

(ii) Conforming to Religious or Sect Norms

“(ii) Conform to certain other norms (religious/ that of a sect/ creed/ caste etc.) ... there might be sects which one joins voluntarily while many acts done within the social realm may not be a product of voluntary membership.” □cite□turn0file0□

Analysis

1. **Voluntary vs. Involuntary Membership:** Religion or sectarian affiliation can be a source of moral rules—dietary laws, worship obligations, charitable giving—that one follows either by birth (involuntary) or by conversion (voluntary).
2. **Overlap with Social Norms:** Religious norms often overlap with broader social norms but can be stricter or differ in specifics. For instance, dietary restrictions during Lent in some Christian traditions or the avoidance of certain foods in other religions.

A real-life example might be a person fasting during Ramadan. They could do so out of personal religious conviction, communal tradition, or both. The obligation is partly internal (faith) and partly social (family and community expect it).

(iii) Producing the Best Consequences

“(iii) Those are the ones that produce the best consequences.” □cite□turn0file0□

Analysis

1. **Consequentialism:** This point directly refers to moral theories like utilitarianism, which argue that the right action is the one that yields the greatest good for the greatest number.
2. **Practical Assessment:** Acting on this principle involves evaluating outcomes and choosing the path that maximizes overall well-being or minimizes harm.

An example: If you have to decide how to allocate a limited budget in a public health system, a consequentialist approach would attempt to save the greatest number of lives or maximize health benefits for the population.

(iv) Conforming to Norms of Reason

“(iv) They conform to norms of reason.” □cite□turn0file0□

Analysis

1. **Rationalist Traditions:** This connects to philosophers (like Kant) who argue that moral duty is grounded in rational consistency. For instance, you do not lie because lying cannot be universalized without contradiction.
2. **Consistency & Universality:** The phrase “norms of reason” implies acting on principles you can logically will for everyone—creating a moral law that is consistent, not contradictory.

An example might be refusing to break a promise because you realize that if *everyone* broke promises, the concept of promise itself would become meaningless.

(v) Actions That "Good People" Do

"(v) These are the actions that good people do. (we conform to certain standards of goodness)"

□ cite □ turn0file0 □

Analysis

1. **Virtue Ethics:** This idea resonates strongly with virtue ethics, which focuses on the character of the moral agent. Here, the question is less "What should I do?" and more "What kind of person should I be?"
2. **Imitation of Role Models:** We look at individuals we regard as moral exemplars—saints, heroes, mentors—and strive to do what they would do.

For instance, if we admire a humanitarian like Mother Teresa, we might volunteer at shelters or donate to charitable causes because "that's what good, compassionate people do."

(vi) Mutual Agreement, Promises, or Contracts

"(vi) Because we mutually agreed to act in certain ways (promises/contracts)" □ cite □ turn0file0 □

Analysis

1. **Social Contract Theory:** Philosophers like Thomas Hobbes or John Locke proposed that moral and political obligations arise from a (real or hypothetical) contract that people make to escape a "state of nature."
2. **Interpersonal Reliability:** On a personal scale, this also applies to everyday agreements: "I promised I would help you move your furniture on Saturday, so I'm obligated to do so."

A common example is signing a lease agreement: both tenant and landlord promise certain behaviors (paying rent on time, providing a livable space). Morally, one feels an obligation to honor that agreement because it was freely entered.

(vii) Caring for Someone

"(vii) Because we care for someone" □ cite □ turn0file0 □

Analysis

1. **Ethics of Care:** This taps into moral theories emphasizing relationships, empathy, and the emotional bonds we form with others (often associated with feminist ethics).
2. **Personal Attachment:** Unlike contractual or universal principles, caring for someone suggests a personal, emotional commitment. When you look after an elderly parent, you do so not because it's necessarily the best universal outcome or an explicit promise, but because you love them.

A day-to-day example would be cooking a meal for a friend who is sick. You do it because you care, which is reason enough to feel morally "obligated" to help them.

(viii) Sympathy/Empathy

"(viii) Because we feel sympathetic/empathetic towards them." □cite□turn0file0□

Analysis

1. **Emotional Basis:** Closely related to the previous point, this source of moral action highlights the power of empathy—feeling another's pain or situation as if it were your own.
2. **Immediate Response:** People often donate to disaster relief after seeing moving images or hearing firsthand accounts of suffering. The impetus is empathy, which can be as strong (or stronger) than rational deliberation.

An example: You see a stray animal injured on the street. Empathy compels you to rescue it or bring it to a vet, even if no contract or explicit rule requires you to do so.

(ix) Acting in Self-Interest Without Harming Others

"(ix) We decide to act in ways that benefit ourselves without harming anyone." □cite□turn0file0□

Analysis

1. **Ethical Egoism (Tempered):** This suggests a version of moral motivation where self-interest is central, but we limit our actions so as not to harm others. It's not purely selfish: it recognizes moral boundaries that keep our self-interest from infringing on others.
2. **Win-Win Situations:** Many routine decisions (like choosing a career path or investing in personal development) fall here. You do something that helps you personally—studying, exercising, building a business—while ensuring it does no harm.

For instance, an entrepreneur might start a company to make a profit (self-interest) but also ensures fair treatment of workers (avoiding harm). The ethical orientation is primarily inward (personal gain), but it's bounded by moral consideration for others' well-being.

5. Bringing It All Together

Collectively, these sources of obligation capture the richness and variety of moral motivations:

- **Social/Cultural Norms (i, ii):** External pressures or teachings shape our sense of right and wrong.
- **Outcome-Oriented (iii):** Evaluating the consequences (best or worst) of actions.
- **Rational Principles (iv):** Acting according to logical consistency or universalizable norms.
- **Virtue/Exemplar (v):** Modeling ourselves on those we consider morally praiseworthy.
- **Agreements (vi):** Fulfilling promises or contracts because we gave our word.
- **Emotional Bonds (vii, viii):** Caring relationships and empathy as drivers of moral action.
- **Respectful Self-Interest (ix):** Pursuing personal benefit but not at the expense of harming others.

The notes remind us that moral decisions may draw on *multiple* sources at once. A person might feed the homeless partly because society applauds charity (i), partly because their religious faith recommends helping the needy (ii), partly because it produces good consequences (iii), and partly because they personally empathize with those in need (viii).

6. Conclusion: Reflective Moral Practice

The driving theme in these notes is **reflective moral practice**. Questions such as “How do we know what is right?” and “In what direction should one be thinking?” (□cite□turn0file0□) push us toward a lifelong process of questioning, analyzing, and refining our moral intuitions. Rather than offering a single prescriptive doctrine, the text outlines *multiple* points of reference—social norms, religion, consequences, reason, virtue, promises, care, empathy, and self-interest (tempered by harm avoidance).

By encompassing these different angles, the notes show that morality is:

- Not *only* about following rules,
- Not *only* about achieving good outcomes,
- Not *only* about caring for others or upholding reason,
- But a dynamic interplay among all these factors.

A final example to unify everything: Imagine volunteering in your local community. You might do so **(i)** out of conformity to a social ideal that “good citizens volunteer,” **(iii)** because helping yields positive outcomes, **(v)** because you emulate role models you consider morally good, and **(vii)/(viii)** because you genuinely care or feel empathy for people in need. Your action is multi-motivated; it draws from overlapping moral commitments.

Key Takeaways

- **Multiplicity of Moral Sources:** There isn’t just one reason why people act morally; it can stem from social, cultural, rational, emotional, or contractual grounds.
- **Training & Reflection:** Moral intuitions and reasoning can be developed. We learn to refine our instincts and apply reason or empathy more consistently.
- **Contextual/Universal Tensions:** The text invites us to consider whether moral truths are universal or context-dependent, highlighting the complexity of real-world ethics.
- **Combination in Practice:** In most real situations, multiple sources of obligation combine to form the mosaic of our moral actions.

In short, “Notes 2” provides a panoramic view of the factors driving moral decision-making. It challenges us to reflect on which combination of factors influences us personally and how we might responsibly cultivate our moral agency in an ever-changing social and cultural environment. □cite□turn0file0□

Notes 3

1. Why Think of Ethics in AI?

The text opens with a crucial question:

“Why think of ethics in AI?” □cite□turn1file0□

This question may appear straightforward, yet it invites us to look closer at AI’s profound impact on human life. AI is no longer an abstract, futuristic technology; it is embedded in everyday decisions—from social media

feeds to credit scoring and healthcare diagnostics. The notes warn that it's naïve to assume AI will inherently be "good" and that ignoring ethical implications might lead to unforeseen harms.

1.1 Challenging Common Assumptions

The notes list four assumptions often made about AI:

1. **"AI is automatically going to be ethical."** (i)
2. **"AI is based on principles of reason so it will be ethical."** (ii)
3. **"AI is never going to be that intelligent to pose ethical challenges."** (iii)
4. **"AI is more objective than humans so it does not require ethics."** (iv)

Each assumption suggests a stance that effectively *disengages* human responsibility. For instance, believing AI is "automatically ethical" might lead engineers or policymakers to pay less attention to potential bias in data or to the exploitative ways a system could be used. Likewise, attributing perfect "objectivity" to AI overlooks how data—collected, processed, and trained upon—often carries embedded human biases.

Analysis & Example

- **Embedded Values:** Even a simple recommendation algorithm for music streaming can favor certain genres or artists, reflecting hidden assumptions about what "good" music is. This situation demonstrates that "objectivity" can be an illusion if the underlying data or model design is skewed.
- **Developmental Complexity:** AI's complexity can surpass the immediate comprehension of its designers. This calls into question the assumption that "it's never going to be *that* intelligent," because modern AI systems (like large language models or advanced reinforcement learning) can behave in unexpected ways.

2. Fundamental Questions in AI Ethics

The notes present a set of foundational inquiries about moral properties and human-machine comparisons:

"The problem of intrinsic moral properties: Does AI have an intrinsic moral property? Is an intrinsic property required for an agent to be ethical?" □cite□turn1file0□

2.1 Intrinsic Moral Properties vs. Interactional Morality

One question is whether an AI system must *itself* possess moral qualities or if ethical considerations arise purely because of how it interacts with us. In other words:

- If an entity lacks emotions, consciousness, or a moral sense, can it still be bound by ethical constraints?
- Or does the entire question of ethics simply emerge when an AI's actions affect human well-being?

Analysis & Example

Suppose an AI system is used in medical diagnoses. Even if the AI has no "intrinsic" moral sense, doctors and hospital administrators have ethical concerns about how it might misdiagnose or prioritize certain patients over others. The moral conversation here focuses on the interaction—patient outcomes, fairness, and accountability—rather than the AI's interior moral state.

2.2 Agency, Autonomy, and Intelligence

"Is agency and autonomy of machines the same as that of humans? Are we using the same concepts to define humans and machines?" □cite□turn1file0□

This part of the notes questions whether philosophers and computer scientists define terms like *intelligence* and *autonomy* in the same way. Philosophers might link autonomy with the capacity for free will or rational reflection; computer scientists might measure autonomy by an AI's ability to operate without human intervention.

Analysis

- **Philosophical Autonomy:** Often tied to free will, moral responsibility, or consciousness.
- **Technical Autonomy:** Tied to self-sufficiency in performing tasks without direct oversight.

A simple example is a self-driving car. It displays *technical* autonomy by navigating roads. But does it have *philosophical* autonomy? Probably not in any robust sense. This dissonance in usage can complicate ethical debates.

3. Where Does the Question of Ethics Arise in AI?

The notes extensively detail contexts in which ethical issues emerge:

"The question of Impact: Is it because AI has impact on humans?" □cite□turn1file0□

3.1 The Impact Question

AI's ability to affect large segments of society—decisions about insurance coverage, job applications, policing, war, or health—forces ethical examination. If an AI system denies someone a job opportunity based on a spurious correlation, the injustice stems from that system's *impact*. Likewise, in warfare, using autonomous drones raises questions of accountability and the moral calculus of using lethal force without a human "in the loop."

Analysis & Example

- **Predictive Policing:** In some cities, AI predicts high-crime areas. This can lead to biased policing if the training data reflect historical biases. The "impact" is direct—targeted communities might face disproportionate surveillance.

3.2 The Question of Knowing

"Do we know how the machine makes the decision? Can we predict the decision? Can we explain the decision?" □cite□turn1file0□

Ethical challenges are compounded by the "black box" nature of many AI models. If even the system's creators struggle to interpret how it arrives at certain outputs, transparency and accountability suffer.

Analysis & Example

- **Explainable AI:** There is a growing field dedicated to making AI decisions interpretable. A medical AI system might generate a conclusion about a patient's risk of developing a disease. If the patient or doctor can't understand *how* the system arrived at that conclusion, can they ethically trust it?
-

3.3 Is It the Machine or the Human?

"Does ethics arise because of the use the machine is put to? Or is it because who puts it to use?"

□cite□turn1file0□

This poses a fundamental question: is the ethical dilemma located in the *technology* itself or in the *human context of its usage*? The text further asks:

"Is the question: what end the machine is serving or whose end the machine is serving?"

□cite□turn1file0□

Analysis

- **Ends vs. Means:** This resonates with Immanuel Kant's moral principle: "Treat humanity...never merely as a means to an end, but always at the same time as an end." If humans use AI purely as a means to exploit others (e.g., invasive data gathering), the moral failing may rest on those human intentions.
- **Ethical Distribution of Benefits and Burdens:** The text asks how benefits and burdens from AI are shared in society. An AI that automates tasks might create profits for some while displacing workers. This raises ethical questions about wealth inequality, responsibility, and compensation.

3.4 Speed of Development

"AI develops faster, ethics always trails. Our ability to create, innovate, and process data has outstripped our control of Data." □cite□turn1file0□

This point captures the *tech-ethics lag*: new technologies often outpace regulatory frameworks and ethical guidelines. Innovations appear so quickly that society struggles to shape them responsibly.

Example

Social media platforms introduced deepfake technology before any robust ethical or regulatory consensus formed. Consequently, deepfakes spread misinformation, and only after their proliferation did governments and institutions scramble to address the ethical implications.

3.5 Superintelligence & the Problem of Control

"Is it because of superintelligence and the problem of control: General intelligent machines could be faster, maybe able to replicate, may have values where it wants more copies of its own self."

□cite□turn1file0□

Though it may sound futuristic, the possibility of AI surpassing human intelligence (artificial general intelligence, or AGI) raises existential risk concerns: how do we ensure it aligns with human values if it becomes more capable than us? The text hypothesizes scenarios where an advanced AI might replicate or have goals contrary to human well-being.

Analysis

- **Value Alignment Problem:** Researchers discuss how to “teach” advanced AI to share human values. If it optimizes for a misaligned goal, the result could be catastrophic (a classic example is the “paperclip maximizer,” an AI whose single goal is making paperclips, inadvertently wreaking havoc on humanity to accomplish this).
 - **Control Dilemma:** Once an AI becomes too advanced, even its creators might struggle to impose constraints. Hence, the question of controlling superintelligence is not merely a technical puzzle but also an urgent ethical one.
-

3.6 Epistemic Reasons

“AI is altering the way we interpret and interact with the environment and how we know the world.”

□cite□turn1file0□

The text also notes that AI changes the *epistemic* foundations—our ways of knowing. It challenges traditional scientific methods by introducing massive data analytics, predictive models, and sometimes non-intuitive correlations that upend conventional theories.

Example

Consider climate modeling. AI can process huge datasets and find complex patterns that humans could never see manually. If these models drive policy decisions, we must wrestle with how to interpret and trust emergent “knowledge” that lacks a straightforward chain of human reasoning. The ethics come into play when deciding how to use these AI-derived insights—especially if they remain opaque to human experts.

3.7 Time

“The direction and the development of AI is unpredictable. Can we use the ethical values of now to predict the use of machines in the future?” □cite□turn1file0□

We face temporal challenges: today’s moral frameworks might not remain suitable as technology evolves. The text suggests a tension between the rapidly changing technological landscape and moral theories that might need constant recalibration.

Analysis

- **Ethical Flexibility:** Philosophical systems often rely on stable principles (e.g., utilitarianism’s “greatest good” or Kant’s categorical imperative). If the context changes drastically (e.g., AI drastically shifts labor markets or redefines intelligence), we may need to adapt or reinterpret these principles.
 - **Resource Allocation:** The text hints: “Would we have the resources to cater to the needs of the development of AI?” (□cite□turn1file0□). Building robust ethical oversight infrastructure might require significant funding, global cooperation, and time—none of which is guaranteed.
-

3.8 Nature of Ethics: Universal vs. Contextual

"Is ethics universal? Contextual? Relative? Would the problems of AI in India be different from other countries?" □cite□turn1file0□

By raising this question, the notes draw attention to cultural and societal differences. Ethical frameworks that work in one context (e.g., Western liberal democracies) might not translate seamlessly elsewhere. Additionally, the text observes that our ethical intuitions and theories have developed over centuries of human-to-human interaction—and that might shift dramatically if AI surpasses human intelligence.

Example

- **Facial Recognition:** Some societies might be more tolerant of widespread surveillance to maintain public order. Others find it a violation of privacy rights. The same AI tool can spark different ethical dilemmas depending on cultural attitudes toward privacy, security, and individual liberty.

4. Ethical Challenges and Open Questions

Finally, the notes list explicit challenges:

"How do we model ethics? Should we use universal frameworks?" □cite□turn1file0□

Modeling Ethics: There is an ongoing debate on whether we should attempt to encode moral principles into AI (top-down approach) or if AI should "learn" ethical behavior by observing social norms (bottom-up approach). The text hints at potential pitfalls in either route.

4.1 Universal Frameworks vs. Cultural Differences

"Societal/cultural frameworks (Is moral machine experiment correct?)" □cite□turn1file0□

The "Moral Machine" experiment, popularized by MIT, asked people worldwide how a self-driving car should respond in life-and-death traffic scenarios. Responses varied significantly across cultures, suggesting that imposing a single, universal set of moral rules might alienate or misrepresent some societies.

4.2 Mathematical Modeling

"What is the way to model ethics mathematically? Are ethical theories amenable to mathematical modelling?" □cite□turn1file0□

This is a key philosophical and technical question. Some frameworks, like utilitarianism (maximizing overall well-being), might appear more straightforward to translate into an algorithm than virtue ethics or deontological principles, which revolve around character and duties. But even utilitarian calculation can become intractably complex in real-world scenarios—who defines what "well-being" means, and how do we quantify it?

4.3 Conceptual Discrepancies in Intelligence, Autonomy, Agency

"The concepts of intelligence, autonomy, and agency used and understood by AI may be completely different than one used in Ethics." □cite□turn1file0□

AI researchers may treat "intelligence" as pattern recognition, problem-solving, or learning capability. Ethical frameworks may treat "intelligence" as the capacity to understand moral principles and reflect upon them.

This conceptual mismatch can hinder conversations about AI's moral responsibilities or rights.

5. Bringing It All Together

From the text, we can see that *ethics of AI* is not simply about adding a final “safety net” to an otherwise neutral technology. Instead, ethics weaves through every stage—from AI's inception to its deployment, from its cultural context to its potential future developments. As the notes emphasize:

1. Common Assumptions Must Be Questioned

Believing AI is inherently ethical or purely objective overlooks how deeply human biases and intentions shape these systems.

2. Defining Moral Agency is Complex

Deciding whether AI can be said to have *agency* or *moral standing* involves comparing philosophical and technical conceptions.

3. Impacts on Real People

Ethics emerges most tangibly when AI's decisions affect individuals' livelihoods, freedoms, and well-being—creating a clear impetus to question fairness, accountability, and transparency.

4. Future Gazing and Superintelligence

We must grapple with hypothetical but potentially monumental scenarios where AI may exceed human capabilities—and how we'd control or align it with human values.

5. Cultural Variations and Evolving Norms

Ethics in AI is not a purely universal puzzle. Cultural contexts matter. Moreover, technology evolves so quickly that ethical standards need constant re-evaluation.

Concluding Thoughts

“Notes 3” compels us to examine the multi-faceted nature of AI ethics:

- **Moral Foundations:** Intrinsic properties vs. interactional ethics.
- **Human vs. Machine Responsibility:** Tools reflect the values and uses imposed by their creators and operators.
- **Speed and Scale:** The rapid development of AI can outstrip ethical guidelines, leading to reactive rather than proactive moral oversight.
- **Global and Cultural Dimensions:** Ethics cannot remain a siloed conversation; it depends on societal contexts and normative frameworks.
- **Future Directions:** Issues like superintelligence, accountability, and the changing epistemic landscape underscore that AI ethics is not static; it must evolve alongside the technology.

In short, these notes make clear that ethics is *central*, not *peripheral*, to AI. They demand that we ask hard questions about how we design, deploy, and ultimately live with increasingly intelligent, autonomous systems.

□cite□turn1file0□

Notes 4

1. Defining Big Data

The text provides multiple definitions and dimensions of “big data,” indicating that size is only part of the story:

“‘Big data’ can be defined as research that represents a step change in the scale and scope of knowledge about a given phenomenon” (Schroeder and Cows, 2014). □cite□turn2file0□

“It is about a capacity to search, aggregate, and cross-reference large data sets.” (Boyd and Crawford, 2012; 663). □cite□turn2file0□

Analysis

- **Scale vs. Capability:** A core characteristic of big data is its *capacity*—the ability to process massive sets quickly, correlate them, and discover patterns. This goes beyond mere volume. It implies a paradigm shift in how we approach empirical research.
- **Pattern Recognition:** The notes emphasize that big data “allows for pattern recognition or analysis across different data sets” □cite□turn2file0□. This means we can find relationships that might be unobservable with smaller, more traditional datasets.

Example

A retail giant might track billions of shopping transactions and correlate them with weather patterns, social media sentiments, and online browsing behaviors. Identifying correlations (e.g., a spike in hot beverage purchases when certain hashtags trend) can be profitable, but the process can also introduce biases if not contextualized properly.

2. Assumptions Underlying Big Data

The notes highlight a set of assumptions that often accompany big data analytics:

“Data exists out there and it exists prior to the investigation, exists for the object under study, and exists in an atomised or divisible form that allows for collection.” □cite□turn2file0□

Analysis

1. **Pre-Existing Data:** The assumption is that data is “out there,” waiting to be collected. This ignores how data generation is often shaped by social, political, or economic processes (e.g., who has internet access, who is more likely to fill out surveys).
2. **Atomization:** Treating data as a set of discrete points makes it easier to store and analyze, but it risks stripping away context.

Think of social media posts: they are captured as text strings, hashtags, metadata, etc. But the context (the user’s mood, the cultural moment, possible sarcasm) can be lost in translation.

3. Big Data and the Limits of Knowing

"Big data represents a challenge to how we know – tends more towards probability and prediction rather than causality and explanation." □cite□turn2file0□

This line underscores one of the biggest philosophical shifts in big data usage: a focus on high-level correlations at the expense of in-depth causal understanding.

3.1 Probability vs. Explanation

- **Predictive Power:** Many big data practitioners argue, "It works, so why should we bother?" regarding causation. If you can predict an outcome with decent accuracy, do you *need* to know the *why*?
- **Ethical Trade-Off:** Foregoing causal understanding can be ethically dangerous. For instance, if a predictive model identifies certain neighborhoods as "high risk" for insurance, it might reflect systemic biases or historical discrimination without uncovering the root causes that lead to higher claims.

Example

In credit scoring, a machine learning model might simply identify a correlation between late-night browsing and loan defaults. The model "works"—it might accurately predict who will default—but it can't explain why. This lack of explanation can be ethically fraught if it inadvertently penalizes people with limited internet access or unusual work schedules.

3.2 The Role of Theory

"Another assumption that governs big data is that we do not need theory to understand – patterns are sufficient." □cite□turn2file0□

However, the text points out that classification and target variable selection often require human judgment and a theoretical lens:

"...this is a subjective process – data mining can only sort out problems that can lead to formalisation—sometimes there is a need to create new classes and this requires an employment of judgement..."
□cite□turn2file0□

Analysis

- **Subjectivity in Data Work:** People still choose what variables matter (e.g., "creditworthiness," "good employee"). These categories are not purely objective; they reflect human values and biases.
- **Danger of Spurious Correlations:** Without theoretical grounding, big data can produce "surprising" but meaningless patterns—like ice cream sales predicting stock prices. They might correlate but have no causal link.

3.3 Objectivity vs. Human Involvement

"Big data is objective as it eliminates the human aspect—both design and selection are as much part of interpretation and involve a theory." □cite□turn2file0□

This statement captures the *myth* of big data's objectivity. The text clarifies that human decisions pervade every step, from data collection methods to which results are accepted or discarded.

Example

A facial recognition system might claim high accuracy, yet its underlying training data could exclude certain ethnic groups, leading to disproportionate errors on those faces. The so-called “objective” model emerges from subjective human choices about data collection, labeling, and evaluation metrics.

4. The Problem of Context

“Big data is devoid of the context in which a certain data is generated. In the quest for ‘bigness’ what is lost is the specificity of the context.” □cite□turn2file0□

Analysis

- **Contextual Nuance:** Reducing social media posts to discrete data points might ignore the local slang, socio-political climate, or personal histories. The text notes the assumption “that a data generated in a certain context (say a tweet) represents accurately the sentiment of the individual” □cite□turn2file0□. In reality, that tweet could be satire, or the user could be joking.
 - **Prescriptive Uses:** Targeted advertising or political messaging often treat data points as direct representations of user preferences. This can lead to manipulative practices, especially if the “context” is absent (e.g., circumstances under which a user posted certain content).
-

5. The Problem with Correlation

“When correlation displaces causality or explanation, ... a particular combination of eating habits, weather patterns, and geographic location correlates with a tendency to perform poorly in a particular job or susceptibility to a chronic illness...” (Andrejevic 2014: 1681). □cite□turn2file0□

Analysis

- **Unintuitive Pairings:** Big data can reveal bizarre correlations (e.g., types of browser usage predicting job performance). Yet these insights may rest on tenuous links rather than direct causation.
- **Ethical Implications:** If such correlations become bases for decisions—hiring, insurance, medical coverage—people can be unfairly penalized for innocuous lifestyle factors or happenstance associations (like the weather in their region).

Example

Imagine a new policy refusing job interviews to individuals who frequent a certain online forum correlated with high employee turnover. The correlation might be genuine in the dataset, but the reason behind it could be entirely unrelated to work performance (e.g., that forum is more popular in a region with high turnover for unrelated economic reasons).

6. Big Data and the Digital Divide

“Divide between those who use big data and those who generate it. There is a systemic opacity in the use and handling of big data.” □cite□turn2file0□

This segment highlights *inequalities* in data collection and exploitation:

1. **Producers vs. Consumers:** Ordinary individuals generate data (social media posts, online searches, smartphone usage), but large corporations or governments have the resources to analyze and profit from it.
2. **Opacity:** People often have little understanding of how their data is handled, sold, or repurposed. They may not even realize they're "generating" data when simply browsing.

6.1 Impact on Decision-Making

"Those who are affected by the decisions of big data are not always in the position to understand it or challenge it." □cite□turn2file0□

When a banking algorithm denies a loan, the individual rarely has the ability to see why or how the decision was made (lack of transparency). This power asymmetry can undermine autonomy and fairness.

Real-World Example

A job applicant is rejected by an AI-driven screening platform. The candidate cannot easily appeal or understand which specific data points or correlations led to that rejection. This lack of recourse exemplifies the digital divide in action: the *algorithm's owners* have all the power.

6.2 Salient Examples from the Text

"Consider, for instance, the finding that 'people who fill out online job applications using browsers that did not come with the computer . . . but had to be deliberately installed (like Firefox or Google's Chrome) perform better and change jobs less often' (Andrejevic 2014: 1681)." □cite□turn2file0□

- **Browser Choice:** This correlation might exist in specific datasets, but it raises serious ethical and methodological questions: Is it fair or accurate to use someone's browser choice as a proxy for conscientiousness or technological savvy?
- **Socioeconomic Bias:** People who only have access to public computers—where they can't install anything—might be unfairly penalized.

7. Broader Ethical Implications for AI

All these big data challenges—lack of context, reliance on correlation over causation, and opaque decision-making—feed into the larger ethics of AI:

1. **Accountability and Transparency:** Who is responsible for decisions made by automated systems that rely heavily on big data correlations?
2. **Bias and Discrimination:** Data inevitably reflect social biases. AI can amplify these if not carefully managed.
3. **Consent and Privacy:** Individuals generating the data often do so unwittingly, raising concerns about informed consent.
4. **Regulatory Gaps:** Fast-moving technology outstrips policy measures, which struggle to keep pace with data analytics capabilities.

8. Concluding Reflections

From *Notes 4*, we see that big data's real power is in unveiling patterns at scale. Yet:

- **Contextual Understanding** is crucial: Data alone doesn't capture the *why* behind a pattern.
- **Theory & Judgment** remain integral: Despite claims of objectivity, human interpretation guides how data is collected, categorized, and used.
- **Ethical Tensions** emerge when correlation replaces explanation, potentially harming individuals who cannot challenge algorithmic decisions.
- **Inequities** persist between data collectors (governments, corporations) and data producers (ordinary citizens) who are subject to opaque analytics and decisions.

Ultimately, the text underscores that big data does not eliminate the ethical dimension—rather, it reshapes it, demanding new forms of scrutiny, regulation, and social dialogue. As big data continues to underpin many AI systems, recognizing and grappling with these ethical concerns becomes an integral part of designing and deploying technology responsibly. [□cite□turn2file0□](#)

Notes 5

1. Why Think About Algorithmic Accountability?

1.1 The Opaque Nature of Algorithmic Decisions

"The nature of decisions taken by algorithms are often opaque. There may be correlations that are not understandable to even those who are using it to arrive at decisions." [□cite□turn3file0□](#)

Algorithms, especially those using machine learning, frequently generate results that can be difficult to interpret. For instance, a neural network that screens job applicants may reject certain candidates without yielding human-readable explanations for *why* it identified them as unsuitable. This "black-box" effect can create tension between efficiency and the need for accountability.

- **Implication:** If decision-makers can't explain their model's outcomes, how can individuals contest unfair or harmful decisions? The note suggests a growing ethical expectation that subjects of algorithmic decisions *have the right* to a clear explanation or justification.

1.2 Biased Data and Embedded Values

"The data on which the algorithms are trained may be biased. Biased data may end up reproducing existing inequalities and patterns of discrimination." [□cite□turn3file0□](#)

Bias in AI can derive from historical or systemic inequalities embedded in data. For example, a credit-scoring model trained on past lending decisions might perpetuate discrimination if it learned from data reflecting racial or gender bias.

- **Key Question:** Should we treat machine learning outcomes as "useful heuristics" rather than "definitive knowledge"? The notes ask whether we can view these outcomes as context-dependent tools, rather than authoritative truths.

1.3 The Need for Explicit Values

“Even if the model is not trained in ethical data it still embeds certain values that is needed to be made explicit.” □cite□turn3file0□

Algorithms are not *value-neutral*. Whether the values come from the dataset itself or from designers' choices about objectives, thresholds, and definitions, these values can shape social outcomes. Making them explicit helps stakeholders grasp how a model prioritizes, for instance, accuracy over fairness—or how it defines “success.”

2. The Rationale Behind Transparency

2.1 Observability and Knowledge

Ananny and Crawford (2017), as quoted, argue:

“Transparency concerns ... rest on an epistemological assumption that ‘truth is correspondence to, or with, a fact.’ The more facts revealed, the more truth that can be known through a logic of accumulation.” □cite□turn3file0□

This view sees **transparency** as a means to gather enough factual details—e.g., design parameters, training data characteristics—to hold systems accountable. If the processes behind decisions are exposed, observers can judge whether the system is fair, accurate, or aligned with public values.

2.2 Transparency as Performative

“Transparency is thus not simply ‘a precise end state in which everything is clear and apparent,’ but a system of observing and knowing that promises a form of control.” □cite□turn3file0□

Here, transparency is more than revealing information; it's a *performative* act. Publicly disclosing aspects of how an AI system functions can build trust (or the semblance of trust) and convey that the responsible party is open to scrutiny. However, the notes also point out that transparency often assumes “audiences are competent, involved, and able to comprehend” the disclosed information. If those conditions aren't met, transparency might not yield the intended accountability.

- **Example:** A company might release a technical white paper detailing its recommendation algorithm's architecture. Despite this, if the average consumer or regulator doesn't possess the expertise to interpret the details, can we really claim meaningful transparency?
-

3. Three Forms of Opacity

Drawing on Burrell (2016), the notes outline three distinct forms of opacity:

“(1) Opacity as intentional corporate or institutional self-protection and concealment ... (2) Opacity stemming from the current state of affairs where writing (and reading) code is a specialist skill ... (3) An opacity that stems from the mismatch between mathematical optimization ... and the demands of human-scale reasoning and styles of semantic interpretation.” □cite□turn3file0□

1. **Intentional Concealment:** Companies may withhold critical details about algorithms for competitive advantage or to protect intellectual property.
2. **Technical Expertise Gap:** Even if details are shared, specialized coding or machine-learning knowledge might be necessary to interpret them properly.
3. **Mathematical vs. Human Reasoning:** Deep-learning models can operate in high-dimensional spaces, generating solutions beyond intuitive human comprehension. This is a structural opacity: the system is simply too complex for human minds to fully parse.

Ethical Tension: If decision-makers themselves do not fully understand how or why a model reached a conclusion, can they ethically delegate life-altering decisions (such as hiring or loan approvals) to that system?

4. Defining Algorithmic Accountability

Binns (cited in the notes) offers a definition:

"Party A is accountable to party B with respect to its conduct C, if A has an obligation to provide B with some justification for C, and may face some form of sanction if B finds A's justification to be inadequate." □ cite □ turn3file0 □

4.1 Justification, Sanction, and Transparency

For accountability to be robust, two conditions must be met:

1. **Justification:** The decision-maker (or system operator) must offer clear, reasoned explanations for how an outcome was reached.
2. **Enforcement:** If the explanation is found lacking or if harm is detected, there must be a tangible mechanism to sanction or correct the decision.

In algorithmic contexts, accountability means:

- **Disclosure of system design:** What data was used? How was it labeled?
 - **Disclosure of operational logic:** How does the model weigh variables?
 - **User recourse:** If you believe you've been treated unfairly, do you have the right to an appeal or a second review?
-

5. What Kind of Justifications "Count"?

"Does any justification count? Or only justifications that are not arbitrary in nature, that are reasonable, that are acceptable, and those that are public?" □ cite □ turn3file0 □

5.1 Reasonable vs. Acceptable Justifications

- **Reasonableness:** A justification might be based on a coherent, data-driven principle, yet still be controversial or unethical if it overlooks critical contextual factors (e.g., "We selected candidates based on personality tests that systematically disadvantage certain demographics").
- **Acceptability to Affected Parties:** The notes raise the question of whether the justification must be subjectively acceptable to those affected by it. This approach aligns with principles of **procedural**

justice, where outcomes are deemed more legitimate if the decision-making process is transparent and respectful of stakeholder input.

5.2 Universally Accepted Norms?

“Should the justifications be based on certain universally accepted norms that we cannot reasonably reject? What should those norms be?” □cite□turn3file0□

This question implies a search for moral or legal standards that transcend cultural or individual differences—perhaps akin to human rights frameworks. For instance, you might say a justification is legitimate if it does not discriminate based on protected traits (race, gender, religion, etc.), reflecting widely accepted anti-discrimination norms.

6. Strategies for Algorithmic Accountability

The notes outline practical measures:

“—Remove biases in data and the code that results from data. — Develop the possibility of offering explanations for the decisions. — Important to clarify what epistemic standards ... are required for the case at hand.” □cite□turn3file0□

1. **Bias Audits:** Evaluate datasets and algorithms for potential discrimination.
2. **Explainable AI:** Implement methods that provide interpretable models or post-hoc explanations to help stakeholders understand outputs.
3. **Epistemic Standards:** Clarify whether a system needs robust causal explanations or if correlation-driven predictions suffice, depending on context (e.g., high-stakes medical decisions might require a stronger causal basis than a movie recommendation system).

7. Is Opacity Always a Problem?

7.1 The “Neural Net Produced It” Defense

“In cases where decision-makers can provide no other explanation for a decision than that, say, a neural net produced it, we may decide that their justification fails by default.” □cite□turn3file0□

If an organization cannot articulate *any* reason for a decision other than the opaque workings of an AI, that may be deemed inadequate. Public reason demands at least a baseline explanation. For example, if an algorithm denies medical care coverage, simply stating “the neural network’s output was negative” will likely not meet accountability thresholds.

7.2 Contextual Goals vs. Outputs

“What matters will not be how a system arrived at a certain output, but what goals it is supposed to serve.” □cite□turn3file0□

Sometimes, the broader objectives or constraints under which the AI operates are more important than the exact method. **Example:** A search engine’s objective might be to rank results by popularity vs. relevance. Users may care more about that policy than the nitty-gritty of the ranking algorithm’s code.

- **Implication:** The notes suggest an “output focus” in some contexts—knowing whether a system is optimized for fairness or purely efficiency may suffice to hold it accountable at a macro level.

7.3 The Case Against Building the System

“If a system is so complex that even those with total views into it are unable to describe its failures and successes, then accountability models might focus on whether the system ... should be built at all.” (Ananny & Crawford, 2017). □cite□turn3file0□

In extreme cases, if the complexity leads to irreducible opacity, ethical deliberation might conclude that the risk of harm is too great. Or, at minimum, the system should operate under strict regulations or with built-in governance features to mitigate potential damage.

8. Concluding Reflections

Bringing all these points together:

1. **Accountability** in AI involves **justification** plus **mechanisms to enforce** that justification. Opacity complicates accountability, especially when algorithms cannot be easily explained.
2. **Bias** is not simply a technical glitch; it is entangled with historical inequalities and designers’ own values. Reducing it requires ongoing scrutiny and an ethical framework that includes fairness and non-discrimination as design goals.
3. **Transparency** isn’t a monolithic solution. While it can enable external scrutiny, it also assumes that the relevant audiences can interpret complex data. Full technical disclosure may still leave systems opaque if the scale or complexity of machine learning surpasses normal human comprehension.
4. **Context Matters.** High-stakes domains (healthcare, criminal justice, credit) demand more robust explanatory and accountability frameworks. Other domains might prioritize different forms of transparency—like clearly stated objectives or user recourse policies.
5. **Public Reason.** Even if technical explanation is elusive, we can still hold systems accountable by focusing on the *goals* of the system, the *data* it uses, and the *real-world consequences* it produces. If none of these can be adequately justified, building or deploying the system may be ethically questionable.

In short, *Notes 5* places accountability at the heart of ethical AI. Whether through transparency, bias mitigation, or alternative justifications, the ultimate goal is to ensure that people affected by algorithmic decisions have an avenue to understand, challenge, and seek redress when needed. When these routes are blocked by opaque design or impenetrable complexity, it raises deep ethical questions about whether such systems should be deployed in the first place. □cite□turn3file0□

Notes 6

1. Why Think of Transparency?

The notes start by revisiting the importance of **transparency**:

“Transparency concerns ... rest on an epistemological assumption that ‘truth is correspondence to, or with, a fact’ ... The more that is known about a system’s inner workings, the more defensibly it can be governed and held accountable.” (Ananny and Crawford 2017) □cite□turn4file0□

1.1 Logic of Accumulation

Transparency is often championed on the premise that revealing more information yields better oversight. If regulators or the public understand how an AI system arrives at its decisions, they can evaluate whether it is biased, fair, or functioning as intended. The underlying view is:

- **Observation → Insight → Knowledge → Accountability**
- **Implication:** The more facts we accumulate about an AI system (source code, training data, design parameters), the closer we get to “the truth” of its operation.

However, it’s worth noting that more information does not always guarantee *meaningful* understanding. Specialized knowledge might be required to interpret complex models, leading to potential gaps in lay comprehension.

1.2 Performative Aspect of Transparency

“Transparency ... includes an affective dimension, tied up with a fear of secrets ... This autonomy-through-openness assumes that ‘information is easily discernible and legible; that audiences are competent ...’” (Christensen and Cheney, 2015) □cite□turn4file0□

Transparency isn’t merely about dumping technical details into the public sphere. It’s also a *performance* of openness, which can build trust—or at least the *appearance* of trustworthiness. Yet this assumes that stakeholders have the requisite expertise and motivation to act on the information provided.

2. Three Forms of Opacity

Drawing on **Burrell (2016)**, the notes identify three kinds of opacity:

“(1) Opacity as intentional ... self-protection and concealment; (2) ... writing (and reading) code is a specialist skill; (3) ... mismatch between mathematical optimization ... and human-scale reasoning.”

□cite□turn4file0□

1. **Intentional Concealment:** Corporations may keep algorithms secret to guard intellectual property, or simply to avoid scrutiny.
2. **Technical Expertise Gap:** Even if source code is published, the average person can’t easily interpret thousands of lines of code or complex neural network architectures.
3. **Mathematical vs. Human Reasoning:** Modern AI, especially deep learning, often operates at a scale beyond human comprehension—“explanations” that might be mathematically valid are still opaque to non-experts.

Example

A face-recognition algorithm might have millions of parameters. Even an open-source release of the model’s code might not help a lay user understand why it misidentifies certain ethnic groups more often than others.

3. Accountability in Social Contexts

3.1 Social Embedding of AI

"We need to think of systems being embedded in the social contexts ... embody the hierarchies, exclusions, marginalization, power dynamics ... technology does not operate in a vacuum."

□cite□turn4file0□

An AI credit-scoring system might replicate systemic biases (e.g., redlining in housing loans) if trained on historically biased data. Accountability, then, requires analyzing **how** that data was generated and **why** it might reflect societal hierarchies. Merely examining the algorithmic code isn't enough.

3.2 Purpose and Impact

"A tool is being designed for a certain purpose. ... Who does it impact? How does it impact those whom it impacts?" □cite□turn4file0□

Accountability also hinges on clear goals:

- **Intended vs. Actual Use:** A system meant for benign tasks (like filtering spam) can be repurposed in harmful ways (like political censorship).
- **Differential Impact:** Even well-intended AI can unequally affect different demographic groups. This raises the question: is such differential treatment justified or ethical?

4. Kinds of Responsibility

The notes outline different responsibility concepts:

Causal Responsibility: "Did you play a contributory role in the wrong?"

Culpable Responsibility: "Could you have reasonably been aware of the wrong your contribution would cause?" □cite□turn4file0□

4.1 The Problem of Many Hands

Complex AI systems involve multiple actors: data collectors, model developers, testers, etc. Each might have partial responsibility for resulting harms. When a predictive-policing algorithm discriminates, blame could be diffused across multiple roles:

- **Data scientist:** Provided the training set.
- **Software engineer:** Implemented the classification logic.
- **Project manager:** Approved deployment.
- **End user:** Interpreted the results in a biased manner.

Ethical Challenge: How do we distribute responsibility fairly? Are we holding the correct people accountable, or do they each bear partial responsibility?

4.2 Oversight Mechanisms

"... we need to think of oversight mechanisms that are able to trace the responsibility chain ..."

□cite□turn4file0□

Oversight may require auditing logs, version control, or system design decisions that reveal who contributed which parts. If the chain of responsibility is transparent, we can more effectively identify *where* biases or design flaws entered the process.

5. Oversight as Operationalizing Accountability

Kroll (2020), cited in the notes:

“Building AI systems that support accountability ... necessitates designing those systems to support robust oversight. ... Accountability is tied directly to the maintenance of records.” □cite□turn4file0□

5.1 Evidence and Record-Keeping

Accountability depends on structured record-keeping:

- **Version Histories:** Capturing changes in the model or data over time.
- **Documentation:** Recording rationales for parameter choices and known limitations.
- **Decision Logs:** Tracking input data and outputs for each key decision, enabling after-the-fact audit.

5.2 Contextual Norms

“The oversight entity ... tie(s) the actions described in those records to consequences.” □cite□turn4file0□

Oversight isn't one-size-fits-all: *ethical norms vary* across contexts (e.g., a medical AI system's oversight might differ from a social media recommendation engine). The system's context sets the bar for what's permissible, and oversight ensures compliance with those norms.

6. Algorithmic Accountability Defined (Binns)

“Party A is accountable to party B ... if A has an obligation to provide B with some justification ... B may sanction A if the justification is inadequate.” (Binns 544) □cite□turn4file0□

6.1 Justifications and Sanctions

- **Obligation to Justify:** The system's creators or operators must *explain* how it made a particular decision.
- **Potential Consequence:** If that explanation fails to meet a standard of reasonableness or fairness, a sanction should follow—ranging from fines and retractions to shutting down the system.

Key Insight: Accountability loses its force if no penalty exists. Without sanctions, we have “responsibility without accountability,” which rarely compels meaningful change.

6.2 Answerability and Outcome Responsibility

“Individuals or organizations can be made to answer for outcomes of their behavior ... or the behavior of tools they make use of. ... Ties actions or outcomes to consequences.” □cite□turn4file0□

In practice:

1. **Explain:** The decision maker must clarify how the AI was used.
 2. **Assess:** Stakeholders judge the explanation's sufficiency.
 3. **Enforce:** If flawed or harmful, those responsible face consequences (financial, legal, reputational).
-

7. What Kind of Justifications Count?

"Does any justification count? Or only justifications that ... are reasonable, acceptable, and public?"
 □cite□turn4file0□

7.1 Criteria for Valid Justifications

- **Non-Arbitrariness:** Justifications can't be random or purely self-serving.
- **Public Reason:** They should be understandable in a public forum, not cloaked in excessive jargon.
- **Acceptability to Affected Parties:** If the justification isn't comprehensible or legitimate to those impacted, accountability falls flat.

Example: A credit-scoring system might say, "Your loan application was rejected due to a proprietary algorithm's negative score." That's a minimal explanation. But if the applicant cannot grasp *why* the algorithm assigned that score—and no further clarifications are provided—that fails as a justification.

7.2 Universally Accepted Norms

"Should the justifications be based on certain universally accepted norms ...?" □cite□turn4file0□

This points to a broader philosophical debate. Ethical systems often rely on either:

- **Deontological Norms:** E.g., "Thou shalt not discriminate based on race, gender, etc."
- **Consequentialist Norms:** "Actions are justified if they produce the best outcomes overall."
- **Discourse Ethics:** "Justifications are valid if no stakeholder can *reasonably reject* them."

In the global arena, some norms (like non-discrimination) approach universal acceptance, though cultural differences complicate how they're interpreted.

8. Conclusion: Tying It All Together

Notes 6 weave together four core themes:

1. **Transparency:** A potential pathway to accountability, but not a silver bullet, especially given opacity's multiple causes.
2. **Social Context:** AI is built and deployed in socially stratified environments, so it inevitably inherits biases and power structures.
3. **Responsibility & Oversight:** Effective accountability requires clearly delineating causal and culpable responsibility across all who contribute to AI systems. Oversight structures (audits, record-keeping) help trace how decisions were made and by whom.
4. **Justifications & Norms:** Ultimately, accountability hinges on providing robust, understandable justifications aligned with norms that stakeholders (and society at large) deem legitimate. If an AI's operators cannot or will not meet that standard, the system may be deemed unfit for deployment.

In short, as AI becomes more influential in high-stakes scenarios—banking, employment, policing, healthcare—**accountability** is no longer optional. We need a holistic approach that considers the *technical* opacity of AI, the *social* context in which it operates, and the *ethical frameworks* that legitimize or reject certain decisions. This ensures that algorithmic decisions not only serve efficiency but also uphold fairness, responsibility, and respect for human agency. □cite□turn4file0□