

→ Chilling Effect: Enacting policies to hide existence of problem via fear, ~~ex~~ instead of solving ~~set~~ the problem

Relevance: AI companies claim sentiment understanding via visual media ⇒ Bogus smiles faked in South Korea by employees

Snake Oil ⇒ Exploiting fear to sell goods ⇒ AI snake oil

→ Massive effort to influence public opinion ⇒ Successful

→ AI { ^{→ crude differentiation} Genuine progress: ~~works, does~~ & should ^(public ~~bad~~, can be harmful)
Useful, but not much progress: Works but bad performance
Fraud: Doesn't, Won't, Harmful ^{perception}

→ Genuine progress: Facial recognition, speech to text, Medical diagnosis from scans, Deepfake
Automated ~~judgment~~ judgement ^(moral dilemma)

→ Improving: Spam detection, Copyright detection, Automated essay grading
↓
Can be ^{done} ~~done~~
but far from perfect in. Predicting social outcome

→ Fraud: Predicting, policing, terror risk, risky kids, job performance, criminal recidivism, ~~how~~ anything related to human behavior

~~Isolated corpus - Gutenberg~~

~~submitted~~

→ ~~Members of UNs~~ need to acknowledge Human rights

→ 13000 parameters based model \approx 4 parameter based linear regression
(in social prediction)

→ Accuracy of recidivism prediction:

137 features: COMPAS tool: $65\% \pm 1\%$ (slightly better than ^{stand})

~~logistic regression~~: 2 features: Logistic Regression: $67\% \pm 2\%$
↑
Age & No. of priors

★ ~~Na~~ Narayana claims AI is shit at predicting social outcomes

→ Harms of AI in social prediction :- Hunger for personal ^{from domain experts & workers} data (privacy harm)

- Massive power transfer to unaccountable tech company

- Lack of explainability - Veneer of accuracy - Distracts from interventions

Taken Takeaways:

→ AI shit at social prediction

→ Funding is a cheap way of ^{solving} ~~half~~ problems

→ Sometimes, Manual > Automation
testing

Philosophy → Mind
→ Ethics

→ What is this

→ Ontology: What is reality composed of?

→ Epistemology: What is truth? Is it real or my mind's playing tricks?
↓
How do I know what it is? What is knowledge?

→ Ethics/Moral Philosophy

Ethics

- Moral philosophy: What is ethics, is it different from morality?
- Normative Ethics: ^{Good/Bad Quality} Virtue ethics: Quality of a person (not concerned with actions)
- Deontological ethics: ^(German) "Majority of human civil." Focuses on the action's nature, action's intention
- Immanuel Kant (POI) Certain things are bad (no ifs, ands & buts) { No regard for consequences or cause.
- Categorical Imperative
- Critique of Reasoning
- John Rawls
- Veil of Ignorance
- Remove everything, become completely impartial (British)
- Consequentialism: Focus on the impact
- Jeremy Bentham → Utilitarianism (criticised) → Contextual (By Karl Marx)
- State consequentialism (reaction (kind of))
- Welfarism (not ~~defined~~ clear here)
- Pragmatism: ^{also Indian} American way, representative of everyone, mutual decision
- John Dewey ^{student} → B.R. Ambedkar Dr.

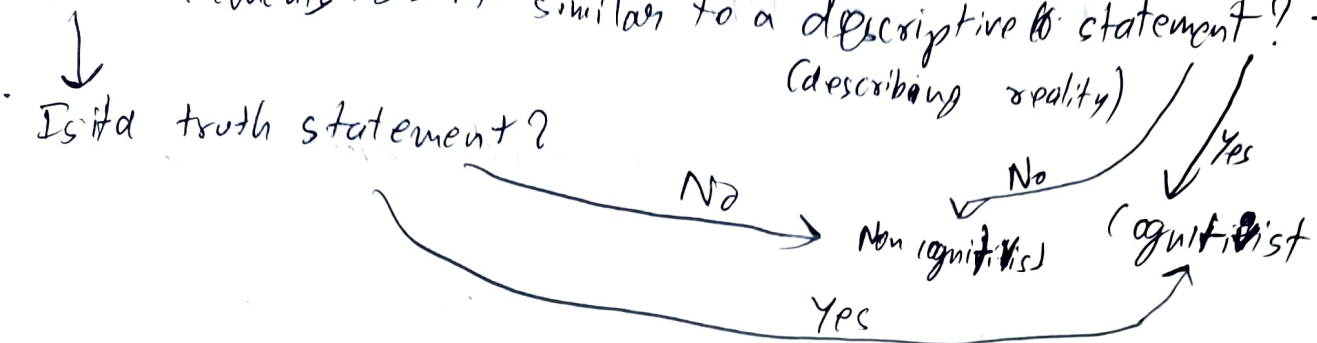
Applied Ethics

- ~~AI control problem~~ → Machine ethics: Instructor believes machines shouldn't even be in discussion with ethics
- AI control: AI taking harmful decisions for task completion (murder for winning a game)
- Economic justice cannot come at the cost of social justice

Metaethics

Policy brief: List of demands from policymakers

Moral statements: Is it similar to a descriptive statement? → Ambiguous answer
(describing reality)



Cognitivist: Moral statements are truth-apt

→ Error theory: All moral statements are false (are cognitivist)
(theorist) (All statements have an ans., it is false)

→ Some are true/false

→ Objective, i.e. exist in reality: Moral Realist/objectivist

→ Only a social construct, subj.-only: Moral subjectivist/relativist

Non cognitivist: Moral statements are not truth-apt

→ Emotivist: Morality is an emotion guided.

→ Prescriptivist: Moral statements are not true/false, they're wills/orders.

→ Quasi-realist: Not real, humanity adopted it as pseudo-real stuff to live life

→ Realized by David Hume

→ Is-ought gap: Most moral args start with logical statements and then go to an ought statement (should, ought). Where does this description change to prescription?

Normative statements come somewhere in all arguments, the theory doesn't question its existence, rather it is concerned with how, why & where does this transition occur.

Ethics

→ Ethics initiatives

- Gov By trans-national orgs, ^{like} UN, IEE

→ Ethics boards

- Companies future proofed themselves to avoid investment in "unethical" activities

→ Lack of actionability

- No one can do anything if "unethical" or non compliance with certain standards, due to ethics being entirely theoretical & ^{hence} no legal basis for the compliance of standards by AI companies

→ Project Maven

- Idea by Pentagon after 9/11, to track everything of every person on Earth by satellites ~~(theoretical)~~ ^(ambition)
- Pitched to ~~Google~~ Google, 1st time a strike occurred? & employees left (due to ethical concerns); hence Google had to refuse, ^{Pentagon}
- ~~Some other~~ Palantir picked it up

→ Lack of Transparency

Issues with Tech Ethics

→ Individualist solutions

- Solutions are not ^{preventing} ~~avoiding~~ the generation of ethical issues, they're targeting individual problems, which will always keep popping up.
Ex: Project maven stopped by Google, when it shouldn't have been thought of.

→ Corporate logics

- Company won't let anything impact its bottom line

→ Contradiction & Vagueness

→ Ethics washing

- Ethics suffers from these | - Describe practices

→ Determinism & Solutionism

- ~~Solution~~ Everything ought to be solved ^{shown by}
↳ Wrong, not only morally, but because it is reductionist / facile as

What is a mind?

- what changes bet/bw consciousness & unconsciousness? We don't know
- Chinese Room: Analogy for AI, ~~NN~~neural network, computer (originally)

↓
An AI, a computer can't ~~ever~~ have a mind.

→ Dualism

- Mind is separate from matter

→ Monism

- No such distinction b/w mind & matter.
- Idealism: Use mind to think about "matter", hence only ideas & its realm is what is thought about
- ~~Matter~~ materialism: Reality can be measured and known, not just ideas in mind.

Philosophy of Mind

- Can't know other ppl have minds
- If matter is what makes up the mind & it's governed by physics, so is the mind, and if so, does free will exist?
- ~~think~~ Can think of things with which there is 0 material contact.

Materialism

- Behaviorism: Behaviour has a 1 to 1 connect with the mind, hence behaviour determines mind.

↓
Wrong

- Physicalism / Identity theory
 - Type Identity theory, ~~Mental~~
 - Token I I : Black box

- Qualia: Subj. conscious experiences, No way to know if everyone has the same or diff. qualia

EAI

→ Explainability ↔ Transparency. Explainability in general & in to an engineer is not the same.

Procedural & Technical expl.
expl.

Issues

→ Ignorance due to no outward looking & considering the society

→ Menamara's philosophy says can't be measured ⇒ not knowledge
↓

Menamara fallacy ↔ ignore

→ Epistemic contestation: which kind of knowledge should have privilege?

→ Political epistemic contestation: Does't count for how tech affects politics & economics.