

Sept 3

→ snake-oil reading:-

- AI → personality assessment } → fraud:
 - ↳ no AI can tell you what a person is or suitable/not suitable
- prediction } only what is observable.
- thinking & suitability } are not observable } not tangibles.
- funding? → I would pay good money if it inspires terror in my employees.
- "chilling effect" → ~~enacting~~ (enacting certain policies not with the agenda of solving anything but instead scares into discipline)
- "SK → companies recommend plastic surgeries & how to smile so that office AI does not target them"
- (AI on twitter)
- Software increases productivity X → terror makes us ^{work} hard
- productivity in econ :- same input → more output
here more " ← more "
- Why so much snake oil in AI?-
 - ↳ In olden America — people claim magically medicine.
 - Snake oil → selling using fear & greed of people
- * public opinion:-
 - ↳ thinks AI can be powerful } dangerous
 - ↙ incentivises AI products even though not useful.

AI works & well & genuine progress.
AI works - part bad & careful.
AI fraud & harmful

face recognition
deep fakes
speech to text
works

no interpretation / reasoning
only identification } harmful

Spam detection, automatic essay grading.
Content recommendation
↳ not always good performing.

don't work :-
but improving

ethical concerns

high accuracy.

fundamentally
dubious :-

errors

dividing social outcome / involving telepathy
emotion detection or predicting human behavior
doesn't work.

don't work

complex society

fundamentally
pseudo scientific polys.

fragile families :-

list of features to predict child well being

457 scientists } 12,942 variables for 4,242 families

results \cong with 7 features

linear regression

COMPAS :-

problem

too many people in jail

jail vs prison

undertrial

guilty

lot need mental help

help all

help most victimised people with
AI ! } compas tool.

predicts
which
prisoners
most

vulnerable

guilty dangerous

decided after

Process is punishment }?

COMPAS for bail/parole or recidivism
(prob of reoffending)
↓
doesn't work but America uses it
↓
giving parole/not

→ fairport :- COMPAS → 65% ± 1%.
2-featured logistic regression → 67% ± 2%.

also biased.

→ would you use it if it has 99%? → No
better 100 guilty than
1 innocent punished.

→ COMPAS has no morality.

— Hazards:-

Smart city in India } certain things use digital automated solns
given by some companies
→ problem not actually
solved but
seems to everyone
people can't
vote on these
& they're unaccountable.

Paper 5 reading:- bird's eye view of how people think about
AI ethics.

Week 6 slides:- F/A/T → Transparency
fairness → accountability

→ Philosophy to understand ethics.

(P's 144)

→ deals with what quality is? what humans are - etc

Science = Natural philosophy

Moral philosophy } how humans should act

(economics is a part)

philosophy encompasses everything

pillars of philosophy

literal statement
not scholastic

ontology
(what is reality?)

how do I know
right/truth
(Epistemology)

what do we need to
do?

(Ethics / Moral
philosophy)

(Mind & Ethics)

In common world

morality → religious
ethics → secular

to philosophers
it's same

3 questions -

① How do we decide it's ethical? } Normative ethics

② What is ethics? } Meta ethics

③ How do I translate to real life ethics? } Applied ethics

Normative:- what's right thing to do?

→ how people decide right/wrong

not consistent
widely changed
over time & space

① ancient greek/china/India } virtue ethics
individual

(does a person have good qualities?)
doesn't look at actions.

Ex:- Hercules in greek → burris countries but is
hero.

② Deontological ethics } → what was the intention
behind the action.

looks at fund intention of action → good/bad.

"The critique of pure reason"
↓
limit } → categorical imperative
*(Emmanuel)

doesn't care about
the situation
Ex:- Murder is bad

* John rogers } → Veil of ignorance

first remove all context & bias, then we
(desire... etc) & be clear I can tell
if action is good

Natural rights

} criticism → didn't look at cause & consequence of action.

② Consequentialism: → only look at the impact of action.

Bentham } → what is good? → what makes people happy
utility of action

problem:- who decides utility

↓
State consequentialism

problem:- if state is partial

Carl Max? → what you call universal utility is
for 1 Shopkeeper (Britain?)
} shopkeeper mentality

Still a consequentialist
Welfarism } improve life of people.

④ Pragmatic: → very american way of ethics
→ get everybody of different social background
to decide mutual.

Indian law } deontological
constitution } pragmatic

Applied Ethics: -

↳ Machine ethics? can machines be moral agents?

AI control problem: -

do what we don't want it to do?
→ is this a real problem? why does this have control over others?

Altruism & longtermism: -

→ use AI to make as much money & distribute to people best?