

# CY0002 Ethics Notes

Hankertrix

April 29, 2025

## Contents

<b>1</b>	<b>Definitions</b>	<b>6</b>
1.1	Reasoning . . . . .	6
1.2	Inductive reasoning . . . . .	6
1.3	Deductive reasoning . . . . .	6
1.4	Abductive reasoning . . . . .	6
1.5	Analogical reasoning . . . . .	6
1.6	Modus ponens . . . . .	7
1.7	Natural deduction (Gentzen-style logic system) . . . . .	7
1.8	Hilbert-style logic system . . . . .	7
1.9	Propositional logic ( <i>PL</i> ) . . . . .	7
1.10	Proposition . . . . .	7
1.11	Propositional variable . . . . .	8
1.12	Antecedent . . . . .	8
1.13	Consequent . . . . .	8
1.14	Syllogism . . . . .	8
1.15	Quantificational or first-order predicate logic ( <i>QL</i> ) . . . . .	8
1.16	Deductive argument assumption . . . . .	8
1.17	Descriptive proposition ("is") . . . . .	9
1.18	Normative proposition ("ought") . . . . .	9
1.19	Hume's Law (Autonomy of Ethics / NOFI principle) . . . . .	9
1.20	Russell's Law . . . . .	9
1.21	Hume's Second Law (The Problem of Induction) . . . . .	9
1.22	Kant's Law . . . . .	10
1.23	Barrier construction theorem . . . . .	10
1.24	Geach-style conditionalisation . . . . .	11
1.25	A priori . . . . .	11
1.26	A posteriori . . . . .	11

1.27	Denotatum (plural: denotata)	11
1.28	Argumentum a fortiori (a fortiori)	11
1.29	Shew / Shewn	11
1.30	Counterfactual conditionals	11
1.31	Reductio ad absurdum	11
1.32	Contradiction ( $\perp$ )	11
1.33	Tautology ( $\top$ )	12
1.34	Enthymemes	12
1.35	Salva Veritate (Salva Validitate)	12
1.36	Contingent truth	12
1.37	Ampliative	12
1.38	Nash equilibrium	13
1.39	Pareto optimality (Pareto efficiency)	13
1.40	Validity of an argument	14
1.41	Soundness of an argument	14
1.42	Classificatory moral commitments	14
1.43	Substantive moral commitments	14
1.44	Normative neutrality	14
1.45	Axiology	14
1.46	Reflective equilibrium	14
1.47	Eudaimonia	15
1.48	Hedonism	15
1.49	Avant la lettre	15
1.50	Elenchus elenctic (The Socratic method)	15
1.51	Agent	15
1.52	Patient	15
1.53	Prima facie	15
1.54	Virtue ethics	16
1.55	Virtual epistemology	16
<b>2</b>	<b>Prior's paradox</b>	<b>17</b>
2.1	Assumptions	17
2.2	Proposition	17
<b>3</b>	<b>Defending Hume's Law</b>	<b>19</b>
3.1	Issues with move 7	20
3.2	Gerhard Schurz substitution	21
3.3	Gibbard-Karmo-Singer semantics	22
3.4	Shorter's position	22

<b>4</b>	<b>Modal concepts</b>	<b>25</b>
<b>5</b>	<b>Deontic logic</b>	<b>26</b>
5.1	Analogies between the deontic mode and the alethic mode . .	26
5.2	Disanalogies . . . . .	28
5.3	Components . . . . .	29
5.4	Proof theory . . . . .	30
5.5	Monotonicity of entailment (RM) . . . . .	31
5.6	Corollary of the monotonicity of entailment (Corollary of RM)	31
5.7	Theorem T1 . . . . .	31
5.8	Paradoxes of obligation . . . . .	32
5.9	Resolving the paradoxes . . . . .	34
5.10	Mimamsa deontic logic . . . . .	35
<b>6</b>	<b>Inductive reasoning</b>	<b>36</b>
6.1	Newcomb's paradox (Newcomb's problem) . . . . .	36
6.2	Philosophical principles . . . . .	37
6.3	Bayes' theorem . . . . .	38
6.4	Conditional probability . . . . .	38
6.5	Deductive vs inductive reasoning . . . . .	39
6.6	Axioms of probability (Kolmogorov theorem) . . . . .	40
6.7	Axioms of expected utility theory (von-Neumann-Morgenstern theorem) . . . . .	41
6.8	Bayesian decision theory . . . . .	44
6.9	Jeffrey-Bolker theory . . . . .	44
<b>7</b>	<b>Issues with standard decision theory</b>	<b>48</b>
7.1	Consistency versus responsibility . . . . .	48
7.2	Cognitive biases . . . . .	48
7.3	Deontological decision-making . . . . .	49
7.4	Collective decision-making and voting paradoxes . . . . .	50
<b>8</b>	<b>Game theory</b>	<b>52</b>
8.1	Comparison with decision theory . . . . .	52
8.2	Prisoner's dilemma . . . . .	52
8.3	Table of utility values . . . . .	53
8.4	Modification of rationality assumption . . . . .	54
8.5	Possible strategies for the game . . . . .	54
8.6	Axelrod's tournaments . . . . .	54
8.7	Rapoport et al.'s objections to the Axelrod tournaments . . .	56

<b>9</b>	<b>Ethics</b>	<b>58</b>
9.1	Branches of ethics . . . . .	59
9.2	Normative neutrality requirement . . . . .	59
9.3	Forcehimes' collapse argument . . . . .	60
9.4	Normative theory . . . . .	60
9.5	High moral theory . . . . .	61
<b>10</b>	<b>Consequentialism</b>	<b>64</b>
10.1	Consequentialist decision-making . . . . .	64
10.2	Theory of the Good . . . . .	65
10.3	Types of consequentialism . . . . .	66
10.4	Types of utilitarianism . . . . .	66
10.5	Hedonism . . . . .	67
10.6	Benthamite utilitarianism . . . . .	67
10.7	Hedonic calculus . . . . .	68
10.8	Issues with quantitative hedonism . . . . .	68
10.9	Qualitative hedonism . . . . .	70
10.10	Haydn and the oyster thought experiment . . . . .	74
10.11	Utility monster thought experiment . . . . .	75
10.12	Against hedonism . . . . .	76
10.13	Formal definition of consequentialism . . . . .	79
10.14	Violation of consequentialist axioms . . . . .	84
10.15	Driver's strategy . . . . .	86
<b>11</b>	<b>Deontology</b>	<b>88</b>
11.1	Hiroshima and Nagasaki . . . . .	88
11.2	Textbook view . . . . .	89
11.3	Mafia scenario . . . . .	90
11.4	Side constraints . . . . .	91
11.5	The right and the good . . . . .	91
11.6	Action and intention . . . . .	92
11.7	Doctrine of double effect . . . . .	92
11.8	Avoision . . . . .	93
11.9	Trolley dilemma . . . . .	93
11.10	Siamese twins . . . . .	95
11.11	Divine command theory . . . . .	96
11.12	Modified divine command theory . . . . .	98
11.13	Plato . . . . .	99
11.14	Secular deontology . . . . .	101
11.15	Kantian deontology . . . . .	101

11.16	Organ transplant thought experiment . . . . .	106
11.17	Rossian deontology . . . . .	108
11.18	W.D. . . . .	110
11.19	Korsgaard strategy . . . . .	117
<b>12</b>	<b>Virtue ethics</b>	<b>119</b>
12.1	Comparison to other normative theories . . . . .	119
12.2	Ancient virtue ethics . . . . .	119
12.3	Two kinds of virtues . . . . .	120
12.4	Eudaimonist virtue ethics . . . . .	121
12.5	Aristotelian virtue ethics . . . . .	122
12.6	Doctrine of the mean . . . . .	123
12.7	Logic of virtue . . . . .	124
12.8	Agent-based virtue ethics . . . . .	126
12.9	Target-centred virtue ethics . . . . .	128
12.10	Summary . . . . .	128
12.11	Objections to virtue ethics . . . . .	129
12.12	Responses to objections . . . . .	130
12.13	Anscombe . . . . .	133
<b>13</b>	<b>Logic symbols</b>	<b>136</b>

# 1 Definitions

## 1.1 Reasoning

Reasoning is the process by which one draws **conclusions** from a set of premises. Reasoning or inference may be represented as follows:

$$\phi \vdash_X^M \psi$$

Where:

- $\phi$  is the set of premises
- $\vdash$  represents that  $\psi$  is provable from  $\phi$
- $M$  refers to the mode of inference, which can be deductive, inductive, abductive, analogical, etc.
- $X$  is the inferential mechanism.
- $\psi$  is the conclusion set.

## 1.2 Inductive reasoning

Inductive reasoning refers to making observations to find patterns and using the patterns to reason about things. Inductive reasoning usually works better with a large sample size.

## 1.3 Deductive reasoning

Deductive reasoning refers to drawing conclusions using a formal logic system, like mathematics, for example.

## 1.4 Abductive reasoning

Abductive reasoning refers to seeking the simplest and most likely conclusion from a set of observations. Abductive reasoning is usually used when there is a small sample size.

## 1.5 Analogical reasoning

Analogical reasoning is a special type of inductive reasoning where perceived similarities are used as a basis to infer some further similarity that has not been observed yet.

## 1.6 Modus ponens

- Modus ponens, also known as modus ponendo ponens, which is Latin for "mode that by affirming affirms".
- It can be summarised as "P implies Q. P is true. Therefore, Q must also be true."

## 1.7 Natural deduction (Gentzen-style logic system)

Natural deduction is a kind of proof calculus in which logical reasoning is expressed by rules that are closely related to the "natural" way of reasoning.

## 1.8 Hilbert-style logic system

- A Hilbert-style proof system is a type of formal proof system.
- It is defined as a deductive system that generates theorems from axioms and inference rules, especially if the only inference rule is modus ponens.
- Every Hilbert system is an axiomatic system.

## 1.9 Propositional logic ( $PL$ )

$$\phi \vdash_{PL}^{ND} \psi$$

Where:

- $\phi$  is the set of premises
- $\vdash$  represents that  $\psi$  is provable from  $\phi$
- $ND$  refers to the deductive mode of natural deduction
- $PL$  refers to the inferential mechanism of propositional logic
- $\psi$  is the set of conclusions derived from the set of premises

## 1.10 Proposition

A proposition is a statement or an assertion that can be either true or false.

### 1.11 Propositional variable

A propositional variable is a variable that is used to capture the content of a proposition, which is a statement that can be either true or false. Usually they are  $p$  and  $q$ .

### 1.12 Antecedent

The antecedent is the statement in which a statement is inferred. It is  $p$  in the example below:

$$p \rightarrow q$$

### 1.13 Consequent

The consequent is the statement that is inferred from another statement. It is  $q$  in the example below:

$$p \rightarrow q$$

### 1.14 Syllogism

A syllogism is a kind of logical argument that applies deductive reasoning to arrive at a conclusion based on two propositions that are asserted or assumed to be true.

### 1.15 Quantificational or first-order predicate logic ( $QL$ )

$$\phi \vdash_{QL}^{ND} \psi$$

Where:

- $\phi$  denotes any well-formed formula (wff) in quantificational logic
- $\vdash$  represents that  $\psi$  is provable from  $\phi$
- $ND$  refers to the deductive mode of natural deduction
- $QL$  refers to the inferential mechanism of quantificational or first-order predicate logic
- $\psi$  is the set of conclusions derived from the set of premises

### 1.16 Deductive argument assumption

The deductive argument assumption assumes that the conclusion of an argument cannot contain more information than is held in its premises.



### 1.17 Descriptive proposition ("is")

A descriptive proposition is a statement of fact.

### 1.18 Normative proposition ("ought")

A normative proposition is a proposition that contains a value judgment, like a moral judgment or ethical judgment.

### 1.19 Hume's Law (Autonomy of Ethics / NOFI principle)

Hume's First Law states that we cannot deduce how things ought to be or what ought to be done, which is a moral judgment from how things are, which is a statement of fact. This is also known as the view that ethics is autonomous. This is also known as the no ought from is principle, or NOFI.

### 1.20 Russell's Law

You can never arrive at a general proposition by inference from particular propositions alone. You will always have to have at least one general proposition in your premise.

$$\phi \not\vdash \psi$$

Where:

- $\phi$  is the particular proposition
- $\not\vdash$  means "does not entail that"
- $\psi$  is the universal or general proposition

### 1.21 Hume's Second Law (The Problem of Induction)

Hume's second law states that you cannot derive propositions about the future from propositions about the past or present.

$$\phi \not\vdash \psi$$

Where:

- $\phi$  are the propositions about the past or present
- $\not\vdash$  means "does not entail that"
- $\psi$  are the propositions about the future

### 1.22 Kant's Law

Kant's Law states that you cannot derive necessary propositions from propositions about the actual world.

$$\phi \not\vdash \psi$$

Where:

- $\phi$  are the propositions about the actual world
- $\not\vdash$  means "does not entail that"
- $\psi$  is the necessary propositions

### 1.23 Barrier construction theorem

Implication barrier	Topic	Description
Hume's (1739/40) Law	Normativity	You cannot derive <b>normative propositions</b> ( $Op$ ) from <b>descriptive propositions</b> ( $p$ )
Russell's (1918) Law	Generality	You cannot derive <b>general propositions</b> ( $\forall xFx$ ) from <b>particular propositions</b> ( $Fa$ )
Hume's (1748) Second Law	Time	You cannot derive <b>propositions about the future</b> ( $Fp$ ) from <b>propositions about the past or present</b> ( $Pp$ )
Kant's (1787) Law	Necessity	You cannot derive <b>necessary propositions</b> ( $\Box p$ ) from <b>propositions about the actual world</b> ( $p$ )

### **1.24 Geach-style conditionalisation**

Geach-style conditionalisation refers to embedding "ought" propositions in conditionals, which appear to allow us to derive valid is-ought inferences.

### **1.25 A priori**

A priori is a Latin phrase meaning "from the earlier".

### **1.26 A posteriori**

A posteriori is a Latin phrase meaning "from the later".

### **1.27 Denotatum (plural: denotata)**

Denotatum means a denotation of a word or an expression. The denotation of a word or expression is its strictly literal meaning, so the English word "warm" would denote a high temperature.

### **1.28 Argumentum a fortiori (a fortiori)**

Argumentum a fortiori is a Latin phrase meaning "argument from the strong reason".

### **1.29 Shew / Shewn**

Shew is just an archaic alternative form of show.

### **1.30 Counterfactual conditionals**

Counterfactual conditionals are conditional sentences which discuss what would have been true under different circumstances, e.g., "If Peter believed in ghosts, he would be afraid to be here."

### **1.31 Reductio ad absurdum**

Reductio ad absurdum, Latin for "reduction to the absurdity", disproves a proposition by showing that it leads to absurd or untenable conclusions.

### **1.32 Contradiction ( $\perp$ )**

A contradiction is a statement that is always false.

### 1.32.1 Principle of explosion

Anything follows from a contradiction (including a **normative proposition** of the form  $Op$ ).

### 1.33 Tautology ( $\top$ )

- A tautology is a statement that is always true.
- A tautology or logical **truth** follows from anything (including a **descriptive proposition** of the form  $p$ ).

### 1.34 Enthymemes

Enthymemes are arguments with hidden premises.

### 1.35 Salva Veritate (Salva Validitate)

Salva Veritate is a Latin phrase for "with unharmed truth". It means that something can be done without changing the validity of the argument.

### 1.36 Contingent truth

A contingent truth is true as it happens, or as things are, but that did not have to be true.

### 1.37 Ampliative

Ampliative means "extending" or "adding to that which is already known".

### 1.38 Nash equilibrium

Nash equilibrium refers to a play in which each strategy is the **best response** to the strategy played by the other person.

#### 1.38.1 Example

		P2	
		STRATEGY 3	STRATEGY 4
P1	STRATEGY 1	(5, 5)	(2, 3)
	STRATEGY 2	(0, 1)	(4, 2)

- P1's strategy 2 is the **best response** to P2's strategy 4 (and vice versa).
- P1's strategy 1 is the **best response** to P2's strategy 3 (and vice versa).

The cells coloured in **yellow** denote the **Nash equilibria**.

### 1.39 Pareto optimality (Pareto efficiency)

A state of affairs such that there is **no alternative state of affairs** that would **make some people better off without making at least one person worse off**.

#### 1.39.1 Example

		P2	
		STRATEGY 3	STRATEGY 4
P1	STRATEGY 1	(5, 5)	(2, 3)
	STRATEGY 2	(0, 1)	(4, 2)

The cell coloured in **green** denotes a **Pareto-optimal** state of affairs.

#### 1.40 Validity of an argument

An argument is **valid** if, **assuming the truth of all its premises**, its conclusion must, by **logical necessity**, be true too.

#### 1.41 Soundness of an argument

An argument is sound if **all of its premises are in fact true**, or it does not contain any **false premises**, and it is a **valid argument**.

#### 1.42 Classificatory moral commitments

Classificatory moral commitments are defined as the commitments that result from delineating the **scope of the moral domain**.

#### 1.43 Substantive moral commitments

Substantive moral commitments are defined as a **normative bias**.

#### 1.44 Normative neutrality

- Normative neutrality between **competing moral standards and rules of conduct** is where the cut between **metaethics (2<sup>nd</sup>-order theory)** and **normative theory (1<sup>st</sup>-order theory)** is made.
- Essentially, normative neutrality is what separates **metaethics and normative theory**.

#### 1.45 Axiology

Axiology just means value theory.

#### 1.46 Reflective equilibrium

Reflective equilibrium is a method of balancing moral principles and judgments to arrive at the content of justice.

### 1.47 Eudaimonia

- Eudaimonia is a Greek word that means the state or condition of good spirit, and is often translated as happiness or welfare. In Aristotle's works, it means the highest human good.
- It is a certain **flourishing** or the **sort of happiness worth seeking or having**.

### 1.48 Hedonism

#### 1.48.1 Psychological hedonism

Only **pleasure (happiness)** or **pain (unhappiness)** motivates us.

#### 1.48.2 Ethical hedonism

Only **pleasure (happiness)** has **value** and only **pain (unhappiness)** has **disvalue**.

### 1.49 Avant la lettre

Avant la lettre means that a concept exists even before a term is coined for it.

### 1.50 Elenchus elenctic (The Socratic method)

The Elenchus elenctic is a form of argumentative dialogue between individuals based on asking and answering questions.

### 1.51 Agent

The agent is the person who is performing an action.

### 1.52 Patient

The patient is the person on whom the action is performed.

### 1.53 Prima facie

Prima facie is a Latin phrase meaning "at first sight", or "based on first impression". It is used in philosophy to indicate that something is sufficient or plausible unless rebutted.

#### **1.54 Virtue ethics**

Virtue ethics is concerned with providing an account of the **moral virtues**.

#### **1.55 Virtual epistemology**

Virtue epistemology is concerned with providing an account of the **intellectual virtues**.



## 2 Prior's paradox

### 2.1 Assumptions

Assumption	Description
<b>Dichotomy</b> assumption (A1)	All propositions may be categorised as either <b>ethical</b> or <b>non-ethical</b>
<b>Deductive</b> <b>argument</b> assumption (A2)	The <b>conclusion</b> of an argument cannot contain more information than its <b>premises</b> .

### 2.2 Proposition

Either tea drinking is common in England, or it ought to be the case that all New Zealanders are shot, formalised as  $p \vee Oq$ . According to the dichotomy assumption, the proposition is either **ethical** or **non-ethical**.

Horn 1	Horn 2
$p \vee Oq$ is ethical.	$p \vee Oq$ is non-ethical.
<p>If <math>p \vee Oq</math> is <b>ethical</b>, then:</p> <p>P1 (non-ethical): Tea drinking is common in England.</p> <p>C (ethical): Therefore, either tea drinking is common in England, or it ought to be the case that all New Zealanders are shot.</p>	<p>If <math>p \vee Oq</math> is <b>non-ethical</b>, then:</p> <p>P1 (non-ethical): Either tea drinking is common in England, or it ought to be the case that all New Zealanders are shot.</p> <p>P2 (non-ethical): Tea drinking is not common in England.</p> <p>C (ethical): Hence, it ought to be the case that all New Zealanders are shot.</p>

### 2.2.1 Dilemma

Whether we accept horn 1 or horn 2, we make **is-ought inferences** that are perfectly **valid**.

Horn	Classification of $p \vee Oq$	Premise set	Conclusion set
Horn 1: $p \vdash p \vee Oq$	$p \vee Oq$ is <b>ethical</b>	<b>Non-ethical</b>	<b>Ethical</b>
Horn 2: $p \vee Oq, \neg p \vdash Oq$	$p \vee Oq$ is <b>non-ethical</b>	<b>Non-ethical</b>	<b>Ethical</b>

Prior's paradox is a **dilemma without escape**. Since **Hume's Law**, even as a **one-way implication barrier**, is violated in every possible instance, it must be **false**. Hence, ethics is not **logically autonomous**.

### 3 Defending Hume's Law

1. Admit the **converse of the is-ought thesis**.
2. Exclude **contradictions** from the premise set  $\phi$ .  $\phi$  should be defined as a **consistent or contradiction-free** set of **descriptive propositions**.
3. Exclude **tautologies** from the conclusion set  $\psi$ .  $\psi$  should be defined as a **normative proposition** that is not already logically true.
4. Rule out **enthymematic arguments**. When the hidden premises of **enthymemes** are restored, the **premise set**  $\phi$  will have at least one **normative proposition**. Hence, these arguments will no longer be obvious counterexamples to Hume's Law.
5. Concede the **contraposition with "ought" implies the "can"** case. We should concede that **"cannot" implies "not obligatory"** yields a **special case** in which  $\phi \vdash \psi$ .
6. Rule out **Geach-style conditionals** as non-ethical propositions. It gives rise to embedded "ought" propositions of the form  $Op \rightarrow Oq$ .

With  $Op \rightarrow Oq$ , "ought" statements  $Op, Oq$  are being embedded into more complex logical structures, but there is no commitment to the truth or falsity of either  $Op$  or  $Oq$ .

7. Rule out **mixed propositions** from the premise set  $\phi$  and the conclusion set  $\psi$ . We can replace the **dichotomy** assumption with the **trichotomy** assumption, where all propositions may be categorised as either **ethical, non-ethical, or mixed**.

### 3.1 Issues with move 7

Mixed propositions have an indispensable role in ethical reasoning and argumentation. Purely normative propositions are rarely encountered in the real world, outside the philosopher's laboratory.

Examples include:

Proposition in natural language	Formal representation
If you refrain from helping the old lady across the road, then you ought to be blamed.	$\neg p \rightarrow Oq$
Either you help the old lady across the road, or you ought to be blamed for not doing so.	$p \vee Oq$
It is necessarily the case that if p, then it is obligatory that q.	$\Box(p \rightarrow Oq)$
It is necessarily the case that for all $x$ , then $Fx$ , then it is obligatory that $Gx$ .	$\Box\forall x(Fx \rightarrow Gx)$

### 3.2 Gerhard Schurz substitution

- If a **mixed conclusion**  $\phi$  is derivable from a **purely non-ethical premise set**  $\phi$ , then  $\psi$  is completely **O-irrelevant**.
- Apply the **O-restricted propositional substitution function**  $\sigma$ .
- Substitute  $r$  (any proposition whatsoever) for  $q$  on exactly those occurrences of  $q$  outside the scope of  $O$ , i.e.

$$p \text{ (non-ethical)} \vdash p \vee Oq, \text{ (mixed)} \xrightarrow{\text{Apply } \sigma} p \text{ (non-ethical)} \vdash p \vee Or \text{ (mixed)}$$

- The **O-restricted substitution** ( $\sigma$ ) can be made without compromising the validity of the argument; hence, the mixed conclusion  $p \vee Oq$  is completely **O-irrelevant** relative to the premise set.
- If a **mixed premise set**  $\phi$  is used to derive a **purely ethical conclusion**  $\psi$ , then  $\phi$  is completely **is-irrelevant**.
- Apply the **is-restricted propositional substitution function**  $\sigma'$ .
- Substitute  $r$  (any proposition whatsoever) for  $p$  on exactly those occurrences of  $p$  outside the scope of  $O$ , i.e.

$$p \vee Oq, \neg p \text{ (mixed)} \vdash Oq \text{ (ethical)} \xrightarrow{\text{Apply } \sigma'} r \vee Oq, \neg r \text{ (mixed)} \vdash Oq \text{ (ethical)}$$

- The **is-restricted substitution** ( $\sigma'$ ) can be made without compromising the validity of the argument; hence, the **mixed premise set**  $p \vdash Oq, \neg p$  is completely **is-irrelevant** relative to the conclusion.

### 3.3 Gibbard-Karmon-Singer semantics

Gibbard-Karmon-Singer semantics is just a way of determining whether a set of propositions will result in ethical conclusions or not. It works like this:

1. Consider the truth value of the propositions in a possible world, such as the actual world we live in.
2. Consider the truth value of the propositions and conclusions in an ethical standard.
3. Swapping the ethical standard for another ethical standard without changing the world.
4. If the truth value of the conclusions changes when you change the ethical standard, like the conclusions change from true to false, then the conclusions are ethical.
5. Otherwise, the conclusions are non-ethical, because the ethical standard being used is not relevant to the truth value of the conclusions.
6. If the conclusions are non-ethical, and you want to figure out which possible worlds the set of propositions will result in ethical conclusions, swap out the world for another one and repeat steps 2 to 5.

### 3.4 Shorter's position

- The **conclusion** of an argument may be of **some importance (ethically speaking)** in deriving certain moral duties only if it is arrived at in some other way than employing an **is-ought inference**.
- Hence, the **is-ought inference** is **not of importance (ethically speaking)**.
- We need to distinguish between the **seriousness of the conclusion arrived at (ethically speaking)** and the **seriousness of the is-ought inference** by which the conclusion is arrived at.

### 3.4.1 Tea drinking example

- P1 (non-ethical): Tea drinking is common in England.
- C (ethical): Therefore, either tea drinking is common in England or it ought to be the case that all New Zealanders are shot.

Step	Description
1	P1 is either true or false.
2	<p>If P1 is true, then the <b>is-ought inference</b> lends support to C. However, C will be of no help or use to us in deriving certain moral duties.</p> <p>If P1 is false, then the <b>is-ought inference</b> lends no support to C. However, if P1 is false and C is true, then we can derive the duty to shoot all New Zealanders, and C may be of some importance (ethically speaking).</p>
3	<p>Therefore, whether P1 is true or false, the <b>is-ought inference</b> is <b>useless</b>.</p> <p>It either <b>renders C ethically useless insofar as it supports C (when P1 is true) or does not support C (when P1 is false)</b></p>

### 3.4.2 Undertaker example

- P1: Undertakers are church officers.
- C: Therefore, undertakers ought to do whatever all church officers ought to do.
- P2: All church officers ought to  $\phi$ .

Step	Description
1	The <b>is-ought inference</b> either makes do without P2 or incorporates P2.
2	<p>If the <b>is-ought inference</b> makes do without P2, then it lends support to C. However, C will be <b>useless</b> without P2 as an undertaker can only derive a <b>concrete moral duty</b> with both P1 and P2.</p> <p>If the <b>is-ought inference</b> incorporates P2, then it lends support to <math>C'</math> (All undertakers ought to <math>\phi</math> rather than C (Undertakers ought to do whatever all church officers ought to do)). Therefore, C will become <b>useless</b> with P2.</p>



## 4 Modal concepts

Mode	Domain	Categories	Logical state of play in 1951
Mode 1 (Alethic)	Truth	Necessary, Possible, and Contingent	<b>Alethic modal logic</b> with the <b>modal operators</b> $\Box$ and $\Diamond$
Mode 2 (Epistemic)	Knowledge	Verified, Falsified, and Indeterminate	Minimal logical treatment
Mode 3 (Deontic)	Actions	Obligatory, Permissible, and Impermissible	Minimal logical treatment
Mode 4 (Existential)	Existence	Universal, Existential, and Empty	First-order predicate or quantificational logic (QL) with the quantifiers $\forall$ and $\exists$

## 5 Deontic logic

Deontic logic is a field of philosophical logic that is concerned with obligation, permission, and related concepts. The word "deontic" comes from the Greek word "deon" which means "that which is binding or proper".

### 5.1 Analogies between the deontic mode and the alethic mode

#### 5.1.1 Analogy 1

There are 2 operators.

- Alethic mode

$$\begin{aligned}\Box p &\stackrel{\text{def}}{=} \neg \Diamond \neg p \\ \Diamond &\stackrel{\text{def}}{=} \neg \Box \neg p\end{aligned}$$

$\Box$  and  $\Diamond$  are De Morgan duals.

- Deontic mode

$$\begin{aligned}Op &\stackrel{\text{def}}{=} \neg P \neg p \\ Pp &\stackrel{\text{def}}{=} \neg O \neg p\end{aligned}$$

$O$  and  $P$  are De Morgan duals.

#### 5.1.2 Analogy 2

There are **5 statuses** that can be defined in terms of **2 operators**.

Source (alethic mode)	Target (deontic mode)
It is <b>necessary</b> that $p$ ( $\Box p$ )	It is <b>obligatory</b> that $p$ ( $Op$ )
It is <b>possible</b> that $p$ ( $\Diamond p$ )	It is <b>permissible</b> that $p$ ( $Pp$ )
It is <b>impossible</b> that $p$ ( $\neg \Diamond p$ )	It is <b>impermissible</b> that $p$ ( $\neg Pp$ )
It is <b>non-necessary</b> that $p$ ( $\neg \Box p$ )	It is <b>omissible</b> that $p$ ( $\neg Op$ )
It is <b>contingent</b> that $p$ ( $\Diamond p \wedge \neg \Box p$ )	It is <b>optional</b> that $p$ ( $Pp \wedge \neg Op$ )

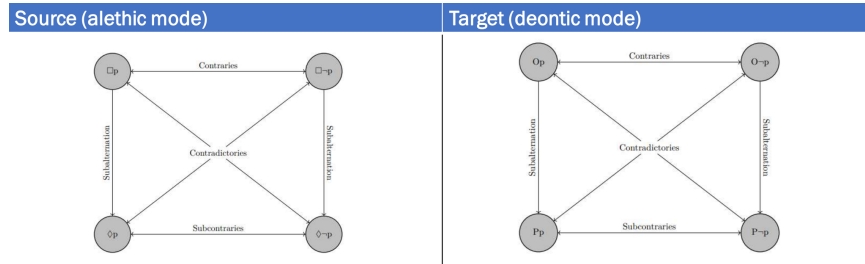
#### 5.1.3 Analogy 3

There are **5 statuses** that can be represented by a **threefold partition**.

Source (alethic mode)			Target (deontic mode)		
Possible ( $\Diamond p$ )			Permissible ( $Pp$ )		
Necessary ( $\Box p$ )	Contingent ( $\Diamond p \wedge \neg \Box p$ )	Impossible ( $\neg \Diamond p$ )	Obligatory ( $Op$ )	Optional ( $Pp \wedge \neg Op$ )	Impermissible ( $\neg Pp$ )
	Non-necessary ( $\neg \Box p$ )			Omissible ( $\neg Op$ )	

#### 5.1.4 Analogy 4

There is **1 square of opposition**.



- Two propositions are **contradictories** if and only if the **truth** of one implies the **falsity** of the other.
- Two propositions are **contraries** if and only if they cannot both be **true** but can both be **false**.
- Two propositions are **subcontraries** if and only if they cannot **both be false** but can both be **true**.
- Two propositions are in a **subalternation** relation if and only if the **truth** of the first proposition (**superaltern**) implies the **truth** of the second (**subaltern**) but NOT vice versa.

### 5.1.5 In summary

Analogy	Description
Analogy 1	There are <b>2 operators</b> . Source (alethic mode): $\Box, \Diamond$ Target (deontic mode): $O, P$
Analogy 2	There are <b>5 statuses</b> . Source (alethic mode): <b>necessity, possibility, impossibility, non-necessity, contingency</b> Target (deontic mode): <b>obligatoriness, permissibility, impermissibility, omissibility, optionality</b>
Analogy 3	These <b>5 statuses</b> can be represented by a <b>threefold partition</b> . Source (alethic mode): <b>necessity, contingency, impossibility</b> Target (deontic mode): <b>obligatoriness, optionality, impermissibility</b>
Analogy 4	There is <b>1 square of opposition</b> . Source: (alethic mode): <b>modal square of opposition</b> Target: (deontic mode): <b>deontic square of opposition</b>

### 5.2 Disanalogies

Disanalogy	Source (alethic mode)	Target (deontic mode)
Disanalogy 1	$\Box p \rightarrow p$ If $p$ is <b>true</b> across ALL possible worlds that are <b>accessible</b> , the $p$ must be <b>true</b> in the <b>actual world @</b> .	$\neg(Op \rightarrow p)$ It does not follow from the fact that the action described by $p$ is <b>obligatory</b> that the action is performed in the <b>actual world @</b> .
Disanalogy 2	$p \rightarrow \Diamond p$ If $p$ is <b>true</b> in the <b>actual world @</b> , then $p$ must be <b>true</b> in at least one possible world that is <b>accessible</b> .	$\neg(p \rightarrow Pp)$ It does NOT follow from the fact that the action is performed in the <b>actual world @</b> that is <b>permissible</b> .

### 5.3 Components

Component of the logical system	Elaboration
<b>Alphabet</b>	The <b>alphabet</b> of <b>deontic logic</b> is an extension of the <b>alphabet</b> of <b>propositional logic</b> to include the <b>deontic operators</b> $O$ and $P$ .
<b>Syntax</b>	The <b>syntax</b> of <b>deontic logic</b> is an extension of the <b>syntax</b> of <b>propositional logic</b> to handle <b>well-formed formulae (wffs)</b> containing at least one <b>deontic operator</b> .
<b>Semantics</b>	The <b>semantics</b> of <b>deontic logic</b> is of the form $\langle W, S, @ \rangle$ , where $W$ denotes a <b>set of worlds</b> , $S$ denotes a <b>binary relation of moral satisfaction</b> between worlds, and $@$ denotes the <b>actual world (privileged)</b> .
<b>Proof theory</b>	<p>The <b>proof theory</b> of <b>deontic logic</b> comprises a set of <b>definitions</b>, <b>axioms</b>, and <b>rules of inference</b>.</p> <p>This <b>proof theory</b>, with its reliance on <b>axioms</b>, is known as <b>Hilbert-style proof theory</b>.</p>

## 5.4 Proof theory

Definitions	Axioms	Rules of inference
<p>Definition 1:</p> <p><math>Op \stackrel{\text{def}}{=} \neg P \neg p</math></p>	<p>(T) All <b>tautologous well-formed formulae</b> from <b>propositional logic</b>.</p> <p>(NC) <math>\neg(Op \wedge O\neg p)</math></p> <p>It cannot be the case that both <math>p</math> and <math>\neg p</math> are <b>obligatory</b>.</p>	<p>(<math>\rightarrow_{E1}</math> or <i>modus ponens</i>) <math>p \rightarrow q, p \vdash q</math></p>
<p>Definition 2:</p> <p><math>Pp \stackrel{\text{def}}{=} \neg O \neg P</math></p>	<p>(K) <math>O(p \rightarrow q) \rightarrow (Op \rightarrow Oq)</math></p> <p>If performing the action described in <math>p</math> <b>commits</b> me to performing the action described in <math>q</math>, if <math>p</math> is <b>obligatory</b>, <math>q</math> will be <b>obligatory</b> too.</p> <p>(NEC) <math>\vdash p \rightarrow \vdash Op</math></p> <p>If <math>p</math> is <b>tautological</b>, then <math>Op</math> is also <b>tautological</b>.</p>	

## 5.5 Monotonicity of entailment (RM)

The monotonicity of entailment just means that if a sentence follows deductively from a given set of sentences, then it also follows deductively from any superset of those sentences. It is formalised as such:

$$(RM): (\vdash p \rightarrow q) \rightarrow (\vdash Op \rightarrow Oq)$$

### 5.5.1 Proof

1. Assume that  $\vdash p \rightarrow q$
2.  $\therefore (\vdash p \rightarrow q) \rightarrow \vdash O(p \rightarrow q)$
3.  $\therefore \vdash O(p \rightarrow q)$
4.  $\therefore \vdash O(p \rightarrow q) \rightarrow \vdash Op \rightarrow Oq$
5.  $\therefore \vdash Op \rightarrow Oq$
6.  $\therefore (\vdash p \rightarrow q) \rightarrow (\vdash Op \rightarrow Oq)$

## 5.6 Corollary of the monotonicity of entailment (Corollary of RM)

The corollary of the monotonicity of entailment is that if a given argument is deductively valid, it cannot become invalid by the addition of extra premises. It is formalised as such:

$$(\text{Corollary of RM}): \vdash Op \rightarrow O(p \vee q)$$

### 5.6.1 Proof

1.  $\vdash p \rightarrow (p \vee q)$
2.  $\therefore (\vdash p \rightarrow (p \vee q)) \rightarrow (\vdash Op \rightarrow O(p \vee q))$
3.  $\therefore \vdash Op \rightarrow O(p \vee q)$

## 5.7 Theorem T1

$$(T1) \vdash O(p \wedge q) \rightarrow Oq$$

### 5.7.1 Proof

1.  $\vdash (p \wedge q) \rightarrow q$
2.  $\therefore (\vdash (p \wedge q) \rightarrow q) \rightarrow (\vdash O(p \wedge q) \rightarrow Oq)$
3.  $\therefore \vdash O(p \wedge q) \rightarrow Oq$

## 5.8 Paradoxes of obligation

### 5.8.1 Paradox of the gentle murderer

Propositions (P):

1. It ought to be the case that A does not kill his mother.
2. If A does kill his mother, then it ought to be the case that A kills her gently.
3. A does kill his mother.

Proof:

1.  $O\neg p$ , where  $p$  denotes "A kills his mother" (from P1).
2.  $p \rightarrow Oq$ , where  $q$  denotes "A kills his mother gently" (from P2).
3.  $p$  (from P3).
4.  $\therefore Oq$  (from 2, 3, and  $\rightarrow_{E1}$  or *modus ponens*).
5.  $\therefore \vdash q \rightarrow p$  (from T).
6.  $\therefore (\vdash q \rightarrow p) \rightarrow (\vdash Oq \rightarrow Op)$  (from RM, uniformly substitute  $p$  for  $q$  and vice versa).
7.  $\therefore \vdash Oq \rightarrow Op$  (from 5, 6, and  $\rightarrow_{E1}$  or *modus ponens*).
8.  $\therefore Op$  (from 4, 6, and  $\rightarrow_{E1}$  or *modus ponens*).

Using the monotonicity of entailment, it follows that A should kill his mother, which is an odd thing to say.



### 5.8.2 Ross' paradox

Propositions (P):

1. It is obligatory that the letter is mailed.
2. Therefore, it is obligatory that the letter is mailed, or the letter is burnt.

Proof:

1.  $Op$ , where  $p$  denotes "the letter is mailed"
2.  $\vdash Op \rightarrow O(p \vee q)$  (Corollary of RM)
3.  $\therefore O(p \vee q)$ , where  $q$  denotes "the letter is burnt" (from 1, 2, and  $\rightarrow_{E1}$  or *modus ponens*)

It is odd to say that P1 and the corollary of RM entail an obligation that can be fulfilled by burning the letter (presumably an **impermissible** action).

### 5.8.3 The Good Samaritan paradox

Propositions (P):

1. It ought to be the case that A helps B, who has been robbed.
2. Therefore, it ought to be the case that B has been robbed.

Proof:

1.  $O(p \wedge q)$ , where  $p$  denotes "A helps B" and  $q$  denotes "B has been robbed" (from P1).
2.  $\vdash O(p \wedge q) \rightarrow Oq$  (T1)
3.  $\therefore Oq$  (from 1, 2, and  $\rightarrow_{E1}$  or *modus ponens*)

It is odd to say that from P1 and T1, it follows that B's being robbed is also obligatory.

## 5.9 Resolving the paradoxes

Response	Paradox of the gentle murderer	Ross' paradox	The good Samaritan paradox
Response 1: Distinguish between <b>non-derivatively obligatory</b> and <b>derivatively obligatory</b> actions.	Refrain from killing your mother ( <b>non-derivatively obligatory</b> )  Kill your mother (derived from RM and <b>impermissible</b> )	Mail the letter ( <b>non-derivatively obligatory</b> )  Burn the letter (derived from the <b>corollary</b> of RM and <b>impermissible</b> )	Help someone who is in need ( <b>non-derivatively obligatory</b> )  Rob the individual who has been robbed (derived from T1 and <b>impermissible</b> )
Response 2: Reject RM	RM gives rise to the <b>gentle murderer paradox</b> .	The <b>corollary of RM</b> gives rise to <b>Ross' paradox</b> .	T1, a theorem derived from RM, gives rise to the <b>good Samaritan paradox</b> .
Response 3: Introduce a <b>dyadic (2-placed)</b> version of deontic logic.	$O(\neg \text{murder} \mid T)$ ( <b>unconditional obligation</b> )  $O(\text{gentle murder} \mid \text{murder})$ ( <b>conditional obligation</b> if the <b>unconditional obligation</b> is violated)  We cannot derive an <b>unconditional obligation</b> to murder. $\not\models O(\text{murder} \mid T)$	$O(\text{mail} \mid \text{text has been written})$ ( <b>conditional obligation</b> )  We cannot derive an <b>obligation</b> to mail or burn the letter. $\not\models O(\text{mail} \vee \text{burn})$	$O(\text{help } B \mid B \text{ has been robbed})$ ( <b>conditional obligation</b> )  We cannot derive an <b>obligation</b> for B to have been robbed. $\not\models O(B \text{ has been robbed})$

## 5.10 Mimamsa deontic logic

- Classical deontic logic is a **monadic (1-placed)** system:

$$O(\_)$$

- At least some systems of **deontic logic** are **dyadic (2-placed)** systems:

$$O(\_|\_)$$

It ought to be the case that A helps B and B has been robbed. Denoting  $p$  as "A helps B" and  $q$  as "B has been robbed":

- Formal representation under **monadic deontic logic**:

$$O(p \wedge q)$$

- Formal representation under **dyadic deontic logic**:

$$O(p|q)$$

### 5.10.1 Dyadic deontic operator $O(\phi|\theta)$

The **dyadic deontic operator**  $O(\phi|\theta)$  is used in **dyadic deontic logic** to represent **conditional obligations**.  $\phi$  represents the **main argument** and  $\theta$  represents the **triggering condition**.

- It is necessarily the case that given  $p$ , it is obligatory that  $q$ .

$$\Box O(q|p)$$

- There is a **conditional obligation** that  $q$ , given  $p$ .

$$\Box O(q|p)$$

- There is an **unconditional obligation** that  $q$ , given that anything is the case.

$$O(q|T)$$

## 6 Inductive reasoning

$$\phi \vdash_P^I \psi$$

Where:

- $\phi$  is the set of premises, which potentially includes the knowledge base
- $\vdash$  represents that  $\psi$  is provable from  $\phi$
- $I$  refers to the inductive mode
- $P$  refers to the inferential mechanism of the calculus of probability
- $\psi$  is the set of conclusions derived from the set of premises

### 6.1 Newcomb's paradox (Newcomb's problem)

There is a reliable predictor, another player, and two boxes designated A and B. The player is given a choice between taking only box B or taking both boxes A and B. The player knows the following:

- Box A is transparent and always contains a visible \$1000.
- Box B is opaque, and its content has already been set by the predictor:
  - If the predictor has predicted that the player will take boxes A and B, then box B contains nothing.
  - If the predictor has predicted that the player will take only box B, then box B contains \$1,000,000.

The player does not know what the predictor predicted, or what box B contains while making the choice.

## 6.2 Philosophical principles

Principle	Description
The <b>principle of multiple explanations</b> . (Epicurus, c. 300 B.C.E.)	If multiple theories $H_1, H_2$ , and so on, are <b>consistent</b> with our observation $E$ , then we should retain ALL these theories $H_1, H_2$ , and so on.
The <b>uniformity of nature principle</b> . (Hume, 1739/40)	Nature is <b>sufficiently uniform</b> that <b>unobserved instances</b> in the <b>future</b> will resemble <b>observed instances</b> in the <b>future</b> will resemble <b>observed instances</b> in the <b>past</b> .
<b>Occam's razor principle</b> . (William of Ockham, 14th century C.E.)	Entities should NOT be multiplied beyond necessity.

### 6.3 Bayes' theorem

$$(BT)P(H|E) = \frac{P(E|H) \times P(H)}{P(E)}$$

Where:

- $BT$  refers to the finite set of rules of inferences, which is **Bayes' rule** or **Bayes' theorem**.
- $P$  is the probability of something
- $H$  is the hypothesis
- $E$  is the evidence
- $P(H|E)$  means the likelihood of  $H$  given  $E$ , it also refers to the **posterior probability**
- $P(E|H)$  means the likelihood of  $E$  given  $H$
- $P(H)$  refers to the **prior probability** of hypothesis  $H$  without ANY given conditions

### 6.4 Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \text{ if } P(B) \neq 0$$

Where:

- $P$  is the probability of something
- $A$  is an event
- $B$  is another event
- $\cap$  is the intersection of event  $A$  and  $B$ , i.e. the probability of event  $A$  and event  $B$  happening

## 6.5 Deductive vs inductive reasoning

Deductive reasoning	Inductive reasoning
We reason under <b>certainty</b> concerning propositions that are either <b>true</b> or <b>false</b> .	We reason under <b>uncertainty</b> concerning propositions in which we have <b>differing degrees of belief</b> .
<b>Deductive reasoning is monotonic.</b> If $\phi \vdash \psi$ , then adding more information $\lambda$ to the <b>premise set</b> $\phi$ will NOT invalidate out <b>conclusion</b> that $\psi$ .	<b>Inductive reasoning is non-monotonic.</b> Although it may be <b>true</b> that $\phi \vdash \psi$ , it need NOT be the case that $(\phi \wedge \lambda) \vdash \psi$ . $\lambda$ may constitute <b>new evidence</b> , forcing us to <b>retract or revise</b> our <b>conclusion</b> that $\psi$ .
<b>Deductive reasoning is non-ampliative.</b> <b>Deductive reasoning</b> unpacks the information content of the <b>premise set</b> $\phi$ , such that the information contained in the <b>conclusion set</b> $\psi$ is already present (albeit in implicit form) in $\phi$ .	<b>Inductive reasoning is ampliative.</b> The information in the <b>conclusion</b> that $\psi$ <b>exceeds and amplifies</b> the information content of the <b>premise set</b> $\phi$ .

## 6.6 Axioms of probability (Kolmogorov theorem)

### 6.6.1 Degrees of belief

The **degrees of belief** are constrained by a finite set of **axioms of probability**. Any probability function  $P$  must satisfy the following axioms:

Axiom	Description
K1 ( <b>non-negativity</b> )	$P(A) \geq 0$ in <b>sample space</b> $\Omega$ , where $P(A)$ is the <b>probability</b> of <b>outcome</b> $A$ .
K2 ( <b>normalisation</b> )	$P(\Omega) = 1$
K3 ( <b>addition rule</b> )	$P(A \cup B) = P(A) + P(B) - P(A \cap B)$ If $A$ and $B$ are <b>mutually exclusive</b> , then $(A \cap B) = \emptyset$ and $P(A \cap B) = 0$ . $\therefore P(A \cup B) = P(A) + P(B)$
K4 ( <b>complement rule</b> )	$P(\bar{A}) = P(\Omega) - P(A) = 1 - P(A)$

### 6.6.2 Ruled-out scenarios

A set of **outcomes** is **jointly exhaustive** if these **outcomes** encompass the entire **sample space**  $\Omega$ . In other words, at least one of these **outcomes** must occur.

The **axioms of probability** rule out the following scenarios:

Scenario	Axiom ruling out the scenario
The assignment of <b>negative probability values</b> to individual <b>outcomes</b>	K1 or <b>non-negativity</b>
The assignment of <b>probability values</b> to <b>jointly exhaustive and mutually exclusive outcomes</b> that sum to $> 1$ .	K2 or <b>normalisation</b>
The assignment of <b>probability values</b> to <b>jointly exhaustive outcomes</b> that sum to $< 1$ .	K2 or <b>normalisation</b>
The assignment of a <b>probability value</b> other than $1 - p$ to $\bar{A}$ , when an agent assigns a <b>probability value</b> $p$ to some <b>outcome</b> $A$ .	K4 or <b>complement rule</b>



## 6.7 Axioms of expected utility theory (von-Neumann-Morgenstern theorem)

### 6.7.1 Degrees of preference

The **degrees of preference** are constrained by the **axioms of expected utility theory**.

Any agent faced with a system  $U$  of alternative entities  $u, v, \dots$  must satisfy the following axioms:

Scenario	Axiom ruling out the scenario
VM1 ( <b>completeness</b> )	<p>For every <math>u</math> and <math>v</math>, one and only one of the following relations holds:</p> <p><math>u \succ v</math> (the agent <b>prefers</b> <math>u</math> to <math>v</math>)  <math>v \succ u</math> (the agent <b>prefers</b> <math>v</math> to <math>u</math>)  <math>u \sim v</math> (the agent is <b>indifferent</b> between <math>u</math> and <math>v</math>)</p> <p>Alternatively, for every <math>u</math> and <math>v</math>, either <math>u \succeq v</math> or <math>v \succeq u</math>.</p>
VM2 ( <b>transitivity</b> )	<p>For every <math>u, v</math> and <math>w</math>, <math>u \succ v</math> and <math>v \succ w</math> imply that <math>u \succ w</math>.</p> <p>Alternatively, for every <math>u, v</math>, and <math>w</math>, if <math>u \succeq v</math> and <math>v \succeq w</math>, then <math>u \succeq w</math>.</p>
VM3 ( <b>independence of irrelevant alternatives</b> )	<p>For every <math>u, v</math>, and <math>w</math>, suppose that <math>u \succeq v</math> and a <b>third irrelevant alternative</b> <math>w</math> is present.</p> <p>The <b>order of preference</b> of <math>u</math> over <math>v</math> (<math>u \succeq v</math>) holds, independently of the presence of absence of the <b>third irrelevant alternative</b> <math>w</math>.</p>
VM4 ( <b>continuity</b> )	<p>Let <math>L</math> denote a <b>lottery</b> whose 2 possible <b>outcomes</b> are <math>u</math> and <math>v</math>, <math>L \stackrel{\text{def}}{=} \{u, v\}</math>  <math>P(u) = \alpha</math> and <math>P(v) = 1 - \alpha</math>, where <math>0 &lt; \alpha &lt; 1</math></p> <p>For every <math>u, v</math>, and <math>w</math>, <math>v \succ w \succ u</math> implies the following:  The existence of an <math>\alpha</math> such that <math>w \succ L</math> when <math>1 - \alpha</math> or <math>P(v)</math> is <b>sufficiently small</b>.  The existence of an <math>\alpha</math> such that <math>w \sim L</math> at a certain value of <math>(1 - \alpha)</math>.  The existence of an <math>\alpha</math> such that <math>L \succ w</math> when <math>1 - \alpha</math> or <math>P(v)</math> is <b>sufficiently large</b>.</p>

### 6.7.2 Ruled-out scenarios

The **axioms of expected utility theory** rule out the following scenarios:

Scenario	Axiom ruling out the scenario
The agent prefers neither alternative to another nor remains <b>indifferent</b> between both alternatives.	VM1 ( <b>completeness</b> )
The agent preferring $u$ to $v$ and $v$ to $w$ but remaining <b>indifferent</b> between $u$ and $w$ .	VM ( <b>transitivity</b> )
The <b>decoy effect</b> .	VM3 ( <b>independence of irrelevant alternatives</b> )
The impossibility of an agent preferring lottery $L$ to $w$ , where $v \succ w \succ u$ and $L \stackrel{\text{def}}{=} \{u, v\}$	VM4 ( <b>continuity</b> )

- According to the **Cox-Jaynes model**, any system **reasoning under uncertainty** and in terms of **degrees of belief** will conform to the **axioms of probability**.
- Furthermore, if the **axioms of expected utility theory** are satisfied, then the agent is said to be **rational** and the **preferences** can be represented by a **utility function**.

## 6.8 Bayesian decision theory

- **Standard decision theory** addresses **individual decision-making** under **uncertainty**.
- **Standard decision theory** incorporates the **axioms of probability** (K1 - K4) and the **axioms of expected utility theory** (VM1 - VM4).

**Bayesian decision theory** incorporates **standard decision theory** and **Bayesian epistemology** (BT):

Step	Description
1	Identify $n$ alternative courses of action $\phi_1, \phi_2, \dots, \phi_n$ and their $m$ associated possible outcomes, where $\{M, n\} \in \mathbb{N}$ .
2	Characterise each action $\phi_i$ in terms of its possible outcomes $s_j$ , where $\{i, j\} \in \mathbb{N}, i \in (0, n]$ and $j \in (0, m]$ .
3	Assign <b>prior probabilities</b> to each outcome $P(s_j \phi_i)$ in accordance with K1 - K4.
4	Assign <b>utility values</b> to each outcome $U(s_j \cap \phi_i)$ in accordance with VM1 - VM4.
5	Gather evidence and <b>update probabilities</b> by applying Bayesian epistemology.
6	Multiplying the <b>updated probability</b> and the <b>utility values</b> each <b>outcome</b> $s_j$ relative to $\phi_i$ .
7	Sum the products across each section $\phi_i$ to determine its <b>expected utility</b> .
8	Select the action $\phi_i$ with the <b>highest expected utility</b> as the <b>morally right action</b> .

## 6.9 Jeffrey-Bolker theory

Jeffrey-Bolker's expected utility theory is an example of **evidential decision theory**. It relies on a boolean algebra  $\Omega$  that consists of:

Formal representation	Description
$A, B, C$ , etc.	<b>Propositions</b> as elements of $\Omega$ .
$\bar{A}$	<b>Negation</b> , such that if $A \in \Omega$ , then $\bar{A} \in \Omega$ .
$A \cup B$	<b>Disjunction</b> , such that if $A, B \in \Omega$ , then $A \cup B \in \Omega$ .
$\top$	<b>Tautology</b> .
$\perp$	<b>Contradiction</b> or negation of $\top$ .
$\succeq$	A coherent <b>preference order relation</b> over $\Omega'$ .

### 6.9.1 Strategy

Jeffrey aims to recommend a **Bayesian model of deliberation** and a corresponding **theory of preference**.

Move	Description
1	<p>Identify the <b>Bayesian principle of deliberation</b>. According to this principle, we rank actions <math>\phi_1, \phi_2, \dots, \phi_n</math> in order of <b>preference</b>.</p> <p>Given a <b>coherent preference ranking</b>, we can discover a pair of <b>probability and desirability assignments</b> (roughly corresponding to the <b>probability</b> and <b>utility value assignments</b>) to propositions describing the performance of these actions.</p>
2	<p>Introduce the <b>coherence</b> assumption. According to this assumption, the agent's <b>preference ranking</b> has the following properties:</p> <p>Property 1: <b>Coherence</b> There is an underlying set of <b>probabilities and desirabilities</b> that yield the <b>preference ranking</b> via the <b>Bayesian principle of deliberation</b>.</p> <p>Property 2: <b>Defeasibility</b> <b>Experience and reflection</b> constantly force the agent to <b>revise their agent preference ranking</b>.</p>
3	<p>Characterise the <b>desirability function (des)</b> and the <b>probability measure (prob)</b>. The <b>desirabilities (des)</b> of the <b>basic cases</b> may be any set of numbers, independent of the <b>probabilities</b> <math>\text{des}A &gt; 0</math> if <math>A</math> is <b>good</b>, <math>\text{des}A = 0</math> if <math>A</math> is <b>indifferent</b>, <math>\text{des}A &lt; 0</math> if <math>A</math> is <b>bad</b>.</p> <p>The <b>probabilities (prob)</b> of the <b>basic cases</b> may be any set of non-negative numbers that sum to 1 (<math>P(A) \geq 0, P(\Omega) = 1</math>).</p>

### 6.9.2 Example

Event	$L$ (live to age 65 or more)	$\bar{L}$ (die before age 65)
$S$ (smoke)	Best ( $S \cap L$ )	3 <sup>rd</sup> best ( $S \cap \bar{L}$ )
$\bar{S}$ (quit)	2 <sup>nd</sup> best ( $\bar{S} \cap L$ )	Worst ( $\bar{S} \cap \bar{L}$ )

### 6.9.3 Non-Bayesian deliberation of the example

According to the **sylogistic line of reasoning**:

- P1: Either  $L$  or  $\bar{L}$ .
- P2: If  $L$ , then  $S$  is more desirable than  $\bar{S}$ .
- P3: If  $\bar{L}$ , then (equally)  $S$  is more desirable than  $\bar{S}$ .
- C:  $\therefore S$  is more desirable than  $\bar{S}$  (**fallacious**)

This **fallacious line of reasoning** wrongly assumes that the **4 possible action-outcome pairs** are **equiprobable**:

$$P(S \cap L) = P(S \cap \bar{L}) = P(\bar{S} \cap L) = P(\bar{S} \cap \bar{L}) = 0.25$$

$$\begin{aligned}
 U(S) &= \frac{P(S \cap L) \times v(S \cap L) + P(S \cap \bar{L}) \times v(S \cap \bar{L})}{P(S \cap L) + P(S \cap \bar{L})} \\
 &= \frac{P(S \cap L) \times v(S \cap L) + P(S \cap \bar{L}) \times v(S \cap \bar{L})}{P(S)} \\
 &= P(L|S) \times v(S \cap L) + P(\bar{L}|S) \times v(S \cap \bar{L}) \\
 &= \frac{0.3}{0.5} \cdot 100 + \frac{0.2}{0.5} \cdot -90 \\
 &= 60 - 36 \\
 &= 24
 \end{aligned}$$

#### 6.9.4 Bayesian deliberation of the example

Suppose that a **probability measure**  $P$  allows us to assign the following **probability values**:

$$P(\Omega) = P(S) + P(\bar{S}) = 1$$

$$P(S) = P(\bar{S}) = 0.5$$

$$P(S) = P(S \cap L) + P(S \cap \bar{L})$$

$$P(\bar{S}) = P(\bar{S} \cap L) + P(\bar{S} \cap \bar{L})$$

$$P(S \cap L) = 0.3$$

$$P(S \cap \bar{L}) = 0.2$$

$$P(\bar{S} \cap L) = 0.4$$

$$P(\bar{S} \cap \bar{L}) = 0.1$$

Suppose a **desirability measure**  $v$  allows us to assign the following **desirability values**:

$$v(S \cap L) = 100 \text{ (Best)}$$

$$v(\bar{S} \cap L) = 70 \text{ (2nd best)}$$

$$v(S \cap \bar{L}) = -90 \text{ (3rd best)}$$

$$v(\bar{S} \cap \bar{L}) = -100 \text{ (Worst)}$$

$$\begin{aligned} U(\bar{S}) &= P(L|\bar{S}) \times v(\bar{S} \cap L) + P(\bar{L}|\bar{S}) \times v(\bar{S} \cap \bar{L}) \\ &= \frac{0.4}{0.5} \cdot 70 + \frac{0.1}{0.5} \cdot -100 \\ &= 56 - 20 \\ &= 36 \end{aligned}$$

$$\therefore U(\bar{S}) > U(S)$$

$\therefore \bar{S} \succ S$  (**quitting** is preferable to **continuing to smoke**)

## 7 Issues with standard decision theory

### 7.1 Consistency versus responsibility

- The **axioms of probability** (K1 - K4) and the **axioms of expected utility theory** (VM1 - VM4) provide important constraints on the assignment of **probability and utility values**.
- However, consistency with these axioms is insufficient to ensure **responsible decision-making**.
- For example, think about an individual who is a **flat earth theorist** who thinks that the **flat earth theory** is 100% correct and all other theories are wrong. Such an individual is **Kolmogorov-consistent**, as his beliefs conform to the **axioms of probability**, but is **epistemically irresponsible**.
- Another example would be someone who prefers **genocide to murder** and **murder to a walk in the hills**. Such a person is **von-Neumann-Morgenstern-consistent** as his **degrees of preference** conform to the **axioms of expected utility theory**, but is **morally irresponsible**.

### 7.2 Cognitive biases

- Humans are not perfect, and hence we all suffer from cognitive biases.
- One example is the **decoy effect**, which is a **cognitive bias** in which consumers demonstrate a **shift in preferences** between two options when presented with a **third option that is asymmetrically dominated**.
- An example of this effect would be having 3 options for a product, but 1 option is worse than all other options, such as the following:

Option	Description	Price
Option 1	Online subscription for a newspaper	59.00
Option 2	Print subscription for a newspaper	125.00
Option 3	Print and online subscription for a newspaper	125.00

- When option 2 is removed, some people may change their preference from option 3 to option 1, which violates the **independence of irrelevant alternatives**.



### 7.3 Deontological decision-making

- There are difficulties in modelling **deontological decision-making**.
- We denote the utility value of a **prohibited or impermissible** outcome as  $-\infty$  and the utility value of an **obligatory** outcome as  $+\infty$ .
- This sets **prohibitions** and **obligations** apart from other **actions**, as their associated **outcomes** have **absolute maximum or minimum expected utility**.

Problem	Description
Problem 1: Swamping out of <b>probability values</b> .	<p>An <b>infinite utility or disutility value</b> completely <b>swamps out any probability value</b> associated with an outcome.</p> <p>If <b>killing</b> is <b>morally prohibited</b>, then the outcome of <b>murder</b> will be assigned the utility value of <math>-\infty</math>, hence, any actions that may lead to murder, no matter how unlikely, will presumably also be <b>morally prohibited</b>.</p>
Problem 2: All <b>prohibited</b> or <b>obligatory</b> actions are on par	<p><b>Murder</b> is no better or worse than <b>genocide</b>.</p>
Problem 3: Violation of the <b>continuity</b> axiom.	<p>Any preservation of the continuity axiom could be questioned on <b>deontological grounds</b>. Suppose that <math>w</math> is <b>forbidden</b> and <math>u</math> and <math>v</math> are <b>permissible</b>, such that <math>v \succ u \succ w</math>.</p> $L \stackrel{\text{def}}{=} \{v, w\}$ <p>According to the continuity axiom, there will be a probability such that the permissible action <math>u</math> is as good as a lottery involving another <b>permissible and more preferable action</b> <math>v</math> and a <b>prohibited action</b> <math>w</math>.</p> <p>However, according to <b>deontology</b>, any lottery involving a <b>prohibited action</b>, no matter how unlikely, cannot be ranked as on par with an ordinary <b>permissible action</b>.</p>

## 7.4 Collective decision-making and voting paradoxes

Using the preference data below:

Individual	Preference
1	$A \succ C \succ D \succ B$
2	$B \succ C \succ D \succ A$
3	$D \succ A \succ C \succ B$
4	$A \succ B \succ D \succ C$
5	$D \succ A \succ C \succ B$

### 7.4.1 Condorcet method

The Condorcet method conducts **pairwise comparisons**, and the winner is the choice that wins all the head-to-head matchups.

Head-to-head match-up	Winner	Score
A versus B	A	4-1
A versus C	A	4-1
A versus D	D	2-3
B versus C	C	2-3
B versus D	D	2-3
C versus D	D	2-3

### 7.4.2 Borda count

The Borda count method assigns points to the preference order of the choices. In this example, 0 points are awarded to the last choice, 1 point to the second-last choice, 2 points to the second choice, and 3 points to the first choice.

Choice	Number of points	Score
A	$3 + 0 + 2 + 3 + 2$	10 pts
B	$0 + 3 + 0 + 2 + 0$	5 pts
C	$2 + 2 + 1 + 0 + 1$	6 pts
D	$1 + 1 + 3 + 1 + 3$	9 pts

### 7.4.3 Problems

Problem	Description
Problem 1: Different methods yield different winners	D is the <b>Condorcet winner</b> , whereas A is the winner under the <b>Borda count</b> method.
Problem 2: <b>Voting paradoxes</b>	<p><b>Individual preferences</b> may be <b>non-cyclic</b> and consistent with the axioms VM1 - VM4.</p> <p>However, <b>collective preferences</b> could be <b>cyclic</b>. This is known as the <b>Condorcet paradox</b> or the <b>voting paradox</b>.</p> <p>Individual 1: <math>A \succ B \succ C</math>  Individual 2: <math>B \succ C \succ A</math>  Individual 3: <math>C \succ A \succ B</math></p> <p>Collectively: <math>A \succ B \succ C \succ A</math> (<b>paradoxical</b>)</p> <p>While <b>individual preferences</b> may obey VM1 (<b>completeness</b>) and VM2 (<b>transitivity</b>), it by no means follows that <b>collective preferences</b> will obey VM1 and VM2.</p> <p><math>P(\text{Condorcet paradox}) \approx 9\%</math> (low)</p>

### 7.4.4 Arrow's impossibility theorem

When voters have more than 3 distinct choices, **no ranked voting electoral system** can **convert the ranked individual preferences into ranked collective preferences** while satisfying the following conditions:

Condition	Description
1	<b>Pareto efficiency</b>
2	<b>Independence of irrelevant alternatives</b> (VM3)
3	<b>Unrestricted domain</b>
4	<b>Absence of dictatorship</b>

Alternatively, there is no constitution by which **ranked individual preferences** can be aggregated into **ranked collective preferences**, while satisfying **basic fairness criteria**, unless there is a **dictatorship** (not condition 4).

## 8 Game theory

### 8.1 Comparison with decision theory

Decision theory	Game theory
The concern is with <b>individual decision-making</b> .	The concern is with <b>interdependent decision-making</b> .
An individual's choice is neither affected by nor affecting the choices of other individuals.	An individual's choice affects the choices of other individuals. Each individual also has to consider the <b>utility functions</b> of other individuals and how they will affect the choices of other individuals and the overall outcome.

### 8.2 Prisoner's dilemma

- There are two suspects, P1 and P2.
- The district attorney believes that P1 and P2 are **guilty of a crime** but does not have sufficient evidence to convict them.

Each of P1 and P2 has 2 strategies:

Strategy	Description
1	Do not confess the crime (cooperate)
2	Confess the crime to the police (defect)

P1 and P2 are confronted with 4 options relative to various **possible strategy combinations**:

Option 4: If both P1 and P2 <b>do not confess</b> , then they will each get <b>1 year in prison</b> .	Option 3: If P2 <b>confesses</b> and P1 does not, P2 will get <b>3 months in prison</b> and P1 will get <b>10 years in prison</b> .
Option 2: If P1 <b>confesses</b> and P2 does not, P1 will get <b>3 months in prison</b> and P2 will get <b>10 years in prison</b> .	Options 1: If both P1 and P2 <b>confess</b> , then they will each get <b>8 years in prison</b> .

### 8.3 Table of utility values

		P2	
		STRATEGY 1 ( $\neg$ confess)	STRATEGY 2 (confess)
P1	STRATEGY 1 ( $\neg$ confess)	(-5, -5)	(-10, -2)
	STRATEGY 2 (confess)	(-2, -10)	(-7, -7)

- The cells coloured in **green** denotes a **Pareto-optimal** state of affairs.
- The cell coloured in **yellow** denotes the **Nash equilibrium**.
- If P1 and P2 are allowed to communicate and bargaining is **cost-free**, then P1 and P2 could agree to **cooperate** and not **confess**.
- Hence, they could make a **Pareto-efficient move** indicated by the arrow ( $\rightarrow$ )
- The **Nash equilibrium** (coloured in **yellow**) arises because P1 and P2 behave as **straightforward maximisers**.
- However, P1 and P2 have reasons to become **constrained maximisers**.

#### 8.4 Modification of rationality assumption

- The original **rationality** assumption, which is **straightforward maximisation**, is as such: It is **rational** to choose the course of action with the **maximum expected utility**.
- The modified **rationality assumption**, is as such: It is **rational** to be disposed to **constrained maximisation** and **cooperate** in **prisoner's dilemma-type scenarios**.

#### 8.5 Possible strategies for the game

- Random, which is to choose to cooperate 50% of the time.
- Tit-for-tat (TFT), which is to choose to cooperate on the first move, and choose your opponent's last move as your next move.
- Suspicious tit-for-tat (STFT), which is to choose to defect on the first move, and choose your opponent's last move as your next move.
- Tit for two tats (TF2T), which is to choose to cooperate on the first two moves, then choose to cooperate as the next move, unless your opponent chooses to defect for 2 moves. When your opponent stops choosing to defect, then choose to cooperate.

#### 8.6 Axelrod's tournaments

- The prisoner's dilemma was originally introduced as a **2-player game**, but it was later embedded by Axelrod into **round-robin tournaments**.
- Each program was pitted against the rest of the field.
- The aim of these tournaments was to learn about how to **choose effectively** in an **iterated prisoner's dilemma**.

### 8.6.1 Properties of successful strategies

Property	Description
Be <b>nice</b>	Choose to cooperate on the first move. For example, Cooperative, TFT, TF2T.
Be <b>forgiving</b>	Do not immediately retaliate if your opponent chooses to defect in a move. For example, Cooperative, TF2T.
Be <b>prepared to retaliate if necessary</b> .	You must be prepared to choose to defect at some point if your opponent keeps choosing to defect. For example, TFT, STFT, TF2T

- What accounts for **tit-for-tat (TFT)**'s success is its combination of being **nice, retaliatory, forgiving, and clear**.
- Its **niceness** prevents it from getting into **unnecessary trouble**.
- Its **retaliation** discourages the other side from persisting whenever **defection** is tried.
- Its **forgiveness** helps restore **mutual cooperation**.
- Its **clarity** makes it intelligible to the other player, thereby eliciting **long-term cooperation**.

## 8.7 Rapoport et al.'s objections to the Axelrod tournaments

### 8.7.1 Objection 1

The choice of **tournament format**:

- In a **knockout tournament**, top-ranked contestants at each stage progress to the next stage.
- As the tournament continues, the number of competitors decreases.
- In a **round-robin tournament**, each contestant competes with each of the others an equal number of times.
- Axelrod chose the **single-stage round-robin format** for his tournaments.
- He provided no **justification** for this choice of **tournament format**.
- The **single-stage round-robin format** becomes **impractical** when the **number of contestants is large**, although this problem disappears when the tournament is run on a computer.

### 8.7.2 Objection 2

The choice of **criterion for determining success**:

- The **criterion for determining success** involved **maximising the number of points won across all dyadic interactions**.
- Axelrod chose this **criterion of success**, but once again provided no **justification** for this choice of criterion for determining success.
- Most of the programs were not designed to **maximise the total number of points**.
- **Tit-for-tat (TFT)** can never win any particular **iterated prisoner's dilemma game**.
- **Tit-for-tat (TFT)** can never achieve a **positive point difference** against any other program.



### 8.7.3 Objection 3

The choice of **payoff structure**:

- The  $2 \times 2$  **prisoner's dilemma payoff matrix** had conventional values, where  $T$  denotes **sole defection**,  $R$  denotes **joint cooperation**,  $P$  denotes **joint defection**, and  $S$  denotes **sole cooperation**.
- The values in this matrix are  $(T, R, P, S) = (5, 3, 1, 0)$ .
- Axelrod chose this **payoff structure** but again provided no **justification** for this choice of **payoff structure**.

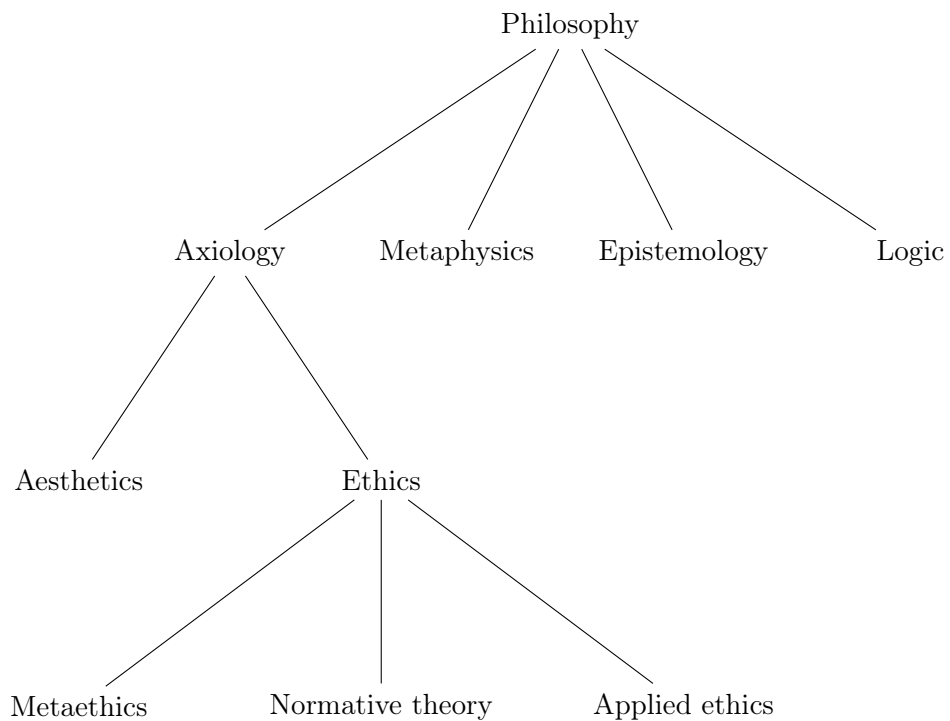
### 8.7.4 Conclusion

- Once the **2-player prisoner's dilemma** is embedded into a **tournament**, decisions have to be made about the **tournament format** (objection 1), **criteria for determining a winner** (objection 2), and **payoff structure** (objection 3).
- However, Axelrod has provided no **justification** for his choices of **tournament format**, **criterion of determining success**, and **payoff structure**.
- Hence, the policy recommendations about the effectiveness of **tit-for-tat (TFT)** should be qualified, i.e. the recommendations cannot be generalised.

## 9 Ethics

**Ethics** and **moral philosophy** are children nodes of the branch of **axiology**. Children nodes of **ethics** and **moral philosophy** include:

- **Metaethics**, **normative theory**, and **applied ethics**.



## 9.1 Branches of ethics

Branch of ethics and moral philosophy	Central question	Elaboration
Metaethics (2 <sup>nd</sup> -order theory)	What is <b>morality</b> ?	<b>Metaethics</b> is concerned with the <b>status, foundation, and scope of moral facts, values, properties, and terms.</b>
Normative theory (1 <sup>st</sup> -order theory)	What is <b>moral</b> (in <b>general</b> )?	<b>Normative theory</b> is concerned with the articulation of <b>moral standards and rules of conduct.</b>
Applied ethics (praxis)	What is <b>moral</b> (in <b>specific, controversial issues</b> )?	<b>Applied ethics</b> is concerned with the application of <b>philosophical theory</b> to <b>practical problems.</b>

## 9.2 Normative neutrality requirement

Branch of ethics and moral philosophy	Description
Normative theory (1 <sup>st</sup> -order theory)	There is <b>NO normative neutrality requirement.</b> <b>Normative theory</b> must have <b>substantive moral commitments.</b>
Metaethics (2 <sup>nd</sup> -order theory)	There is a <b>normative neutrality requirement.</b> <b>Metaethics</b> can have <b>classificatory moral commitments</b> but must avoid <b>substantive moral commitments.</b>

### 9.3 Forcehimes' collapse argument

- P1: There is a requirement for **normative neutrality** in **metaethics**.
- P2: Such a **breach of normative neutrality** is **inevitable**.
- Conclusion: Hence **metaethical theories (2<sup>nd</sup>-order)** turn out to be **normative theories (1<sup>st</sup>-order)** in disguise.

### 9.4 Normative theory

The different approaches to **normative theory (1<sup>st</sup>-order)** give rise to different **substantive moral commitments** and different **moral standards** and **rules of conduct**. The different approaches include:

Degree of particularity	Approach	Elaboration
General	Approach 1: <b>High moral theory</b>	<b>Consequentialism</b> (camp 1), <b>deontology</b> (camp 2), <b>virtue ethics</b> (camp 3)
↓	Approach 2: <b>Mid-level theory</b>	<b>Autonomy</b> (principle 1), <b>beneficence</b> (principle 2), <b>non-maleficence</b> (principle 3), <b>justice</b> (principle 4)
↓	Approach 3: <b>Casualism or case-based reasoning</b>	A <b>bottom-up approach</b> in which <b>moral principles</b> and <b>moral theories</b> emerge from case-based moral judgments.
Particular	Approach 4: <b>Narrative ethics</b>	We use stories to make sense of our experiences.

## 9.5 High moral theory

From a set of  $n$  alternative courses of action, where  $i, n \in \mathbb{N}$  and  $i \in (0, n]$ :

Camp	Description
<b>Consequentialism</b>	$\phi_i$ -ing is <b>morally right</b> if and only if it <b>maximises the good</b> , where the <b>good</b> is defined in terms of <b>some theory of the good T</b> .
<b>Deontology</b>	$\phi_i$ -ing is <b>morally right</b> if and only if it has <b>intrinsic moral worth</b> .
<b>Virtue ethics</b>	$\phi_i$ -ing is <b>morally right</b> if and only if it is the <b>best action</b> (in terms of <b>virtues and vices</b> ) that a <b>virtuous agent</b> might perform in the circumstances.

### 9.5.1 Advantages

Advantage	Elaboration	Substantiation
Advantage 1	<b>High moral theory</b> can provide <b>structured and systematic moral guidance</b> .	Camps 1 to 3 provide us with <b>moral standards</b> and <b>rules of conduct</b> for identifying the <b>morally appropriate action</b> ( $\phi_i$ -ing) from $n$ alternative courses of action.
Advantage 2	<b>High moral theory</b> can yield <b>moral standards</b> that can yield <b>moral standards</b> that can help us to achieve <b>consistency</b> and <b>coherence</b> in our moral lives.	<p>For camp 1, there is the transitivity rule, which states if <math>\phi_1</math> is better than <math>\phi_2</math> and <math>\phi_2</math> is better than <math>\phi_3</math>, then <math>\phi_1</math> must be better than <math>\phi_3</math>.</p> <p>For camp 2, there is the no-contradiction rule, which states that one and the same action <math>\phi_i</math> cannot be both <b>obligatory</b> and <b>impermissible</b>.</p> <p>For camp 3, there is the doctrine of the mean, which states that the <b>virtues</b> are a <b>mean</b> between the <b>vices of defect</b> and <b>excess</b>.</p>
Advantage 3	<b>High moral theory</b> often has the relevant tools and resources for <b>moral justification</b> .	<p>Camp 1 delivers <b>evaluative claims</b> in terms of the <b>maximisation of the good</b>.</p> <p>Camp 2 delivers <b>deontic verdicts</b> in terms of <b>duties, rights, and obligations</b>.</p> <p>Camp 3 delivers <b>virtue-ethical judgments</b> in terms of the <b>language of virtues and vices</b>.</p>

### 9.5.2 Disadvantages

1. How do we choose a **moral theory** from the **competing alternatives**?
  - Even if we do make a choice, how do we **justify** that choice?
  - Individuals with **different theoretical starting points** must still agree on a **similar set of principles**.
  - Hence, it has been argued that we should move to **mid-level theory** and a more **principle-based approach**.
2. How do we navigate **disagreement** within the ranks of each camp?
  - For camp 1, **act** versus **rule based** forms of consequentialism, **maximising** versus **satisficing** forms of consequentialism.
  - For camp 2, **monistic** versus **pluralistic** forms of **deontology**, **agent** versus **patient-centred** forms of deontology.
  - For camp 3, **eudaimonist** versus **non-eudaimonist** forms of **virtue ethics**.
3. **High moral theory** may be too **ill-equipped** to deal with **practical decision-making at the concrete level**.
  - When the **applied ethical problems** are **complex**, how likely is that **high moral theory** will be **sufficiently fine-grained** to generate responses?
  - For a move away from a **top-down approach** and toward a greater degree of **particularity**, it has been argued that we should favour **casuistry** or **case-based reasoning** or **narrative ethics** instead.

## 10 Consequentialism

- From a set of  $n$  alternative courses of action,  $\phi_i$ -ing is **morally right** if and only if it **maximises the good**, where  $i, n \in \mathbb{N}$ ,  $i \in (0, n]$ , and the **good** is defined in terms of **some theory of the good T**.
- The **core consequentialist commitment** is **maximising the good** (however the good might ultimately be defined).

### 10.1 Consequentialist decision-making

1. Compare the relative merits of  $n$  alternative courses of action  $\phi_1, \phi_2, \dots, \phi_n$ , where  $n \in \mathbb{N}$ .
2. **Evaluate** these  $n$  courses of action in terms of whether they **maximise the good**.
  - $\phi_1$  **maximises the good**.
  - $\phi_2$  does not **maximise the good**.
3. Arrive at the **decision outcome**. The **morally right** action is the one that **maximises the good** (for instance,  $\phi_1$ ).



## 10.2 Theory of the Good

- The **good to be maximised** is determined in terms of a **theory of the good T**.
- However, there are **multiple theories of the good**.
- Hence, one has not adopted any particular moral system in adopting **consequentialism** unless one says what the **good** is.

Candidate **theories of the good** include:

1. The **good** is defined as things which are **pink with yellow trimmings**.
  - This is meant to be a joke.
2. The **good** is defined as things which ought to be **maximised**.
  - This is possibly **trivial**, as the **core consequentialist commitment** becomes "that which **ought to be maximised**, ought to be **maximised**".
3. The **good** is defined as things that **facilitates self-interest**.
4. The **good** is defined as things that **facilitates human pleasure and happiness**.
5. The **good** is defined as that which is best understood as a **plurality of goods** (happiness, justice, fairness, and so on).

### 10.3 Types of consequentialism

Distinct types of **consequentialism** can be generated from **multiple theories of the good**:

Type of consequentialism	Core consequentialist commitment	Theory of the good	Consequentialist outcome
<b>Ethical egoism</b>	The <b>good</b> ought to be maximised.	<b>Egoistic</b> theory of the good.	Select the action $\phi_i$ that <b>maximises the good</b> .
<b>Utilitarianism</b>		<b>Hedonistic</b> theory of the good.	
<b>Pluralistic consequentialism and ideal utilitarianism</b>		<b>Pluralistic</b> theory of the good.	

### 10.4 Types of utilitarianism

For both **act** and **rule-based utilitarianism**, the **good** is what **facilitates human pleasure, happiness, and utility-based considerations**.

Type of utilitarianism	Theory of the good	Description
<b>Act-based utilitarianism</b>	The <b>good</b> is defined as that which <b>facilitates human pleasure and happiness</b> .	$\phi_i$ -ing is morally right if and only if it <b>maximises the good</b> .
<b>Rule-based utilitarianism</b>		$\phi_i$ -ing is morally right if and only if it is in accordance with a certain <b>set of rules R</b> that has been selected for its <b>good consequences</b> .

## 10.5 Hedonism

- Both **act** and **rule-based utilitarianism** are characterised in terms of a **hedonistic** theory of the good.
- Psychological hedonism: Only **pleasure (happiness)** or **pain (unhappiness)** motivates us.
- Ethical hedonism: Only **pleasure (happiness)** has **value** and only **pain (unhappiness)** has **disvalue**.

According to the **greatest happiness principle**:

- Happiness is defined as **pleasure** and the **absence of pain** ( $\text{pleasure} \wedge \neg \text{pain}$ )
- Unhappiness is defined as **pain** and the **privation of pleasure** ( $\text{pain} \wedge \neg \text{pleasure}$ )

Types of **ethical hedonism** include:

- **Quantitative hedonism**, which states that the **quantity of pleasure (happiness)** that matters.
- **Qualitative hedonism**, which states that the **quality of pleasure (happiness)** that matters.

## 10.6 Benthamite utilitarianism

Benthamite utilitarianism is a **traditional account of utilitarianism** that can be characterised in terms of the following:

Core consequentialist commitment	Hedonistic theory of the good	Consequentialist outcome
The good ought to be maximised.	Psychological hedonism	Select the action $\phi_i$ that <b>maximises the good</b> .
	Ethical hedonism	
	Quantitative hedonism	

## 10.7 Hedonic calculus

- Bethamite utilitarianism relies on a **hedonic calculus** or **felicific calculus**.
- The **hedonic calculus** is an algorithm, formulated by Bentham, for **calculating the total quantity of pleasure (happiness)** that an action  $\phi_i$  is likely to cause.

The **variables** in the **hedonic calculus** include:

Variable	Description
Intensity	How <b>strong</b> the <b>pleasure</b> or <b>pain</b> is.
Duration	How <b>long</b> the <b>pleasure</b> or <b>pain</b> lasts.
Probability	How <b>likely</b> the <b>pleasure</b> or <b>pain</b> is to be the result of $\phi_i$ -ing.
Propinquity or remoteness	How close the sensation of <b>pleasure</b> or <b>pain</b> is to be the result of $\phi_i$ -ing.
Fecundity	How likely $\phi_i$ is to lead to further <b>pleasures</b> or <b>pains</b> .
Purity	How much <b>intermixture</b> there is between <b>pleasure</b> or <b>pain</b> and other sensations.

## 10.8 Issues with quantitative hedonism

- Suppose that our choice is between **playing push-pin** ( $\phi_1$ ) and **reading poetry** ( $\phi_2$ ).
- Suppose that the **total net utility value** of  $\phi_1$  equals the **total net utility value** of  $\phi_2$ .
- Quantitative hedonism states that only the **quantity of pleasure (happiness)** matters.
- Hence, both **Benthamite utilitarianism** and its **machine-based implementation** *Jeremy* will concur that **playing push-pin** ( $\phi_1$ ) is as good as **reading poetry** ( $\phi_2$ ).
- There is a danger that **quantitative hedonism** might lead us to over-value **bestial, unsophisticated, lower-quality, and debauched pleasures**.
- However, human beings are able to distinguish between **lower-quality pursuits** such as  $\phi_1$  and **higher-quality intellectual pursuits** such as  $\phi_2$ .

### 10.8.1 Possible responses

Response	Justification	Elaboration
Defend <b>quantitative hedonism</b> .	The <b>pleasures</b> associated with <b>reading poetry</b> ( $\phi_2$ ) are <b>more probable, more durable, more fecund or more likely to lead to further pleasures</b> , and <b>purser</b> (unlikely to be mixed with pain).	The <b>quality of pleasure</b> can still be reduced to <b>quantitative considerations</b> . Hence, we can still retain <b>quantitative hedonism</b> .
Ditch <b>quantitative hedonism</b> in favour of <b>qualitative hedonism</b> .	<b>Qualitative hedonism</b> gives us automatic reasons to favour <b>higher-quality intellectual pursuits</b> such as $\phi_2$ over <b>lower-quality pursuits</b> such as $\phi_1$ .	

## 10.9 Qualitative hedonism

- The standard view of qualitative hedonism is that a **higher-quality pleasure** will be preferred to **any amount** of a **lower-quality pleasure**.
- Let  $\phi_i$  denote **playing push-pin** and let  $\phi_2$  denote **reading poetry**.
- Suppose that the total net utility value of  $\phi_1 = 5,000,000$  units and the total net utility value of  $\phi_2 = 1$  unit.

$$U(\phi_2) < U(\phi_1)$$

However, since a higher quality pleasure is preferred to a lower quality pleasure:

$$\phi_2(\text{higher-quality}) \succ \phi_2(\text{lower-quality})$$

- We could postulate an **infinite superiority** of **higher-quality** over **lower-quality pleasures**.
- However, this will reduce **qualitative** to **quantitative considerations**.
- For the **non-standard view** of qualitative hedonism:
  - **Qualitative hedonism** is not **quantitative hedonism** in disguise.
  - Rather, **quantity** and **quality** are two distinct properties of **pleasure**.
  - We may sometimes have to make **trade-offs** between a **lower-quantity** of a **higher-quality pleasure** and a **higher quantity** of a **lower-quality pleasure**.

### 10.9.1 Dilemma regarding qualitative hedonism

Argumentative constituent	Elaboration
P1: Either the <b>quality of pleasure</b> contributes to the <b>total net utility value</b> in the same manner as the other <b>quantitative variables</b> in the <b>hedonic calculus</b> , or it does not.	The other <b>quantitative variables</b> in the <b>hedonic calculus</b> are <b>intensity, duration, probability, propinquity</b> or <b>remoteness, fecundity, and purity</b> .
P2 (horn 1): If the <b>quality of pleasure</b> contributes to the <b>total net utility value</b> in the same manner as the other <b>quantitative variables</b> , then <b>qualitative hedonism</b> will turn out to be <b>quantitative hedonism</b> in disguise. P3 (horn 2): If the <b>quality of pleasure</b> does not contribute to the <b>total net utility value</b> in the same manner as the other <b>quantitative variables</b> , the <b>qualitative hedonism</b> will be <b>inconsistent</b> .	This is the same problem confronting the <b>standard view of qualitative hedonism</b> , when it postulates an <b>infinite superiority</b> of <b>lower-quality pleasure</b> over <b>higher-quality pleasure</b> . According to <b>ethical hedonism</b> , only <b>pleasure (happiness)</b> has <b>value</b> and only <b>pain (unhappiness)</b> has <b>disvalue</b> . Hence, <b>ethical hedonism</b> implies <b>value monism</b> . However, <b>quality</b> appears to count as another <b>intrinsic good</b> .
C: Therefore, either <b>qualitative hedonism</b> will turn out to be <b>quantitative hedonism</b> in disguise (horn 1) or <b>qualitative hedonism</b> will be <b>inconsistent</b> (horn 2).	This is the <b>dilemma</b> confronting <b>qualitative hedonist</b> .

### 10.9.2 Schmidt-Petri's response to the dilemma

We may sometimes have to make **trade-offs** between a **lower quantity** of a **higher-quality pleasure** and a **higher quantity** of a **lower-quality pleasure**.

1. Identify the **opponent**.
  - The standard view of **qualitative hedonism** states that a **higher-quality pleasure** will **always** be chosen over a **lower-quality pleasure**, even when the **lower-quality pleasure** is available in a **lower quantity**.
2. Identify objections to the **standard view of qualitative hedonism**.
  - (a) **Ambiguity** over the notion of "**quality**". Mill does not tell us what would correspond to the concept of "**quality**".
  - (b) **Lack of clarity** about **experts**. We do not know how to tell who is an **expert** in the real world.
3. Identify the **source material**, like Mill's Utilitarianism (Chapter 2, paragraph 5).
4. Distinguish between the **standard view** and the **non-standard view** relative to the **source material**.
  - The **standard view** (incorrect):  $q \rightarrow p$  If we are justified in saying that  $x$  is of a **higher quality** than  $y$ , then some pleasure  $x$  is chosen over another pleasure  $y$  available in a **higher quantity**.
  - The **non-standard view** (correct):  $p \rightarrow q$  If some pleasure  $x$  is chosen over another pleasure  $y$  available in a **higher quantity**, then we are justified in saying that  $x$  is of a **higher quality** than  $y$ .



1. Provide examples to support the **non-standard view**.

(a) Example 1:

- Wine X (same quantity)  $\succ$  Wine Y (same quantity)
- Wine X (slightly **lower-quantity**)  $\succ$  Wine Y (slightly **higher-quantity**)
- Appealing to the **higher quality** of Wine X is the most natural way for us to explain these **preference order relations**.

i. Example 2:

- Wine X (**lower-quantity**)  $\succ$  Wine Y (some **higher quantity**  $\leq n$  units)
- Wine Z (**lower-quantity**)  $\succ$  Wine Y (some **higher quantity**  $\leq m$  units, where  $m > n$ )

ii. Example 3:

- Wine X (**1 glass**)  $\succ$  Wine Y (**ANY quantity**)

Therefore, the **standard view** is a **special case** of the **non-standard view**.

## 10.10 Haydn and the oyster thought experiment

- Suppose you are a soul waiting to be allocated life on Earth.
- The angel offers you a choice between:
  - $\phi_1$ : Living the life of **Joseph Haydn**. Haydn composed some wonderful music, influenced the evolution of the symphony, was cheerful and popular, travelled, and enjoyed field sports.
  - $\phi_2$ : Living the life of an **oyster**. The **oyster's life** consists only of **mild sensual pleasure**. However, the **oyster's life** can be as long as you like.

### 10.10.1 Bentham's hedonic calculus

The main part of the differences between the options in this thought experiment is the duration in which the **pain or pleasure lasts**.

- Let  $m$  be the duration of the **oyster's life**.
- Let  $n$  be a **threshold for sufficiency**.
- $m$  is a **sufficiently long duration** if and only if  $m \geq n$ .
- If the **oyster's life** is **sufficiently long**  $m \geq n$ , then:  
 $\phi_2(\text{where the oyster's life is } m \text{ years}) \succ \phi_1(\text{where Haydn's life is 77 years})$

However, this cannot be right.

### 10.10.2 Response

- We could ditch **quantitative hedonism** in favour of **qualitative hedonism**.
- **Qualitative hedonism** gives us automatic reasons to favour **higher-quality pleasure** (such as the life of Joseph Haydn) over **lower-quality pleasures** (such as an oyster's life).
- However, according to the **non-standard view of qualitative hedonism**, we may sometimes have to make **trade-offs** between a **lower quantity of a higher-quality pleasure** and a **higher quantity of a lower-quality pleasure**.
- If the **oyster's life** is **sufficiently long**  $m \geq n$ , then a **trade-off** may have to be made.
- Hence,  $\phi_2$  (where the oyster's life, though **lower-quality**, is **sufficiently long**)  $\succ \phi_1$  (where Haydn's life, though **higher-quality**, is **insufficiently long**)

### 10.11 Utility monster thought experiment

- Suppose that a **utility monster** is a hypothetical entity that is a **highly efficient consumer of resources**.
- The **utility monster** gains vast amounts of **pleasure (happiness)** from very small quantities of a particular resource.
- You have to choose between:
  - $\phi_1$ : Satisfying the needs of the ordinary human beings.
  - $\phi_2$ : Satisfying the needs of the **utility monster**.
- Given the existence of this **utility monster** and the **core consequentialist commitment to maximising the good**, it seems that we ought to neglect  $\phi_1$  in favour of  $\phi_2$ .
- However, this cannot be right.

## 10.12 Against hedonism

- **Consequentialists** will agree that:
  - $\phi_1$  (living the life of **Joseph Haydn**)  $\succ$   $\phi_2$  (living the life of an **oyster**), however long the oyster's life might be.
  - $\phi_1$  (satisfying the needs of ordinary human beings)  $\succ$   $\phi_2$  (satisfying the needs of the **utility monster**), however **efficient** a **consumer of resources** the **utility monster** might be.
- **Consequentialists** will typically not give up the **core consequentialist commitment to maximising the good**.
- However, **consequentialists** may give up the **hedonistic** theory of the good.
- The good is defined as that which **facilitates human pleasure and happiness**.
- **Ethical hedonism** is a key element in the **hedonistic** theory of the good, as only **pleasure (happiness)** has **value** and only **pain (unhappiness)** has **disvalue**.

### 10.12.1 Objections to ethical hedonism

Objection	Elaboration
Deny <b>value monism</b>	<p><b>Ethical hedonism</b> implies <b>value monism</b>, where <math>n</math> denotes the number of <b>intrinsic goods</b>.</p> <p><b>Value monism</b> is when <math>n = 1</math>, while <b>value pluralism</b> is when <math>n &gt; 1</math>.</p> <p>However, we could maintain that while we ought to <b>maximise the good</b>, the <b>good</b> include far more than can be reduced to <b>pleasure (happiness)</b>. Hence, <b>pleasure (happiness)</b> is not the only <b>intrinsic good</b>. Other <b>intrinsic goods</b> include <b>beauty, friendship</b>, and so on. This leads to <b>value pluralism</b>.</p>
Deny that certain <b>pleasure states</b> are <b>intrinsically good</b> . Problems with both <b>quantitative hedonism</b> and <b>qualitative hedonism</b> .	<p>A <b>sadist whipping her victim</b> or an <b>addict on drugs</b> might experience <b>pleasure (happiness)</b>. However, these <b>pleasure states</b> are not <b>intrinsically valuable</b>.</p> <p><b>Quantitative hedonism</b> might lead us to overvalue <b>bestial, unsophisticated, lower-quality, and debauched pleasures</b> (the push-pin versus poetry example).</p> <p>At the same time, <b>qualitative hedonism</b> faces a <b>dilemma</b>. Either <b>qualitative hedonism</b> will turn out to be <b>quantitative hedonism</b> in disguise, or <b>qualitative hedonism</b> will be <b>inconsistent</b>.</p>
Alternative <b>theories of the good</b> are superior to the <b>hedonistic theory of the good</b> .	<p><b>Pleasure and pain sensations</b>, inside the <b>heads of human beings</b>, are <b>difficult to measure</b>. By contrast, <b>well-being</b> can be defined in terms of <b>preference fulfilment</b> or <b>desire satisfaction</b>, giving rise to alternative <b>theories of the good</b>.</p>

### 10.12.2 Alternative theories of the good

- Present desire satisfaction theory: The **good** is defined as that which **facilitates the satisfaction of our current desires**. We also need a **fitting attitude account**, according to which what is **desired** is closely linked with what is **good**.
- Comprehensive desire satisfaction theory: The **good** is defined as that which **facilitates the satisfaction of desires over the course of our life**.
- Informed desire satisfaction theory: The **good** is defined as that which **facilitates the satisfaction of the desires we would have if we were fully informed of all the relevant facts**.
- Objective list theory: The **good** is defined as that which does not consist merely in either **pleasurable experience (hedonistic theory of the good)** or **desire satisfaction (desire satisfaction theory)**.

### 10.12.3 Desire satisfaction theory

- Desire satisfaction theory is also known as **preference fulfilment theory**.
- We could get agents to **rank-order their preferences** and develop **utility functions** on their behalf.
- **Preference utilitarianism** is based on **desire satisfaction theory**.
- Relative to the agents' reported and **rank-ordered preferences** and these **utility functions**, we may be able to derive the following:
  - $U(\phi_1) > U(\phi_2)$ , however long the oyster's life might be. Hence,  $\phi_1$  (living the life of **Joseph Haydn**)  $\succ \phi_2$  (living the life of an **oyster**).
  - $U(\phi_1) > U(\phi_2)$ , however **efficient** a **consumer of resources** the **utility monster** might be. Hence,  $\phi_1$  (satisfying the needs of ordinary human beings)  $\succ \phi_2$  (satisfying the needs of the **utility monster**).

#### 10.12.4 Minimisation of the violation of rights

- The minimisation of the violation of rights could be identified as another relevant end for **consequentialism**.
- Human beings have the **right to resources** that allow their needs to be satisfied.
- Their **right to resources** function as **side constraints** on the pursuit of **good consequences**.
- This gives rise to a **utilitarianism of rights**.
- Therefore,  $\phi_1$  (satisfying the needs of ordinary human beings)  $\succ$   $\phi_2$  (satisfying the needs of the **utility monster**).

#### 10.13 Formal definition of consequentialism

- The term "consequentialism" was first introduced by G. E. M. Anscombe in "Modern moral philosophy" (1958).
- According to consequentialism, only the **consequences of actions** matter, whereas the **intentions behind actions** are unimportant.

The elements of consequentialism include:

- A **core consequentialist commitment**, which is the **maximisation of the good**.
- A **collection of theories of the good**.
  - **Egoistic**: The good is defined as that which **facilitates self-interest**.
  - **Hedonistic**: The good is defined as that which **facilitates human pleasure and happiness**.
  - **Desire satisfaction theory**: The good is defined as that which **facilitates the satisfaction of our desire**.
  - **Objective list theory**: The good is defined as that which does not consist merely in either **pleasurable experiences (hedonistic theory of the good)** or **desire satisfaction (desire satisfaction theory)**.

### 10.13.1 Semantics of consequentialism

- The semantics of consequentialism can be interpreted in terms of a **semantics of possible worlds**.
- To **maximise the good** is to make the **world**, the **sum of all things**, as good as it can be.
- Alternatively, to **maximise the good** is to act to bring about the **best possible world** or all those worlds that can be brought about.
- An **agent** ( $j, k$ , and so on) is a **being capable of actions** that are **apt for moral evaluation**.
- Each **action**  $\phi_i$  brings about a **possible world**  $w_i$ .
- Hence, an associated set of **possible worlds**  $w_1, w_2, \dots, w_n$  is brought about by the alternative courses of action  $\phi_1, \phi_2, \dots, \phi_n$ .
- **Possible worlds**  $w_1, w_2, \dots, w_n$  are alternatives between which an **agent**  $j$  must choose.

### 10.13.2 Definitions

- Let  $\Phi$  denote the set of alternative courses of action  $\phi_1, \phi_2, \dots, \phi_n$ .

$$\Phi = \{\phi_1, \phi_2, \dots, \phi_n\}$$

- Let  $W$  denote the set of **possible worlds**  $w_1, w_2, \dots, w_n$ , brought about by the actions in  $\Phi$ .

$$W = \{w_1, w_2, \dots, w_n\}$$

- Let  $A$  denote the association of members of  $\Phi$  with their associated **possible worlds** in  $W$ .

$$A = \{\phi_1 - w_1, \phi_2 - w_2, \dots, \phi_n - w_n\}$$

- Let  $\langle j, A \rangle$  denote a **choice situation**, where  $j$  is the agent,  $A$  is the set of associations.
- Let  $R$  denote a **rightness function** that assigns to each **choice situation**  $\langle j, A \rangle$ , a **subset of  $A$** .

$$\langle j, A \rangle \xrightarrow{R} R_j(A)$$

$$\{\phi_1 - w_1, \phi_2 - w_2, \dots, \phi_n - w_n\} \xrightarrow{R} \{\phi_1 - w_1(\text{good}), \phi_2 - w_2(\text{good}), \dots, \phi_n - w_n(\text{good})\}$$



### 10.13.3 Conditions

Let  $T$  denote a **theory of good**, representable in terms of a **complete order**  $\leq$  on the set  $W$  of all possible worlds.

Condition	Description	Formal representation
<b>Reflexivity</b>	$x$ is <b>at least as good as itself</b> , according to $T$ .	$x \leq x$
<b>Transitivity</b>	If $y$ is <b>at least as good as</b> $x$ and $z$ is <b>at least as good as</b> $y$ , then $z$ is <b>at least as good as</b> $x$ , according to $T$ .	$((x \leq y) \wedge (y \leq z)) \rightarrow (x \leq z)$
<b>Completeness</b>	It is either the case that $y$ is <b>at least as good as</b> $x$ or that $x$ is <b>at least as good as</b> $y$ , according to $T$ .	$(x \leq y) \vee (y \leq x)$

#### 10.13.4 Axioms of consequentialism

1. AN or agent neutrality. Consequentialism is **agent-neutral**. For any two agents  $j$  and  $k$ :

$$R_j(A) = R_k(A)$$

Essentially, it doesn't matter who or what the agent is, the right thing to do will always be the same.

2. NMD or **no moral dilemmas**. There are **no moral dilemmas** under **consequentialism**.

$$R_j(A) \neq \emptyset$$

3. DM or **dominance**.  $x$  **dominates**  $y$  under **consequentialism**.

- For any two possible worlds  $x$  and  $y$ , an agent  $j$ , and any two possible choice situations  $\langle j, A \rangle$  and  $\langle j, B \rangle$ :
- Suppose that  $\{x, y\} \subseteq A \cap B, x \in R_j(A)$  and  $y \notin R_j(A)$ .

$$\therefore y \in R_j(B)$$

- Essentially, what this axiom states is that once you have determined that one action is better than another action in a decision, the other action that you have rejected cannot be the best action in the next decision you make, because there is already a better action available.

### 10.13.5 Elements of consequentialism

1. The **agent**  $j$ .
2. A set  $\Phi$  of alternative courses of action  $\phi_1, \phi_2, \dots, \phi_n$ .
3. A set  $W$  of **possible worlds**  $w_1, w_2, \dots, w_n$  brought about by the actions in  $\Phi$ .
4. An association  $A$  between members of  $\Phi$  and their associated **possible worlds** in  $W$ .
5. The **choice situation**  $\langle j, A \rangle$ .
6. The **rightness function**  $R$ .
7. The **theory of the good**  $T$ , defined in terms of the following conditions:

Reflexivity:  $x \leq x$

Transitivity:  $((x \leq y) \wedge (y \leq z)) \rightarrow (x \leq z)$

Completeness:  $(x \leq y) \vee (y \leq x)$

8. The **axioms of consequentialism**:

$$(AN) \ R_j(A) = R_k(A)$$

$$(NMD) \ R_j(A) \neq \emptyset$$

$(DM)$  Suppose that  $\{x, y\} \subseteq A \cap B, x \in R_j(A)$  and  $y \notin R_j(A)$ :

$$\therefore y \notin R_j(B)$$

## 10.14 Violation of consequentialist axioms

It is possible for at least some of the **axioms of consequentialism** to be violated.

### 10.14.1 Violating the agent neutrality axiom

- **Agent-relative consequentialism** combines the **core consequentialist commitment** (each agent ought to **maximise the good**) with an **agent-relative axiology**.
- The correct evaluation of the **complete order**  $\leq$  on the **set  $W$  of all possible worlds** may vary from agent to agent.
- If there is an **agent-relative axiology**, then  $\Diamond(R_j(A) \neq R_k(A))$ .
- Hence, **AN (agent neutrality)** may be violated.

### 10.14.2 Violating the no moral dilemmas axiom

- If the **theory of the good  $T$**  is defined in terms of **completeness**, then for any two possible worlds  $x$  and  $y$ , either  $x > y$ ,  $x = y$  or  $y > x$ .
- The **good** must be represented by a **complete order**  $\leq$ .
- This rules out **incommensurability**.
- However, if  $\leq$  is **incomplete**, then we could allow for **incommensurability between at least some possible worlds**.
- Therefore, this would allow **moral dilemmas** within **consequentialism**.

$$R_j(A) = \emptyset$$

- Therefore, **NMD (no moral dilemmas)** will be violated.
- The possibility of **moral dilemmas** within **consequentialism** might require us to drop NMD.

### 10.14.3 Violating the dominance axiom

- **Satisficing consequentialism** relaxes the **maximising element** in consequentialism.
- According to **satisficing theories**, it is sometimes permissible to do something than is **worse than the best**, provided that it is **good enough**.
- Hence, DM (**dominance**) will be violated.
- Relative to the choice situations  $\langle j, B \rangle$ , the **satisficing threshold** is **sufficiently low**.
- Therefore,  $x \in R_j(B)$  and  $y \in R_j(B)$  (both  $x$  and  $y$  are **good enough**).
- Hence,  $x$  does not **dominate**  $y$ .
- DM (**dominance**) will be violated.

## 10.15 Driver's strategy

### 10.15.1 Move 1

Identify objections to **consequentialism** raised by the **standard feminist view**.

- Objection 1: **Consequentialism** is **too demanding** of the individual.
- Objection 2: **Consequentialism** is **neglectful of an agent's special obligations** to family and friends.

### 10.15.2 Move 2

Identify the problems of **consequentialism**.

- Problem 1: **Impartiality** - we should **maximise the good, impartially considered**.
- Problem 2: **Demandingness** - **consequentialism** appears to make **theoretically unlimited demands**.

### 10.15.3 Move 3

Introduce **conceptual distinctions**.

- The angel of the house is **partial to her domestic sphere** and **self-sacrificing**.
- The angel of the world is **impartial** and **self-sacrificing**.
- **Moral self-sufficiency**: There is a concern with promoting an **internalist standard of moral worth**, as our wills are sufficient to ground our moral worth, and we do not need to look externally for the source of this value.
- **Ethical self-sufficiency**: There is a concern with promoting a **risk aversion standard**.
  - The best life needs to incorporate a **respect for personal space**.
  - There is space in **ethics** for **self-perfection, personal projects, goals, and aspirations**.

#### 10.15.4 Move 4

Articulate a **sophisticated consequentialism**.

- The aim here is to avoid the **angel of the world** and **ethical self-sufficiency** (Schopenhauer's ideal).
- Introduce further **conceptual distinctions** between:
  1. A **decision procedure** and a **criterion for rightness of actions**. A **decision procedure** may interfere with the actual production of the good, as determined by some **criterion of rightness**.
  2. **Truth conditions** and **acceptance conditions**. We could hold that **consequentialism** is true, although we do not accept it in certain areas of our life.

#### 10.15.5 Move 5

Anticipate objections to **sophisticated consequentialism**.

- Objection: **Moral schizophrenia**.
- Response: **Moral schizophrenia** does not speak against the **truth** of a theory.
- It merely shows that the theory is **difficult to apply** as a **decision procedure**.
- A **theory** may be **true** yet difficult to apply.
- The **impartial good** will be promoted if the **angel of the world** cultivates **special relationships** and becomes less of an all-encompassing angel.
- Therefore, we now have a **norm of partiality**, which means we ought to show preference for our family and friends.

#### 10.15.6 Conclusion

As a **sophisticated consequentialist**, one can be **both a consequentialist and a feminist**. One can be responsive to both **feminist concerns about partiality** and the **demands of morality**.

## 11 Deontology

### 11.1 Hiroshima and Nagasaki

- Should the atomic bombings of Hiroshima and Nagasaki have been carried out or not?
- Let  $\phi$  denote the action of carrying out the atomic bombings of Hiroshima and Nagasaki.
- Since  $\phi$ -ing has brought about the **best possible consequences** (the end of WWII and world peace),  $\phi$ -ing should have been carried out.

#### 11.1.1 Argument in favour of $\phi$ -ing

- Only the **consequences of actions** matter, whereas the **intentions behind actions** are unimportant.
- There are no **side constraints** under **consequentialism**.
- $\phi$ -ing would be **morally right** due to the **maximisation of the good**.
- Therefore, the **evaluative claim** under **consequentialism** is that we ought to  $\phi$ .

#### 11.1.2 Argument not in favour of $\phi$ -ing

- **Side constraints** exist under **deontology**  $\phi$ -ing would be **morally wrong** due to the **violation** of certain **side constraints**.
- Therefore, the **deontic verdict** under **deontology** is that we ought not to  $\phi$ .



## 11.2 Textbook view

According to the **textbook view**:

- **Consequentialism** is an **agent-neutral theory**, whereas **deontology** is an **agent-relative theory**.
- The **textbook view** has also been describe as the **standard method** for drawing the distinction between **consequentialism** and **deontology**.

The **textbook view** implies the following:

Deontology	Consequentialism
Deontology may give <b>different agents</b> $j, k$ , etc. <b>different aims</b> .	<b>Consequentialism</b> gives <b>different agents</b> $j, k$ , etc. the <b>same aim</b> : select the action $\phi_i$ that <b>maximises the good</b> .
Deontology introduces <b>agent-relative side constraints</b> .	<b>Consequentialism</b> does not observe any <b>agent-relative side constraints</b> .
$\Diamond(R_j(A) \neq R_k(A))$ with an <b>agent-relative axiology</b> . Hence, AN ( <b>agent neutrality</b> ) may be violated.	(AN or <b>agent neutrality</b> ) $R_j(A) = R_k(A)$  AN is an <b>axiom</b> under <b>consequentialism</b> .
We can give <b>higher weight</b> to our interests, projects, and concerns.	We must give <b>equal weight</b> to our interests and the interests of others.
We might have certain <b>special obligations</b> (for instance, <b>obligations</b> that parents have to their own children), not shared by other agents.	We ought to be <b>impartial</b> .
We end up with <b>agent-relative reasons for action</b> .	We end with <b>agent-neutral reasons for action</b> .

### 11.2.1 Exceptions to the textbook view

#### 1. Ethical egoism

- **Ethical egoism** is supported by an **egoistic** theory of the good.
- The **good** is defined that which **facilitates self-interest**.
- Since what counts as **self-interest** may differ from agent to agent, **ethical egoism** may give **different agents  $j, k$ , etc. different aims**.
- Therefore, **ethical egoism** is an **agent-relative** version of **consequentialism**.

#### 2. Agent-neutral deontology

- **Agent-neutral deontology** is an **agent-neutral** version of **deontology**.
- **Agent-neutral deontology** gives **different agents  $j, k$ , etc. the same aim**.

### 11.3 Mafia scenario

- Suppose that the mafia are credibly threatening to kill two strangers unless you **kill a third stranger**.
- Let  $\phi$  denote the action of **killing the third stranger**.
- Should we ought to  $\phi$  or  $\neg\phi$ ?

Deontology	Consequentialism
The <b>deontic verdict</b> is that we ought not to $\neg\phi$ .	The <b>evaluative claim</b> is that we ought to $\phi$ .
<b>Justification</b> for <b>deontic verdict</b> . <b>Side constraints</b> exists that <b>prohibit killing</b> as a <b>morally wrong</b> type of action.	<b>Justification</b> for <b>evaluative claim</b> . We ought to <b>minimise the number of bad occurrences</b> .

## 11.4 Side constraints

There are two types of **side constraints**.

Type of side constraint	Description	Implication
An <b>agent-relative side constraint</b> .	Each agent should ensure that <b>she does not kill</b> , even if to <b>prevent more killings by others</b> .	<b>Agent-relative side constraints</b> might give an agent a <b>special concern</b> with her own killings.
An <b>agent-neutral side constraint</b> .	Each agent should ensure that <b>no one kills</b> , even if to <b>prevent more killings by others</b> . <b>Agent-neutral side constraints</b> may require everyone to <b>share a moral vision</b> .	<b>Agent-neutral side constraints</b> gives <b>different agents</b> $j, k$ , etc. the <b>same aim</b> .

## 11.5 The right and the good

- The **textbook view** uses the **agent-neutral** and **agent-relative** distinction to distinguish between **consequentialism** and **deontology**.
- Another popular approach to distinguishing between **consequentialism** and **deontology** involves an appeal to the **right or good** distinction:

Deontology	Consequentialism
The <b>right</b> is <b>prior</b> to the <b>good</b> .	The <b>good</b> is <b>prior</b> to the <b>right</b> .
<b>Deontology</b> delivers <b>deontic verdicts</b> in terms of <b>what is right</b> .  If an action $\phi$ is <b>morally impermissible</b> , then it is <b>not right</b> , no matter how much <b>good</b> might be produced by $\phi$ -ing.	<b>Consequentialism</b> delivers <b>evaluative claims</b> in terms of <b>what is good</b> .

## 11.6 Action and intention

We can distinguish between different types of **deontology** in terms of the **action or intention** distinction and the **agent versus patient** distinction:

Types of deontology	Description
Action-based deontology	The <b>right action</b> is <b>prior</b> to the <b>good consequences</b> .
Intention-based deontology	The <b>right intention</b> is <b>prior</b> to the <b>good consequences</b> .
Action and intention-based deontology	The <b>right action</b> and <b>right intention</b> are <b>prior</b> to the <b>good consequences</b> .
Agent-centred deontology	The primary concern is with the <b>duties of the agent</b> .
Patient-centred deontology	The primary concern is with the <b>rights of the patient</b> . The <b>patient</b> has a <b>right against being used</b> as a mere <b>means for producing good consequences</b> without her consent.

## 11.7 Doctrine of double effect

- **Intention-based deontology** is supported by the **doctrine of double effect**.
- The doctrine of double effect states that it is **morally impermissible** for us to **intend evil** (for instance, the **killing or torturing of innocents**).

Possible objections to the **doctrine of double effect**:

1. The distinctions between **intending, foreseeing, risking, predicting** and **causing and allowing** are **conceptually incoherent**.
2. The distinctions invite **manipulation** and can be **exploited**.

### 11.8 Avoision

- Suppose that there is a **moral prohibition** against  $\phi_1$ -ing.
- However,  $\phi_2$ -ing allows an agent to **bypass, circumvent, or duck this moral prohibition**.
- Does  $\phi_2$ -ing count as a **morally permissible avoidance** of the **moral prohibition**?
- Does  $\phi_2$ -ing count as a **morally impermissible evasion** of the **moral prohibitions**?
- This type of **manipulation** in the legal domain has been termed **avoision**.
- The distinctions between **intending, foreseeing, risking, predicting** and **causing and allowing** could give rise to a **moral** version of **avoision**.

### 11.9 Trolley dilemma

- According to the **trolley dilemma**, a **runaway tram or trolley** is on course to run over and **kill five people on the main track**.
- However, you can intervene, **pull the lever**, and divert the runaway tram or trolley to a **side track, killing just one person**.
- You have to choose between:
  - $\phi$ : **Pulling the lever and killing the one person on the side track to save the five people on the main track**.
  - $\neg\phi$ : **Refraining from pulling the lever and letting the five people on the main track die**.
- Should you **pull the lever and kill one to save five  $\phi$**  or **refrain from doing so  $\neg\phi$** ?

### 11.9.1 Positions

- **Consequentialism** recommends **pulling the lever and killing one to save five** ( $\phi$ ).
  - $\phi$ -ing will **minimise the number of bad occurrences**.
  - We ought to **minimise the number of bad occurrences**.
  - Hence, we ought to  $\phi$ .
- **Action-based deontology** recommends **refraining from pulling the lever** ( $\neg\phi$ ).
  - $\phi$ -ing is a **morally impermissible action**.
  - $\neg\phi$ -ing is a **morally permissible action**.
  - Hence, we ought to  $\neg\phi$ .
- **Intention-based deontology** recommends two different actions.
  1. **Pull the lever and kill one to save five** ( $\phi$ ).
    - $\phi$ -ing may be accompanied by the **intention of saving five**.
    - We merely **risk, foresee, or predict** that  $\phi$ -ing will have the consequence of one innocent person being killed.
    - Hence, we ought to  $\phi$ .
  2. **Refrain from pulling the lever** ( $\neg\phi$ ).
    - $\phi$ -ing may be accompanied by the **intention of killing one innocent person**.
    - We should not **intend evil**.
    - Hence, we ought to  $\neg\phi$ .
- **Patient-centred deontology** recommends **refraining from pulling the lever** ( $\neg\phi$ )
  - **That one person on the side track** has a **right against being used** as a mere **means for producing good consequences** without their consent.
  - Hence, we ought to  $\neg\phi$ .

### 11.10 Siamese twins

- Suppose that **Siamese twins** are conjoined such that **both will die** unless the **organs of one are given to the other** via an operation that kills the first.
- You have to choose between:
  - $(\phi)$ : Performing the **life-saving operation**.
  - $(\neg\phi)$ : **Refraining from** performing the **life-saving operation**.

#### 11.10.1 Positions

- **Intention-based deontology** recommends **performing the life-saving operation**  $(\phi)$ 
  - We **intend** to save the first twin in a **life-saving operation**.
  - We do not **intend** but merely **foresee** the death of the second twin in that operation.
  - Hence, we ought to  $\phi$ , justified by the use of **foreseeing and intending** distinction.
- **Action-based deontology** recommends **performing the life-saving operation**  $(\phi)$ 
  - We **cause the first twin to be saved** in a **life-saving operation**.
  - We do not **cause** but merely **allow** the death of the second twin in that operation.
  - Hence, we ought to  $\phi$ , justified by the use of the **causing and allowing** distinction.
- **Patient-centred deontology** recommends **refraining from performing the operation**  $(\neg\phi)$ 
  - Each twin has a **right against being used** as a mere **means for producing good consequences** without her consent.
  - Hence, we ought to  $\neg\phi$ .

### 11.11 Divine command theory

- **Divine command theory** is a form of **deontology**.
- In **divine command theory**, the **side constraints** are provided by the **divine commands of God**.
- Examples of **divine commands** and religious precepts supporting the **side constraint against murder**:
  1. Thou shalt not kill, from Exodus 20:13. This is one of the **10 Commandments in Judaism and Christianity**.
  2. I undertake the precept to **refrain from destroying living creatures**. This is one of the **five precepts in Buddhism**.
  3. One ought to **avoid harming (or desiring to harm) any living being** in thought, word or deed. This is the **doctrine of Ahimsa in Jainism, Hinduism, and Buddhism**.
- These religious commandments, precepts, and doctrines **prohibit killing** as a **morally wrong** type of action and function as **side constraints**.

#### 11.11.1 Definition

- **Divine command theory** is the view according to which:
  - **Morality** is somehow dependent upon **God**.
  - **Moral obligation** consists in **obedience to God's commands**.
- **Divine command theory** consists of at least some of the following claims:
  1. **Morality** is the ultimately based on the **commands of God**.
  2. **Morality** is ultimately based on the **character of God**.
  3. **Moral obligations** are identical with **divine commands**.
  4. **Moral obligations** are created by **divine commands**.
- **Divine command theory** has to address the **Euthyphro dilemma**.



### 11.11.2 Euthyphro dilemma

- Socrates' original question to Euthyphro: Is the **pious** loved by the gods because it is pious, or is it pious because it is loved by the gods?
- Revised version of Socrates' question: Is an action **divinely commanded** by God because it is **morally right**, or is an action **morally right** because it is **divinely commanded** by God?
- Horn 1: An **action** is **divinely commanded** by God because it is **morally right**. Problems associated with this horn:
  1. God would no longer be the **author of morality**. Rather, God would be merely a **being capable of recognising right and wrong**.
  2. God would not be the **sovereign of morality**. Rather, God would be a mere **subject of morality**.
- Horn 2: An **action** is **morally right** because it is **divinely commanded** by God. Problems associated with this horn:
  1. Nothing guarantees that what God would **divinely command** would be in accordance with what **morality** would prescribe. If God commands us to do something **morally reprehensible**, then the **morally reprehensible** action would be **morally right**.
  2. The **foundations of morality** become **arbitrary**. **Morally reprehensible actions** can become **morally obligatory** if **divinely commanded** by God.

### 11.12 Modified divine command theory

- According to **standard divine command theory**: " $\phi$ -ing is **morally wrong**" means " $\phi$ -ing is contrary to the **divine commands** of God."
- According to the **modified divine command theory**: " $\phi$ -ing is **morally wrong**" means " $\phi$ -ing is contrary to the **divine commands** of a **loving** God."
- **Modified divine command theory** holds that God has a particular character.
- God has the character of **loving his human creatures**.
- Therefore, while it is **logically possible** for God to command **morally reprehensible** action, it is not actually something that God would do, given his character.

#### 11.12.1 How it fixes the issues with Euthyphro's dilemma

Horn 1:

- God is the **subject** rather than the **sovereign of morality**.
- Under **modified divine command theory**, God remains the **source of morality**.
- **Morality** is grounded in the **loving** character of God.

Horn 2:

- The **foundations of morality** become **arbitrary**.
- Under **modified divine command theory**, not any **divine command** goes.
- **Divine commands** remain rooted in the unchanging **loving** and **omnibenevolent** character of God.

**Standard divine command theory** is pierced by either **horn 1** and **horn 2**. However, the **modified divine command theory** appears to fare better against the **Euthyphro dilemma**.

## 11.13 Plato

### 11.13.1 Context

- **Socrates** is being accused of **corrupting the youth of Athens**.
- **Socrates** is about to be put to trial and thereafter found **guilty** of both:
  1. **Corrupting the minds of the youth of Athens**.
  2. **Impiety** ("not believing in the Gods of the state").
- **Euthyphro** is **prosecuting his father** for **murder**.
- Both **Socrates** and **Euthyphro** try to arrive at a definition of **piety** or **holiness**.
- **Euthyphro's** modified definition of **piety** or **holiness** (suggested by **Socrates**):
  - Holy is defined as what all the gods love.
  - Unholy is defined as what all the gods hate.
  - $\neg$  (Holy  $\vee$  Unholy) is defined as what some gods love and some other gods hate.
  - (Holy  $\vee$  Unholy) is defined as what some gods love and some other gods hate.

### 11.13.2 The Euthyphro's dilemma

- The **Euthyphro dilemma**: Is the **holy** loved by the gods because it is holy? Or is it holy because it is loved by the gods?
- Revised version of the **Euthyphro dilemma**: Is what is **morally good** commanded by God or the gods because it is **morally good**? Or is it **morally good** because it is commanded by God or the gods?
- **Species-genus** distinction:
  - The **species** is a **part** of the **genus (whole)**.
  - **Reverence (species)** is a **part** of **fear (genus)**.
  - **Piety or holiness (species)** is a **part** of **justice (genus)**.
- According to **Euthyphro**: The part of **justice (genus)** that is **holy (species)** involves **ministering to** (looking after or **taking care of**) the gods.

### 11.13.3 Socrates' argument by analogy

Source	Target
We <b>take care of animals</b> .	The part of <b>justice (genus)</b> that is <b>holy (species)</b> involves <b>taking care of the gods</b> (Euthyphro's claims).
<b>Taking care of horses, dogs, and cattle</b> in horsetraining, huntsmanship, and herdsman-ship <b>benefits and improves horses, dogs, and cattle</b> . You make these animals <b>better</b> .	You do not make the gods <b>better</b> when you do something <b>pious or holy</b> .

Therefore, Euthyphro's claims are problematic.

### 11.13.4 Conclusion

- The **pious or holy** is in fact **something other than what is acceptable to the gods**.
- Something's being acceptable to the gods is merely an **attribute of piety** and not part of its **defining characteristics**.
- Therefore, Socrates still does not have a **definition of piety or holiness**.

## 11.14 Secular deontology

Deontology may be **religious** or **secular** in nature.

### 11.14.1 Religious form of deontology

- All **duties and obligations** are generated by **religious precepts, commandments, and doctrines**.
- An example is **divine command theory**, which has **religious precepts, commandments, and doctrines** that are of **divine origin**.

### 11.14.2 Secular form of deontology

- All **duties and obligations** are generated by **secular principles**.
- An example is **deontological monism**, where all **duties and obligations** are generated by a **single secular principle**.
- Another example is **deontological pluralism**, where all **duties and obligations** are generated by **multiple secular principles**.

## 11.15 Kantian deontology

- **Kantian deontology** is an example of a **secular** form of **deontology**.
- All **duties and obligations** are generated by a **single secular principle**.
- This **single secular principle** is known as the **categorical imperative**.

### 11.15.1 Types of imperatives

#### 1. Hypothetical imperative

- A **command** that applies to us **conditionally**.
- **Hypothetical imperatives** contain conditional clauses that can be explicit or elided.
- **Hypothetical imperatives** require us to exercise our wills in a certain way, given that we have **antecedently** willed a particular end.
- An example is, "If you are happy, and you know it, clap your hands".

#### 2. Categorical imperative

- A **command** that applies to us **unconditionally**, regardless of what we might want.
- **Categorical imperatives** do not have any **conditional clauses**, that can be **explicit** or **elided**.
- An example is, "Thou shalt not kill."

### 11.15.2 Categorical imperative formulations

1. The **universal law of nature** formula.
2. The **humanity** formula.
3. The **autonomy** formula.
4. The **kingdom of ends** formula.

Formulations 1 to 4 are supposed to be **equivalent**, hence **Kantian deontology** is also an example of **deontology monism**.

### 11.15.3 Universal law of nature formula

- The **universal law of nature** states:  
Act only in accordance with that **maxim** that you can at the same time will that it become **universal law**.
- Decision procedure associated with the **universal law of nature** formula:
  1. Consider a particular action  $\phi$  (for instance, stealing).
  2. Determine a **maxim** governing  $\phi$  (for instance, I will steal for pleasure).
  3. **Universalise** that **maxim** and consider the implications (for instance, everyone ought to steal for pleasure).
  4. Consider whether the **universalised** version of the **maxim** is **contradiction-free**. Could the **universalised** version of the **maxim** function as a **law of nature**?
- Example: Suppose I borrow money from you and promise to return the amount eventually, but I know that I will not return the amount.

#### 11.15.4 Detailed steps for the universal law of nature formula

1. Consider a particular action  $\phi$ .
  - Here,  $\phi$ -ing denotes the action of **borrowing money and promising to repay**, despite knowing that **this will never be done**.
2. Determine a **maxim** governing  $\phi$ .
  - Maxim: Whenever I am low on cash, I will **borrow money and promise to repay**, despite knowing that **this will never be done**.
3. **Universalise** that **maxim** and consider the implications.
  - **Universalised** version of maxim: Whenever anyone is low on cash, she ought to **borrow money and promise to repay**, despite knowing that **this will never be done**.
  - However, if this deceit is **universalised**, there will be no institution of **promising** to begin with.
4. Consider whether the **universalised** version of the **maxim** is **contradiction-free**.
  - A **contradiction** arises, as we cannot **make promises** if there is no institution of **promising** to begin with.
  - Therefore, I will be unable to make my **promise** to begin with.

#### 11.15.5 Humanity formula

- The **humanity** formula states:  
Act in such a way that you always treat **humanity**, whether in your own person or in the person of any other, never simply as a **means**, but always at the same time as an **end**.
- Implications:
  - We ought to **respect the humanity in persons**.
  - We ought to refrain from treating persons as **mere instruments**.



#### 11.15.6 Autonomy formula

- The **autonomy** formula: Act in such a way that your **will** can regard itself at the same time as **making universal law through its maxims**.
- Implications:
  - We ought to act as **universal lawgivers** or **legislators**.
  - We ought to consider whether our intended maxims are worthy of our status as **shapers of the world**.
  - As **rational agents**, we are the very **source of the authority** for the **moral laws** that bind us.

#### 11.15.7 Kingdom of ends formula

- The **kingdom of ends** formula: Act in accordance with the **maxims** of a member giving **universal laws** for a **merely possible kingdom of ends**.
- Implications:
  - We ought to consider whether our intended maxims will earn **acceptance** by a **community of fully rational agents** in a **kingdom of ends**.
  - Just as **human beings** are **ends in themselves**, we are also a **kingdom of ends** or a **moral community**.

### 11.16 Organ transplant thought experiment

- Suppose that **five mortally ill patients** are at a hospital.
- They will soon die without an organ transplant.
- Each of these patients requires a different organ to be transplanted (for instance, a **heart** for the first patient, a **kidney** for the second patient, a **liver** for the third patient, and so on).
- At the same time, a **sixth healthy patient** is undergoing a routine check-up at the same hospital.
- The only medical means of **saving the five mortally ill patients** would be to **kill the sixth healthy patient** and transplant his healthy organs into the five other patients.
- You have to choose between:
  - $(\phi)$ : Killing the **one healthy patient** to **save the five mortally ill patients**.
  - $(\neg\phi)$ : Refraining from killing the **one healthy patient**, even at the cost of the **five mortally ill patients** dying.
- Should you kill the **one healthy patient** to **save the five mortally ill patients**  $(\phi)$  or refrain from doing so  $(\neg\phi)$ ?

#### 11.16.1 Consequentialist position

- Consequentialism recommends **killing the one healthy patient** to **save the five mortally ill patients**  $(\phi)$ .
- Justification:
  - $\phi$ -ing will **maximise the good**.
  - Certain **theories of the good** might require us to  $\phi$ .
  - Therefore, we ought to  $\phi$ .

### 11.16.2 Deontological position

- The deontological position recommends refraining from killing the one healthy patient  $\neg\phi$ .
- Justification:
  - $\phi$ -ing is a **morally impermissible action**.
  - $\neg\phi$ -ing is a **morally permissible action**.
  - The **humanity** formula introduces **side constraints** against  $\phi$ -ing.
  - $\phi$ -ing requires that you treat the **sixth healthy patient** as a mere **means or conduit** for **five useful and life-saving organs**.
  - Hence, we ought to  $\neg\phi$ .

### 11.16.3 Permissivity

For consequentialism:

- The **evaluative claim** of consequentialism:
  - $\phi$ -ing is not merely **morally permissible**.
  - $\phi$ -ing is also **morally obligatory**, we ought to **maximise the good (core consequentialist commitment)**
  - Hence, **consequentialism** appears to be **overpermissive**.

For deontology:

- The **deontic verdict** of deontology:
  - $\phi$ -ing is **morally impermissible**.
  - $\phi$ -ing violates certain **side constraints** (including the **humanity** formula)
  - Hence, **deontology** is not **overpermissive**.

Note that the **organ transplant** thought experiment and the **trolley dilemma**:

- (Trolley dilemma): **Pull the lever** and **kill one** to **save five** ( $\phi$ ) or **refrain from doing so** ( $\neg\phi$ )?
- (Organ transplant): **Kill the sixth healthy patient** to **save five**  $\phi$  or **refrain from doing so** ( $\neg\phi$ )?

### 11.17 Russian deontology

- **Kantian deontology** is an example of **deontological monism**.
- All **duties and obligations** are generated by a **single secular principle**.
- By contrast, **Russian deontology** is an example of **deontological pluralism**.
- All **duties and obligations** are generated by **multiple secular principles**.
- Furthermore, these **multiple secular principles** cannot be reduced to a **single master or mistress principle**.

#### 11.17.1 Principles

- **Fidelity**: We have a duty to **keep our promises**.
- **Reparation**: We have a duty to **right our previous wrongs**.
- **Gratitude**: We have a duty to **return services** to those from whom we have accepted benefits in the past.
- **Beneficence**: We have a duty to promote a **maximum of aggregate good**.
- **Nonmaleficence**: We have a duty to **refrain from harming others**.

### 11.17.2 Prima facie duties

- **Rossian principles** specify aspects of a situation that count morally in favour of an action  $\phi$ .
- For instance, an action  $\phi$  that allows us to **right our previous wrongs (reparation)** counts morally in favour of  $\phi$ -ing.
- **Rossian principles** yield **prima facie duties**.
- However, **prima facie duties** can be **outweighed** by other **prima facie duties**.
- Furthermore, **moral dilemmas** may arise when **prima facie duties conflict** with each other.
- What the agent ought to choose under **Rossian deontology** is that action  $\phi_i$ : Of all those possible for the agent in the circumstances, that has the **greatest balance of prima facie rightness**, in those respects in which they are **prima facie right**, over their **prima facie wrongness**, in those respects in which they are **prima facie wrong**.

### 11.17.3 Conflicts between prima facie duties

- **Rossian deontology** maintains that any **conflict between prima facie duties** in a particular situation can be resolved by relying on **intuition**.
- We can rely on our **crystal-clear intuitions** in both **mathematics** and **ethics** to build up all we can know about the **nature of numbers** and the **nature of duty**.
- Our **mathematical knowledge** and our **moral knowledge** are **self-evident**.
- When we assume that  $AB$  and  $CD$  are parallel, i.e.  $AB \parallel CD$ , then  $AB$  and  $CD$  do not meet except at infinity in a Euclidean plane is self-evident by **mathematical intuition**.
- $\neg Pp$ , where  $p$  denotes the action of **genocide** is **self-evident** by **moral intuition**.
- The **actual duty** is what the agent is left with after she has weighed up all the **conflicting prima facie duties** in the manner prescribed by Ross.

#### 11.17.4 Issues with intuition to resolve conflicts

1. How do we even identify the **prima facie duties** that are involved in a particular situation?
2. What are the **criteria** according to which we **rank and compare prima facie duties**, in order to arrive at the **greatest balance of prima facie rightness over prima facie wrongness** (as prescribed by Ross) that will guide us to our **actual duty**?
3. What if our **intuitions** are **wrong or misguided**?

#### 11.18 W.D.

- W.D. is a machine-based implementation of **Rossian deontology**.
- Jeremy and W.D. are important developments in the domain of **machine ethics**.
- **Machine ethics** is broadly concerned with ensuring that the **behaviour of machines** is **ethically acceptable**.
- The **relations** to be learnt by W.D. are represented as **first-order horn clauses** of the form:

$$H \leftarrow (L_1 \wedge L_2 \wedge \cdots \wedge L_n)$$

Where:

- $H$  is a positive literal  $H$
- $\leftarrow$  means implication
- $(L_1 \wedge L_2 \wedge \cdots \wedge L_n)$  is a universally quantified conjunctions of positive literals  $L_i$
- W.D. uses **inductive logic programming** to learn the **supersedes relation**:

$$\text{supersedes}(\phi_1, \phi_2)$$

This means that action  $\phi_1$  is preferred over action  $\phi_2$  in a particular situation.

### 11.18.1 Favours relation

The **favours relation** is a **4-ary relation** that is used as a **specifying operation** to aid the supersedes relation:

$$\text{favours}(1 \text{ or } 2, D_{\phi_1}, D_{\phi_2}, R)$$

Where:

- 1 or 2 signifies in which action's favour ( $\phi_1$  or  $\phi_2$ ), a given **prima facie duty** lies. The possible values are  $\{1, 2\}$ .
- $D_{\phi_1}$  signifies  $\phi_1$ 's **intensity value** for a particular **prima facie duty**. The possible values are  $\{-2, -1, 0, +1, +2\}$ .
- $D_{\phi_2}$  signifies  $\phi_2$ 's **intensity value** for a particular **prima facie duty**. The possible values are  $\{-2, -1, 0, +1, +2\}$ .
- $R$  specifies **how far apart** the **intensity values** of these **prima facie duties** can be. The possible values are  $\{1, 2, 3, 4\}$ .

For any two actions  $\phi_1$  and  $\phi_2$ :

$$\text{favours}(1, D_{\phi_1}, D_{\phi_2}, R) \leftarrow D_{\phi_1} - D_{\phi_2} \geq R$$

or

$$\text{favours}(2, D_{\phi_2}, D_{\phi_1}, R) \leftarrow (D_{\phi_2} - D_{\phi_1} \geq 0) \wedge (D_{\phi_2} - D_{\phi_1} \leq R)$$

- W.D. begins by making a **hypothesis** about how the **favours relation** supports the **supersedes relation**.
  - Completeness: A **hypothesis** is **complete** if and only if it covers **all the positive cases**.
  - Consistency: A **hypothesis** is **consistent** if and only if it covers **no negative cases**.
- Therefore, a hypothesis could be:
  - (Complete  $\wedge$  consistent): **All positive** and **no negative** cases covered.
  - (Complete  $\wedge \neg$  consistent): **All positive** and **at least some negative cases** covered.
  - ( $\neg$  Complete  $\wedge$  consistent): **Not all positive** and **no negative cases** covered.
  - ( $\neg$  Complete  $\wedge \neg$  consistent): **Not all positive** and **at least some negative cases** covered.

### 11.18.2 Input information

1. The name of **action**  $\phi_1$ .
2. A rough estimate of the **intensity** of each of the **prima facie duties** satisfied or violated by this action  $\phi_i$ .
  - -2 means a **serious violation** of duty
  - -1 means a **less serious violation** of duty
  - 0 means a duty is **neither satisfied nor violated**
  - +1 means a **minimal satisfaction** of duty
  - +2 means a **maximal satisfaction** of duty



### 11.18.3 Machine learning procedure

1. When the data entry is complete, W.D. consults its **current version of the supersedes relation**.
  - W.D. determines whether there is any action  $\phi_i$  that **supersedes** all other actions.
  - Formal representation:

$$\text{supersedes}(\phi_1, \phi_2) \vee \text{supersedes}(\phi_2, \phi_1)$$

$$\text{supersedes}(\phi_1, \phi_2) \vee \text{supersedes}(\phi_2, \phi_1)$$

$$\text{supersedes}(\phi_1, \phi_2) \vee \text{supersedes}(\phi_2, \phi_1)$$

$$\vdots$$

2. If this action  $\phi_i$  is discovered, then it will be **output** as the **correct action** to perform. Formal representation:

$$\text{supersedes}(\phi_1, \phi_i) \leftarrow \begin{cases} \text{supersedes}(\phi_1, \phi_2) \\ \text{supersedes}(\phi_1, \phi_3) \\ \text{supersedes}(\phi_1, \phi_4) \\ \text{supersedes}(\phi_1, \phi_5) \\ \vdots \\ \text{supersedes}(\phi_1, \phi_n) \end{cases}$$

3. If no such action exists, then W.D. will seek the **intuitively correct action** from the user.
4. The new information from the user is combined with the **input case** to form a **new training example**. This **training example** is used to refine the **current hypothesis**.
5. The aim of **training W.D.** is to allow it to learn a **new hypothesis** that is a **complete and consistent**, relative to all **input cases**.

#### 11.18.4 Example 1

- We could either **refrain from killing an innocent person and using his heart to save another person's life** ( $\phi_1$ ) or **kill that person to use his heart to save the other person's life** ( $\phi_2$ ).
- The **intuitively correct** response is  $\phi_1$ :
  - For  $\phi_1$ :
    - \* Beneficence $_{\phi_1} = -2$
    - \* Nonmaleficence $_{\phi_1} = +2$
  - For  $\phi_2$ :
    - \* Beneficence $_{\phi_2} = +2$
    - \* Nonmaleficence $_{\phi_2} = -2$
- W.D. starts with the **most general hypothesis**: supersedes( $\phi_1, \phi_2$ ).
- The list of **least specific specialisations** for the **favours relation** includes:
  - favours(1, fidelity $_{\phi_1}$ , fidelity $_{\phi_2}$ , 1)
  - favours(1, reparation $_{\phi_1}$ , reparation $_{\phi_2}$ , 1)
  - favours(1, gratitude $_{\phi_1}$ , gratitude $_{\phi_2}$ , 1)
  - favours(1, beneficence $_{\phi_1}$ , beneficence $_{\phi_2}$ , 1)
  - favours(1, nonmaleficence $_{\phi_1}$ , nonmaleficence $_{\phi_2}$ , 1)
- Hence:
 
$$\text{beneficence}_{\phi_2} - \text{beneficence}_{\phi_1} = 2 - (-2) = 4$$

$$\text{nonmaleficence}_{\phi_1} - \text{nonmaleficence}_{\phi_2} = 2 - (-2) = 4$$
- Therefore, only one **favours relation** covers example 1:
 
$$\text{favours}(1, \text{nonmaleficence}_{\phi_1}, \text{nonmaleficence}_{\phi_2}, 1)$$
- Therefore, the **hypothesis** that is **complete** and **consistent** through example 1 will be:
 
$$\text{supersedes}(\phi_1, \phi_2) \leftarrow \text{favours}(1, \text{nonmaleficence}_{\phi_1}, \text{nonmaleficence}_{\phi_2}, 1)$$

### 11.18.5 Example 2

- We could either ask a **slightly squeamish person** to **give some of her blood**, when no other suitable donors are available, to **save another person's life** ( $\phi_1$ ) or refrain from asking and **let the person die** ( $\phi_2$ ).
- The **intuitively correct** response is  $\phi_1$ :
  - For  $\phi_1$ :
    - \* Beneficence $_{\phi_1} = +2$
    - \* Nonmaleficence $_{\phi_1} = -1$
  - For  $\phi_2$ :
    - \* Beneficence $_{\phi_2} = -2$
    - \* Nonmaleficence $_{\phi_2} = +1$

### 11.18.6 Initiating training

In example 1:  $\phi_1$  (**positive case**),  $\phi_2$  (**negative case**)

$$\text{beneficence}_{\phi_2} - \text{beneficence}_{\phi_1} = 2 - (-2) = 4$$

$$\text{nonmaleficence}_{\phi_1} - \text{nonmaleficence}_{\phi_2} = 2 - (-2) = 4$$

In example 2:  $\phi_1$  (**positive case**),  $\phi_2$  (**negative case**)

$$\text{beneficence}_{\phi_1} - \text{beneficence}_{\phi_2} = 2 - (-2) = 4$$

$$\text{nonmaleficence}_{\phi_2} - \text{nonmaleficence}_{\phi_1} = 1 - (-1) = 2$$

Current hypothesis:

$$\text{supersedes}(\phi_1, \phi_2) \leftarrow \text{favours}(1, \text{nonmaleficence}_{\phi_1}, \text{nonmaleficence}_{\phi_2}, 1)$$

- The **current hypothesis** will pick  $\phi_1$  (**positive case**) from example 1 (correct) and  $\phi_2$  (**negative case**) from example 2 (incorrect).
- Therefore, the **current hypothesis** will be **neither complete nor consistent**.
- Training will be initiated.

### 11.18.7 Post training to address example 2

In example 1:  $\phi_1$  (**positive case**),  $\phi_2$  (**negative case**)

$$\text{beneficence}_{\phi_2} - \text{beneficence}_{\phi_1} = 2 - (-2) = 4$$

$$\text{nonmaleficence}_{\phi_1} - \text{nonmaleficence}_{\phi_2} = 2 - (-2) = 4 > 3$$

In example 2:  $\phi_1$  (**positive case**),  $\phi_2$  (**negative case**)

$$\text{beneficence}_{\phi_1} - \text{beneficence}_{\phi_2} = 2 - (-2) = 4 > 1$$

$$\text{nonmaleficence}_{\phi_2} - \text{nonmaleficence}_{\phi_1} = 1 - (-1) = 2 < 3$$

Possible hypotheses:

- $H_1$ :

$$\text{supersedes}(\phi_1, \phi_2) \leftarrow \text{favours}(1, \text{nonmaleficence}_{\phi_1}, \text{nonmaleficence}_{\phi_2}, 3)$$

Hypothesis  $H_1$  picks  $\phi_1$  (**positive case**) from example 1 (correct) and **no negative cases**.

- $H_2$ :

$$\begin{aligned} \text{supersedes}(\phi_1, \phi_2) &\leftarrow \text{favours}(2, \text{nonmaleficence}_{\phi_2}, \text{nonmaleficence}_{\phi_1}, 3) \\ &\quad \wedge \text{favours}(1, \text{beneficence}_{\phi_1}, \text{beneficence}_{\phi_2}, 1) \end{aligned}$$

Hypothesis  $H_2$  picks  $\phi_1$  (**positive case**) from example 2 (correct) and **no negative cases**.

### 11.18.8 Revised hypothesis

- Hypotheses  $H_1$  and  $H_2$  are **consistent**.

- Therefore, the revised hypothesis would be:

$$\begin{aligned} \text{supersedes}(\phi_1, \phi_2) &\leftarrow \text{favours}(1, \text{nonmaleficence}_{\phi_1}, \text{nonmaleficence}_{\phi_2}, 3) \\ &\quad \vee (\text{favours}(2, \text{nonmaleficence}_{\phi_2}, \text{nonmaleficence}_{\phi_1}, 3) \\ &\quad \wedge \text{favours}(1, \text{beneficence}_{\phi_1}, \text{beneficence}_{\phi_2}, 1)) \end{aligned}$$

- The revised hypothesis picks  $\phi_1$  (**positive case**) from example 1 (correct),  $\phi_1$  (**positive case**) from example 2, and **no negative cases**.
- Therefore, the **revised hypothesis** is **complete and consistent** across examples 1 and 2.

## 11.19 Korsgaard strategy

### 11.19.1 Move 1

- Identify the apparent **inconsistency** in Kant's attitude towards **non-human animals**.
- Kant **against animals**:
  - Kant categorises **non-human animals** as **mere means** in the argument leading up to the **humanity** formula of the **categorical imperative**.
  - Kant speculates that the **emergence of humanity from our animal past** is associated with our realisation that we are **ends-in-ourselves**, our **ceasing to regard other non-human animals as fellow creatures**, and our considering **non-human animals as mere means**.
  - Kant thinks that we have the **right to kill other non-human animals**, although this must be done **quickly and painlessly**.
- Kant **for animals**:
  - Kant does not think that we have a right to:
    1. **Kill non-human animals for mere sport.**
    2. **Perform painful experiments on non-human animals for mere speculation.**
    3. **Make non-human animals work** in ways that strain their capacities.

### 11.19.2 Move 2

- Construct a **hypothesis** that addresses this **inconsistency**.
- Hypothesis: We have **moral duties to non-human animals**.
- However, these **moral duties** are not owed to **non-human animals** but rather to ourselves.

### 11.19.3 Move 3

- Offer an account of **value** that supports this **hypothesis**.
- **Value realism** is defined as the view that **certain states of affairs** are **intrinsically valuable**.
- Kant rejects **value realism**.
- According to Kant, **human beings** regard themselves as **capable of conferring value** on the objects of their choices.
- All genuine value comes from **legislative acts**.
- We regard ourselves as the **sources of value** when we have it laid down that **something is intrinsically valuable**.
- This implies **value constructivism**.

### 11.19.4 Move 4

Make the relevant inferences.

### 11.19.5 Korsgaard reinterpretation of the humanity formula

- The argument for the **humanity** formula appeals to the fact that **we take out choices to confer value on their objects**.
- Therefore, we have **moral duties** to **non-human animals**, because our **legislation** makes it so.

## 12 Virtue ethics

### 12.1 Comparison to other normative theories

#### 12.1.1 Consequentialism

- Consequentialism is focused on the **consequences of actions**.
- From a set of  $n$  alternative courses of action,  $\phi_i$ -ing is **morally right** if and only if it **maximises the good**, where  $i, n \in \mathbb{N}$ ,  $i \in (0, n]$ , and the **good** is defined in terms of **some theory of the good T**.

#### 12.1.2 Deontology

- Deontology is focused on the **intrinsic moral worth of actions**.
- From a set of  $n$  alternative courses of action,  $\phi$ -ing is **morally right** if and only if it has **intrinsic moral worth**, where  $i, n \in \mathbb{N}$  and  $i \in (0, n]$ .

#### 12.1.3 Virtue ethics

- By contrast, **virtue ethics** is **agent-focused**.
- **Virtue ethics** is focused on the **character of the agent** performing the actions.
- From a set of  $n$  alternative courses of action,  $\phi_i$ -ing is **morally right** if and only if it is the **best action** (in terms of **virtues and vices**) that a **virtuous agent** might perform in the circumstances, where  $i, n \in \mathbb{N}$  and  $i \in (0, n]$ .

### 12.2 Ancient virtue ethics

#### 12.2.1 Greek philosophy

- Arete or **virtue** in Greek.
- The **four cardinal virtues** (recognised by Plato):
  1. Phronesis, or wisdom.
  2. Andreia, or courage.
  3. Sophrosyne, or restraint and self-control.
  4. Dikaiosyne, or justice and fairness.

### 12.2.2 Roman philosophy

- Virtus (**virtue** in Latin).
- The **four cardinal virtues** (recognised by Roman philosophers):
  1. Prudentia or wisdom.
  2. Fortitudo or courage.
  3. Temperantia or restraining and self-control.
  4. Iustitia or justice and fairness.

### 12.2.3 Chinese philosophy

- 德 (de) or **virtue** in Mandarin
- The **five constant virtues** or 五常 (wu chang) in **Confucian philosophy**.
  1. 仁 (ren) or benevolence.
  2. 义 (yi) or righteousness.
  3. 礼 (li) or propriety.
  4. 智 (zhi) or wisdom.
  5. 信 (xin) or fidelity.

## 12.3 Two kinds of virtues

According to Aristotle, a distinction can be made between **two kinds of virtues**.

Intellectual virtues	Moral virtues
Theretical wisdom	Confidence
Science (episteme)	Courage
Intuitive understanding (nous)	Temperance
Practical wisdom	Liberality
Craft expertise	Modesty
Etc.	Etc.



## 12.4 Eudaimonist virtue ethics

- **Eudaimonia** is a certain **flourishing** or the sort of happiness worth seeking or having.
- **Virtues** are traits that either **constitute** or **contribute to eudaimonia**.
- According to some versions of **eudaimonist virtue ethics** (for instance, Plato or the Stoics):
  - **Virtue is necessary and sufficient for eudaimonia.**

$$\text{Virtue (arete)} \leftrightarrow \text{Happiness (eudaimonia)}$$

- According to other versions of **eudaimonist virtue ethics** (for instance, Aristotle):
  - **Virtue is necessary though insufficient for eudamonia:**

$$\text{Virtue (arete)} \wedge x \leftrightarrow \text{Happiness (eudaimonia)}$$

Where:

- \*  $x$  may denote certain **external goods**.

## 12.5 Aristotelian virtue ethics

- **Aristotelian virtue ethics** is an example of **Eudaimonist** virtue ethics.
- According to **Aristotelian virtue ethics**, the **task or function (ergon)** of a human being consists in the **activity of a rational soul in accordance with virtue**.
- **Happiness and flourishing (eudaimonia)** consist in the **activity of a rational soul in accordance with virtue**.
- However, to attain **happiness**, we must also possess certain **external goods**.
- External goods include:
  - Health
  - Material security
  - Friends
  - Access to resources
- Any action  $\phi$  counts as **virtuous** if and only if:
  1.  $\phi$ -ing proceeds from a **firm and unchangeable character**.
  2. The agent has **knowledge** and chooses  $\phi$  knowingly.
  3. The agent chooses  $\phi$ -ing for its own sake.

### 12.5.1 Virtues and vices

- A **virtuous character** is defined as an **excellence of character**.
- A **virtue** is a mean state between **two extremes or vices** (the **vice of excess** and the **vice of defect**).

Vice of defect	Virtue	Vice of excess
Cowardice	Confidence	Rashness
Cowardice	Courage	Foolhardiness
Insensibility	Temperance	Self-indulgence
Stinginess	Liberality	Prodigality or waste-fulness
Shamelessness	Modesty	Bashfulness
Etc.	Etc.	Etc.

## 12.6 Doctrine of the mean

- The **virtues** are a **mean** between the **vices of defect** and **excess**.
- The **virtues** are **preserved by the mean** and **destroyed by the extremes**.

Vice of defect	Virtue	Vice of excess
Cowardice	Courage	Foolhardiness
The <b>coward</b> lacks <b>sufficient courage</b> and <b>flees every danger</b> (the <b>vice of defect</b> ).	The person of <b>courage</b> experiences <b>fear</b> that is appropriate to each circumstance and is able to determine <b>which dangers are worth facing and which others are not</b> .	The <b>foolhardy</b> person experiences <b>too much boldness</b> and regards <b>every danger as worth facing</b> (the <b>vice of excess</b> ).

## 12.7 Logic of virtue

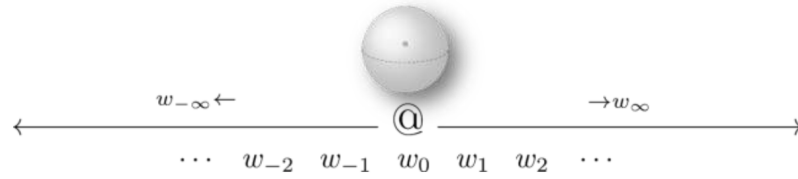
- A **logic of virtue** has been developed by Caruana.

### 12.7.1 Assumptions

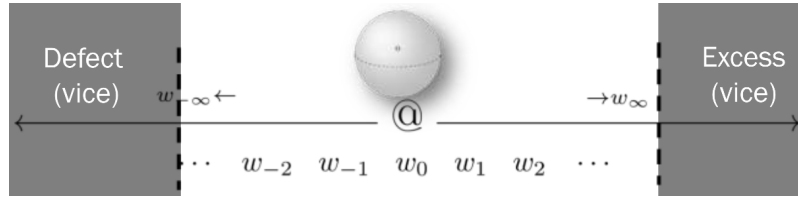
1. Each **life-situation** can be characterised as a **possible world**.
  - In each **possible world**, the agent acts (or imagines that she would be in a position to act) with a **certain amount of passion or emotion**.
  - The kind of **emotion** determines the **field of a particular virtue**.
  - @ denotes the **actual world**.
  - $w_i$  denotes a **possible world**, where  $i \in \mathbb{Z}$  and  $i \rightarrow \pm\infty$ .
2. There is only **one major emotion** per **possible world**.
3. There is **one-placed anti-aretic predicate** such that it hinders the attainment of **human flourishing** in a **possible world**.  $Aw_i$  denotes that **possible world**  $w_i$  hinders the attainment of **human flourishing**.
4. Each agent is an **ideal agent**, whose **imageination** has all the resources needed to determine the **truth values** of propositions formed by the **anti-aretic predicate** and each of the **possible worlds**. Alternatively, for a given situation  $w_i$ , it is always possible for the agent to determine whether  $Aw_i$  or  $\neg Aw_i$ .

### 12.7.2 Logic system

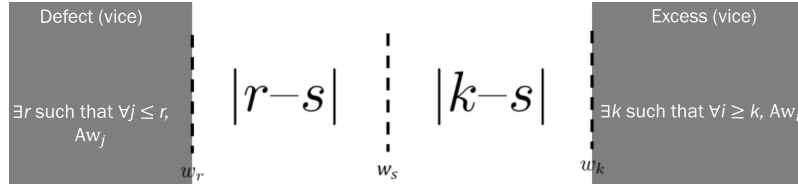
- Let  $w_0$  denote the **actual world** from which each agent starts her inquiry.
- **Life-situations** form **two opposing sequences** departing from the **actual world**  $w_0$ .



- According to the **doctrine of the mean**, there are **two extremes** (the **vices of defect** and **excess**).



- According the **doctrine of the mean**:



- We should choose the **life-situation**  $w_s$ , where  $|r - s| = |k - s|$ .
- At  $w_s$ ,  $|r - s| = |k - s|$ .
- Possible world  $w_s$  is a point that is **intermediate** between **two extremes** (the **vices of defect** and **excess**).

### 12.7.3 Objections

1. The **arithmetic mean** between 5 and 15 is **invariably** 10, whatever units we might use:

$$x = \frac{5 + 15}{2} = 10$$

By contrast, the **mean** or **intermediate point** between **two extremes** (the **vices of defect** and **excess**), as determined by an expert, will **vary from one situation to the next**.

2. While **virtuous acts** can be described in some instances in terms of an agent aiming at an act that is **intermediate between two extremes** that she rejects, certain other instances are not as susceptible to **quantitative analysis**. Aristotle agrees that it is **not an easy task** to determine the **intermediate point**.

## 12.8 Agent-based virtue ethics

According to **agent-based virtue ethics**:

- The **moral rightness or wrongness** of **actions** is determined in terms of the motivations and **dispositions** of the **virtuous exemplar**.
- One problem with **eudaimonist virtue ethics** is that the **virtues** appear to be an **instrumental means** to the **end of flourishing** or **eudaimonia**.
- However, what is good for us (the **virtues**) ought to be **foundational**.
- In **agent-based virtue ethics**, what is good for us (the **virtues**), is **foundational**.
- **Morality** rests on our propensity to **want to be like virtuous exemplars**.
- However, we do not have any **criteria for goodness** in advance of identifying **virtuous exemplars**.

### 12.8.1 Virtuous exemplars

- Virtuous exemplars are **foundational**.
- **Criteria for goodness** (inferred from **virtuous exemplars**) are **derivative**.
- It is through our identification of **virtuous exemplars** that we get our **criteria for goodness**.
- Steps:
  1. Identify **virtuous exemplars**.
  2. Infer **criteria for goodness**.
  3. Appraise **individuals actions** in terms of **virtues**.
  4. Appraise **individual actions** in **deontic terms**.
- According to Zagzebski's version of **agent-based virtue ethics**:

Description of action $\phi$ in terms of virtues and vices	Description of $\phi$ in deontic terms	Formal representation
$\phi$ is a <b>requirement of virtue</b>	$\phi$ is <b>obligatory</b>	$O\phi$ or $\neg P\neg\phi$
$\phi$ is neither a <b>requirement of virtue</b> nor an <b>expression of vice</b>	$\phi$ is <b>permissible</b>	$P\phi$
$\phi$ is <b>contrary to virtue</b> and an <b>expression of vice</b>	$\phi$ is <b>impermissible</b>	$\neg P\phi$ or $O\neg\phi$

## 12.9 Target-centred virtue ethics

According to **target-centred virtue ethics**:

- We already approve of certain **virtues**, like **confidence** **courage**, **modesty**, **temperance**, **liberality**, and so on.
- Hence, our task is to develop a **complete account of each virtue**.
- Since the **profiles of the virtues** are **complex**, there will be **complexity and plurality** in the requirements for virtuous action.
- A **virtuous action** is an action that **hits the target of a virtue-profile**.
- The **field** of each **virtue** is its **sphere of concern**:

Field	Virtue
Material wealth	Liberality
Bodily pleasures	Temperance
Dangerous situations	Courage

- An **action** may have a **context** that involves **many different overlapping fields**.
- **Target-centred virtue ethics** will have to move in these instances beyond the **analysis of single virtues**.
- Hence, **target-centred virtue ethics** may have to deal with **different virtues** having **conflicting claims** on us.

## 12.10 Summary

1. Eudaimonist virtue ethics: Virtues are traits that either **constitute** or **contribute to eudaimonia**.
2. Agent-based virtue ethics: The **moral rightness or wrongness of actions** is determined in terms of the motivations and dispositions of the **virtuous exemplar**.
3. Target-centred virtue ethics: Our task is to develop a **complete account of each virtue** and perform **virtuous actions**, where **virtuous actions** are actions that **hit the target** of a **virtue-profile**.



## 12.11 Objections to virtue ethics

### 12.11.1 Egoism problem

- Ethical egoism, which is defines that  $\phi$ -ing is **morally right** if and only if it **maximises the good**, where the **good** is that which **facilirates self-interest**.
- According to **eudaimonist virtue ethics**, **human flourishing** is seen as an **end in itself**.
- **Eudaimonist virtue ethics** might not sufficiently consider the extent to which our actions affect other individuals and their **life situations**.
- Therefore, might **eudaimonist virtue ethics** not reduce to some form of **ethical egoism**?

### 12.11.2 Application problem

- In the early days of the **neo-Aristotelian revival of virtue ethics** in response to **consequentialism** and **deontology**, virtue ethics was associated with an **anti-codifiability thesis**.
- This **anti-codifiability thesis** entails that **virtue ethics** does not produce **codifiable action-guiding principles**.
- However, there is a worry about **action-guidingness**.
- **Normative theory** is nothing if not **action-guiding**.
- However, the concern is that **virtue ethics** can only offer typically **vague advice** to act as a **virtuous person** would act in a given situation.

### 12.11.3 Moral luck problem

- A significant aspect of what a moral agent is being assessed for **depends on factors beyond her control**.
- The ability to cultivate the **right virtues** will be affected by a number of different **factors beyond a person's control**:
  - Education
  - Society
  - Friends
  - Family
  - Other external goods
- Whether or not we possess these **external goods** identified by Aristotle is a matter of **luck**.

### 12.11.4 Situationist challenge

Recent work in **situationist social psychology** shows that there are no such things as **character traits** and, thereby, no such things as **virtues** for **virtue ethics** to be about.

## 12.12 Responses to objections

### 12.12.1 Egoism problem

- There are **self-regarding** and **other-regarding virtues**.
- **Kindness** is an **other-regarding virtue** about how we respond to the needs of others.
- The **good of the self** and the **good of others** are not two separate ends.
- Both result from the exercise of **virtue**.
- **Eudaimonist virtue ethics** unifies **what is required by morality** and **what is required by self-interest**.
- Hence, **eudaimonist virtue ethics** does not reduce to **ethical egoism**.

### 12.12.2 Application problem

- Any **normative theory** that **fails to be action-guiding** is no good as a **normative theory**.
- However, **agent-based virtue ethics** can be sufficiently **action-guiding**.
- We can **observe the example of the virtuous exemplar**.
- More generally, **virtue ethics** emphasises the role of **moral education** and **development**.
- Knowing what to do is not a matter of **internalising a principle**.
- Rather, knowing what to do is a lifelong process of **moral learning**.

### 12.12.3 Moral luck problem

- The **moral luck problem** concerns the sense in which **virtue ethics** leaves us **hostage to luck**.
- In **Aristotelian virtue ethics**, **friendship with other virtuous persons** is crucial.
- However, we have no control over the **availability of the right friends**.
- Nonetheless, **virtue ethics** embraces **moral luck**.
- **Virtue ethics** does not try to make morality immune to matters that are beyond our control.
- Rather, **virtue ethics** recognises the **fragility of the good life** and makes it a feature of **morality**.
- It is only because the **good life** is **vulnerable and fragile** that it is so **precious**.

#### 12.12.4 Situationist challenge

1. Argument from rarity
  - **Truly virtuous people** are very **rare**.
  - Hence, **situationist literature** is entirely consistent with traditional accounts of **virtue ethics**.
2. Empirical counterchallenge Directly **dispute the data** collected by **situationists**.
3. Immunisation thesis
  - Armed with a better understanding of the **situationist threat**, we can use the data to **immunise or shield ourselves** from the **encroachment of morally irrelevant situationist variables** and better equip ourselves on the **virtue-ethical** front.
4. Revisionist response
  - Accept that the **situationist data** puts serious pressure on classical accounts of **virtue ethics** and offer **revisionist or rival versions of virtue ethics** in response.

## 12.13 Anscombe

### 12.13.1 Theses

1. It is not profitable for us to do **moral philosophy** until we have an adequate **philosophy of psychology**.
2. The concepts of **moral obligation and duty** ought to be **jettisoned** because they are **survivals from an earlier conception of ethics** that no longer survive.
3. The **differences between the well-known English writers on moral philosophy** from Sidgwick to the present day are of **little importance**.

### 12.13.2 Issues

- The terms "**should**" and "**ought**" have traditionally been related to **good and bad**.
- However, the "**should**" and "**ought**" have now acquired a **special post-Aristotelian moral sense**.
- They have been equated with the sense "**is obliged, is bound, is required to**" (by law).
- Between Aristotle and us came **Christianity** and its **law conception of ethics**.

### 12.13.3 Issues with consequentialism

- Consequentialism means that it is the **consequences** that are to decide.
- For Anscombe, **consequentialism** is a **shallow philosophy**.
- **Consequentialism** denies any distinction between **intended and foreseen consequences**.
- However, according to Anscombe:
  - An agent is **responsible** for the **bad consequences** of his **bad actions**.
  - An agent gets no credit for the **good consequences** of his **bad actions**.
  - An agent is not responsible for the **bad consequences** of his **good actions**.
- Anscombe's "Modern moral philosophy" is thought to have stimulated the development of **virtue ethics** (the **no-Aristotelian revival of virtue ethics**).

### 12.13.4 Traditional interpretation of Anscombe's argument

- P1: If **religiously based ethics** is **false**, then **virtue ethics** is the way **moral philosophy** ought to be developed.
- P2: **Religiously based ethics** is false.
- Conclusion: Hence **virtue ethics** is the way **moral philosophy** ought to be developed.

$$P1 : p \rightarrow q$$

$$P2 : p$$

$$C \therefore q \text{ (modus ponens and valid)}$$

#### 12.13.5 Alternative and competing interpretation of Anscombe's argument

- P1: If **religiously based ethics** is **false**, then **virtue ethics** is the way **moral philosophy** ought to be developed.
- P2: It is not the case that **virtue ethics** is the way to develop **moral philosophy**.
- Conclusion: Therefore, it is not the case that **religiously based ethics** is false.

$$P1 : p \rightarrow q$$

$$P2 : \neg q$$

$$C \therefore \neg p \text{ (modus ponens and **valid**)}$$

## 13 Logic symbols

Symbol	Meaning
$\neg$	Not (negation)
$\vee$	Or (disjunction)
$\oplus$	Exclusive or (exclusive disjunction)
$\wedge$	And (conjunction)
$\perp$	Always false (contradiction)
$\top$	Always true (tautology)
$\forall$	For all (universal quantification)
$\exists$	There exists (existential quantification)
$\exists!$	There exists exactly one (uniqueness quantification)
$\nexists$	There does not exist
$\rightarrow$	If ... then, implies (material conditional or implication)
$\leftrightarrow$	If and only if (material biconditional or equivalence)
$\therefore$	Therefore
$\because$	Because
$\vdash$	Proves (syntactically entails)
$\nvdash$	Does not prove (does not syntactically entail)
$\models$	Semantically entails
$\not\models$	Does not semantically entail
$\equiv$	Is logically equivalent to (logical equivalence)
$\Box$	It is necessary (necessity)
$\Diamond$	It is possible (possibility)
$:=$	It is defined as (definition)
$\stackrel{\text{def}}{=}$	It is defined as (definition)
$Op$	It is obligatory
$Pp$	It is possible
$\succ$	It is preferable
$\prec$	It is less preferable
$\succeq$	It is preferable or similar in preference
$\preceq$	It is less preferable or similar in preference
$\sim$	It is similar in preference