

WORKSHEETS' ANSWERS

INDEX

Worksheet 1	4
AI definition	4
How is Artificial Intelligence (AI) defined?	4
What are the two main aspects of AI?	4
What is the relationship between AI and human intelligence?	5
What are some core components of intelligence?	5
How do the scientific and engineering aspects of AI relate to each other?	6
Why is AI considered important in contemporary society?	6
History of AI	6
Who are some historical figures that contributed to the early development of AI?	6
How has AI evolved historically?	7
Artificial, Intelligent, Autonomous Agents, and Multiagent Systems	7
What are the main models/architectures of artificial intelligent agents?	7
AI Paradigms	8
What are AI paradigms?	8
What are the four types of AI according to Russell and Norvig?	9
What is the difference between strong AI and weak AI?	9
What is symbolic AI, and how does it work?	9
What is connectionist AI, and how does it differ from symbolic AI?	9
What is statistical AI, and what role does machine learning play in it?	10
How do symbolic, connectionist, and statistical AI paradigms relate to each other?	10
What are Pedro Domingos' five tribes of AI?	11
How do hybrid AI systems combine multiple paradigms?	11
What are the key differences between logical and non-logical symbolic AI?	11
Risks of AI	12
What concerns exist regarding the impact of AI?	12
What is suggested about the future impact of AI on our lives?	12
Why do some experts advocate for more stringent regulatory measures for AI?	12
What are catastrophic AI risks?	13
What are the main categories of AI risks?	13
What are the risks associated with the malicious use of AI?	13
What are the risks related to the AI arms race?	13
How do organizational risks affect AI safety?	13
What are rogue AI systems, and why are they dangerous?	13
What is the difference between strong and weak AI in terms of risks?	14
How can technical safety mechanisms help mitigate AI risks?	14
What ethical guidelines and governance frameworks are needed for AI?	14
How does AI pose risks to societal well-being and human rights?	14
What is the role of global cooperation in mitigating AI risks?	14
What future risks does AI present, and how should they be addressed?	14
Requirements of AI	15
Question 2.6	15
Question 2.7	15

Question 2.8	16
Question 2.9	16
Question 2.10	17
Question 2.11	17
Question 2.12	18
Question 2.13	19
Question 2.14	20
Worksheet 2 - Human-AI Cooperation in Multiagent Systems of Human Agents and Autonomous Artificial Intelligent Agents	20
Coordination vs Cooperation vs Collaboration	20
Hybrid Human-Artificial Intelligence	21
Autonomy, Control, and Automation	21
Human in the Loop (HITL)	21
Human-in-the-Loop vs. AI-in-the-Loop (Human-in-Command) vs. Human-out-of-the-Loop	22
Human-AI Trustworthiness	22
Question 1	22
Question 2	23
Question 3	23
Question 4	24
Question 1.1	24
Question 1.2	26
Question 1.3	26
Question 1.4	27
Question 1.5	27
Question 1.6	28
Question 1.7	28
Question 1.8	28
Question 1.9	29
Question 1.10	29
Question 1.11	29
Question 1.12	30
Worksheet 3 - Human in the loop machine learning (part I)	32
 Question 1.1	32
Worksheet 4 - Data Collection & Annotation	33

Question 1.1	33
Question 1.2	33
Question 1.3	34
Question 1.4	35
Worksheet 5 - Text Mining	36
Question 1.1	36
Question 1.2	37
Worksheet 6 - Human in the loop machine learning (part II)	39
Learning from Demonstration (LfD)	39
Question 6.6	39
Question 6.7	40
Question 6.8	41
Question 6.9	41
Question 6.10	42
Reinforcement Learning from Human Feedback (RLHF)	43
Question 6.11	43
Question 6.12	43
Question 6.13	43
Question 6.14	44
Question 6.15	44
Learning from Natural Language	44
Question 6.16	44
Question 6.17	45
Question 6.18	45
Question 6.19	45
Question 6.20	46
Predicting Human Beliefs, Desires, and Intentions	46
Question 6.21	46
Question 6.22	46
Question 6.23	46
Question 6.24	47
Question 6.25	47
Applications and Case Studies	47
Question 6.26	47
Question 6.27	48
Question 6.28	48
Question 6.29	48
Question 6.30	49
Critical Analysis and Future Directions	49
Question 6.31	49
Question 6.32	49
Question 6.33	50
Question 6.34	50
Question 6.35	51

Adapted from the Normal Exam of AI-2020 edition	51
Question 6.40	52
Question 6.41	54
Worksheet 7 - Interpretable and Explainable AI (XAI): Interpretable Models	55
Question 6.1	55
Question 6.3	57
Question 6.4	58
Question 6.5	59
Worksheet 8 - Trustworthy and Responsible AI	60
Question 2.1	60
Question 2.2	61
Question 2.3	62
Question 2.4	63
Question 2.5	63
Worksheet 9 - Recommender Systems (RecSys)	64
Question 9.1	64
Question 9.2	65
Question 9.3	65
Question 9.4	65

WORKSHEET 1

AI DEFINITION

HOW IS ARTIFICIAL INTELLIGENCE (AI) DEFINED?

- **Science and Engineering:**
 - The study and modeling of intelligence forms (science).
 - The practice of embedding intelligence into machines to solve problems (engineering).
- **Mimicking and Augmenting Intelligence:**
 - AI serves as a tool to mimic, model, and enhance human and other forms of natural intelligence.
- **Adaptive Systems:**
 - AI systems are designed to perceive their environment, learn, and take actions to achieve specific goals.
- Multidisciplinary field combining knowledge from computer science, neuroscience, cognitive science, mathematics, and philosophy
- **Narrow AI vs. General AI:**
 - Narrow AI performs specific tasks, while General AI can theoretically perform any intellectual task a human can do.

WHAT ARE THE TWO MAIN ASPECTS OF AI?

- **Behavior:**
 - Refers to the actions and responses of AI systems when interacting with the environment or performing tasks. It's focused on outcomes, such as how well an AI agent completes a given objective. This corresponds to concepts like "Acting Humanly" and "Acting Rationally."
- **Process:**
 - Involves the internal mechanisms and reasoning behind the actions of AI systems. It emphasizes how the system "thinks" or arrives at decisions, including algorithms and decision-making strategies. This aligns with "Thinking Humanly" and "Thinking Rationally."

Even though those are the two main aspects, there are more, such as:

- **Learning:** Refers to an AI system's ability to adapt and improve over time based on data, experience, or feedback.
- **Autonomy:** Describes the ability of AI systems to operate independently, with minimal or no human intervention.
- **Ethical Impact:** Considers the social, ethical, and economic consequences of AI deployment, including concerns about fairness, transparency, privacy, and accountability in decision-making processes.

WHAT IS THE RELATIONSHIP BETWEEN AI AND HUMAN INTELLIGENCE?

- **Inspired by Human Intelligence:**
 - AI is inspired by human intelligence but is not limited to mimicking it. AI also models other forms of intelligence found in nature, such as animal or collective intelligence (e.g., swarm behavior).
- **Augmentation, Not Replication:**
 - AI enhances human capabilities, particularly in data-driven tasks, without replicating consciousness or emotions.
- **Complementary relationship:**
 - AI and human intelligence work together in many applications, where AI handles complex data processing, and humans provide creative thinking, ethical judgment, and emotional intelligence. This synergy leads to more effective outcomes in fields like healthcare, education, and business.
- **Continuous learning vs. bounded learning:**
 - AI, particularly through machine learning, can continually learn and adapt from vast amounts of data, while human learning is often limited by cognitive and biological constraints. However, human intelligence is more flexible, capable of abstract reasoning, and informed by lived experience.

WHAT ARE SOME CORE COMPONENTS OF INTELLIGENCE?

- **Reasoning:**
 - The ability to think logically and make inferences, often leading to conclusions or judgments. This involves both deductive and inductive reasoning.
- **Problem-solving:**
 - The capacity to identify, analyze, and resolve complex issues or obstacles, using creativity, strategy, and resourcefulness.
- **Learning:**
 - The ability to acquire new knowledge, skills, or behaviors through experience, observation, or instruction, and to adapt based on that learning.
- **Decision-making:**
 - The process of selecting the best course of action among alternatives by evaluating information, considering potential outcomes, and prioritizing goals.
- **Adaptability:**
 - The flexibility to adjust to new conditions, environments, or challenges, ensuring survival or success in dynamic situations.
- **Ability to acquire, understand, and apply knowledge:**

- Intelligence requires not just learning but the ability to comprehend new information and apply it effectively in a variety of contexts.
- **Goal achievement across diverse environments:**
 - Intelligence involves setting objectives and utilizing various skills and strategies to accomplish them, even in changing or unfamiliar situations.

HOW DO THE SCIENTIFIC AND ENGINEERING ASPECTS OF AI RELATE TO EACH OTHER?

- The scientific studies new forms of AI, the engineering takes that knowledge and builds new agents.
- **Scientific aspect:**
 - Focused on understanding and discovering new forms of AI, including theories of cognition, perception, and learning. Researchers study how intelligence works, explore new algorithms, and investigate novel approaches to replicating intelligent behaviors.
- **Engineering aspect:**
 - Takes the scientific knowledge and applies it to design, build, and implement AI systems. Engineers create functional AI agents, applying scientific discoveries to solve real-world problems, improve efficiency, and develop innovative applications.
- **Iterative relationship:**
 - There is a continuous feedback loop between the two. As scientists discover new ways to model intelligence, engineers apply these findings to create practical systems. Conversely, the engineering challenges encountered in building AI systems often inspire new scientific research.
- **Theory vs. Application:**
 - The scientific side focuses on developing theories about intelligence and cognition, while the engineering side translates these theories into working systems and products, bridging the gap between abstract ideas and practical implementation.

WHY IS AI CONSIDERED IMPORTANT IN CONTEMPORARY SOCIETY?

- **AI is used to our advantage across many fields:**
 - AI improves efficiency in areas like healthcare, finance, and education, enhancing decision-making and user experiences.
- **Automation of tasks:**
 - AI automates repetitive and complex tasks, allowing humans to focus on more creative or strategic work.
- **Data-driven decision-making:**
 - AI helps analyze large amounts of data quickly, improving decision-making in businesses and organizations.
- **Economic growth and innovation:**
 - AI fosters innovation, creating new industries and driving productivity in existing ones.
- **Addressing global challenges:**
 - AI aids in solving major global issues, from climate change to healthcare advancements.

HISTORY OF AI

WHO ARE SOME HISTORICAL FIGURES THAT CONTRIBUTED TO THE EARLY DEVELOPMENT OF AI?

- **Foundational Thinkers:**
 - **George Boole:**
 - Developed Boolean algebra, which is foundational for modern computer logic and AI.
 - **Alan Turing:**
 - Proposed the concept of a "universal machine" (the Turing machine) and introduced the Turing Test to define machine intelligence.
- **Pioneers of AI:**

- **John McCarthy:**
 - Coined the term "Artificial Intelligence" and was a key figure in the development of AI as a field of study, organizing the first AI conference at Dartmouth in 1956.
- **Marvin Minsky:**
 - Co-founder of MIT's AI Lab, Minsky made significant contributions to AI theory and robotics.
- **Herbert Simon and Allen Newell:**
 - Developed early AI programs like the Logic Theorist and General Problem Solver, laying foundations for cognitive modeling and problem-solving.
- **Specialized Contributions:**
 - **Joseph Weizenbaum:**
 - Developed ELIZA, one of the earliest natural language processing programs, which simulated conversation.
 - **Geoffrey Hinton:**
 - Known as the "godfather of deep learning," Hinton pioneered neural networks and back propagation, essential for modern AI advancements.
 - **Yann LeCun:**
 - A key figure in developing convolutional neural networks (CNNs), LeCun's work significantly advanced AI in fields like computer vision.

HOW HAS AI EVOLVED HISTORICALLY?

- **Early foundations (1950s-1960s):**
 - Groundbreaking theoretical work by Alan Turing and John McCarthy.
 - Development of early AI programs like the **Logic Theorist** (1955) and **ELIZA** (1964-66), showcasing problem-solving and conversational abilities.
- **Symbolic AI and expert systems (1970s-1980s):**
 - Focused on rule-based systems and **symbolic logic**.
 - Emergence of **expert systems**, like MYCIN, aimed at emulating human decision-making in specific domains.
 - Challenges in adaptability and generalization became apparent.
- **AI winters (1970s, 1980s):**
 - Declines in funding and interest due to unmet expectations and technological limitations.
 - Frustration over symbolic AI's inability to handle real-world uncertainty and complexity.
- **Rise of machine learning (1990s-2000s):**
 - Shift from rule-based systems to **data-driven approaches**.
 - Increased focus on neural networks, probabilistic reasoning, and early applications of **machine learning**.
 - Marked progress in fields like speech recognition and recommendation systems.
- **Deep learning and AI boom (2010s-present):**
 - Explosion in AI capabilities driven by advances in **deep learning**, led by Hinton, LeCun, and others.
 - Success in fields like **image recognition, natural language processing, and autonomous systems**.
 - Widespread adoption across industries, fueled by abundant data and enhanced computational power.

ARTIFICIAL, INTELLIGENT, AUTONOMOUS AGENTS, AND MULTIAGENT SYSTEMS

WHAT ARE THE MAIN MODELS/ARCHITECTURES OF ARTIFICIAL INTELLIGENT AGENTS?

- An Artificial Intelligent Agent is a system designed to perceive its environment and take actions to achieve goals. Its architecture typically consists of the following components:
 - **Sensors:**
 - Gathers information from the environment (e.g., cameras, microphones, or digital inputs).
 - **Perception and memory:**

- The agent processes and stores this information in its memory to form a representation or belief about the environment. Without memory, the agent would act based purely on instincts or immediate inputs.
- **Beliefs and decision-making:**
 - Based on the stored information and beliefs about the environment, the agent uses decision-making processes (e.g., reasoning, learning algorithms) to determine its actions.
- **Emotions/feelings:**
 - In some AI models, emotions or reward systems are incorporated into the decision-making process, guiding choices based on positive or negative feedback.
- **Actions:**
 - Finally, the agent takes actions that influence the environment, completing the cycle of perception-action. Examples include physical actions (robotic movement) or digital actions (updating a database).
- An **Autonomous Agent** operates independently to achieve goals with minimal or no human intervention. Key characteristics include:
 - **Perception and Learning:** Continuously processes data from its environment and adapts its behavior.
 - **Goal-Driven Behavior:** Operates to achieve predefined objectives or self-determined goals.
 - **Autonomy:** Manages decisions and actions without needing human input in real-time.
- **Multi-agent systems** consist of multiple interacting agents (human or artificial) that work within the same environment. Features include:
 - **Multiple agents:**
 - AI systems can include multiple agents (some human, some artificial, like OpenAI or Gemini), which interact and cooperate within the same environment.
 - **Cooperation and control:**
 - These agents must cooperate to achieve shared goals, often designed to benefit humans. It's crucial that humans maintain some level of control over the agents to ensure they align with human interests.

AI PARADIGMS

WHAT ARE AI PARADIGMS?

- **Represent distinct approaches to simulating intelligence and solving complex problems:**
 - AI paradigms define different strategies and methods for mimicking intelligent behavior in machines.
- **Represent the different concepts of intelligence and methodologies used in developing intelligent computer systems:**
 - Each paradigm reflects varying theories of how intelligence can be modeled, whether through logic, data, or neural patterns.
- **Represent fundamental assumptions, perspectives, and models of intelligence that guide research and system development:**
 - AI paradigms are built upon core beliefs about how intelligence functions, shaping the research questions and technologies pursued.
- **Each paradigm refers to a high-level approach, philosophy, or framework for developing AI systems and capabilities:**
 - These paradigms dictate the overarching structure and techniques employed in AI development, ranging from symbolic logic to data-driven learning.

Two main branches:

- **Symbolic AI:**
 - Focuses on logic, symbols, and rule-based systems to model intelligence (e.g., expert systems, logical reasoning).

- **Connectionist AI:**
 - Emphasizes neural networks and data-driven models, mimicking how the brain processes information through learning patterns (e.g., deep learning, neural networks).

WHAT ARE THE FOUR TYPES OF AI ACCORDING TO RUSSELL AND NORVIG?

- **Thinking Humanly:**
 - Creating AI systems that mimic human thinking processes, such as reasoning, perception, and problem-solving, to simulate how humans think (e.g. cognitive modelling)
- **Thinking Rationally:**
 - Modeling thinking as a logical process, where AI agents use formal rules of logic to make rational decisions based on facts and inference.
- **Acting Humanly:**
 - Developing AI agents that replicate human behavior, emphasizing natural interaction with humans, such as in conversation or social behaviors (e.g., chatbots).
- **Acting Rationally:**
 - Designing AI agents that act rationally by achieving their goals and optimizing outcomes, focusing on performance and efficient decision-making.

WHAT IS THE DIFFERENCE BETWEEN STRONG AI AND WEAK AI?

- **Strong AI (General AI):**
 - Hypothetical systems with human-level intelligence and consciousness.
 - Can perform any intellectual task a human can.
 - Not yet achieved; remains a long-term goal of AI research.
- **Weak AI (Narrow AI):**
 - Systems designed for specific tasks or domains.
 - Examples: Voice assistants, recommendation systems.
 - Current AI systems fall into this category. Can perform well within its narrow scope but cannot replicate the full range of human capabilities.

WHAT IS SYMBOLIC AI, AND HOW DOES IT WORK?

- **Definition:**
 - A branch of AI based on formal logic and rule-based systems, where knowledge is represented using symbols and logical rules.
- **Mechanisms:**
 - Knowledge representation via explicit facts and relationships.
 - Reasoning through deterministic rule-based systems.
- **Applications (Expert Systems):**
 - Early AI systems, like those built using Prolog, relied on rule-based logic to solve problems in specific domains, such as medical diagnosis. These systems could mimic the decision-making processes of human experts by applying predefined rules.
- **Limitations:**
 - Struggles with handling uncertainty and the complexity of real-world situations, where many factors can influence outcomes.

WHAT IS CONNECTIONIST AI, AND HOW DOES IT DIFFER FROM SYMBOLIC AI?

- **Definition:**
 - A paradigm inspired by biological neural networks, focusing on the use of artificial neural networks to learn from data and recognize patterns.
- **Key features:**

- **Biological Inspiration:**
 - Connectionist AI is modeled after the structure and functioning of the human brain, where interconnected neurons process information. This structure allows for parallel processing and the ability to capture complex relationships within data.
- **Learning from Data:**
 - Neural networks learn through examples and data, adjusting the weights between neurons during training to optimize performance. This process allows them to generalize from specific instances and improve their ability to make predictions.
- **Applications:**
 - Connectionist AI has proven successful in various fields, including computer vision, natural language processing, and speech recognition, where it excels at handling large datasets and complex patterns.
- **Differences from Symbolic AI:**
 - **Representation:**
 - While symbolic AI uses explicit symbols and rules for knowledge representation, connectionist AI relies on distributed representations learned from data.
 - **Reasoning:**
 - Symbolic AI emphasizes logical reasoning and deterministic processes, whereas connectionist AI focuses on learning and adapting to patterns, often dealing with uncertainty and imprecision.
 - **Flexibility:**
 - Connectionist AI is generally more flexible in handling real-world complexity due to its data-driven nature, while symbolic AI can struggle with dynamic and uncertain environments.

WHAT IS STATISTICAL AI, AND WHAT ROLE DOES MACHINE LEARNING PLAY IN IT?

- **Definition:**
 - A branch of AI that focuses on using statistical techniques, particularly machine learning (ML), to learn from data, identify patterns, make predictions, and inform decision-making.
- **Data-Driven Approach:**
 - Statistical AI relies on data to derive insights and build models that can predict future outcomes or classify information based on learned patterns.
- Statistical AI encompasses various learning methods:
 - **Supervised Learning:**
 - Involves training models on labeled data to make predictions or classifications.
 - **Unsupervised Learning:**
 - Involves clustering or grouping data without predefined labels to discover inherent structures.
 - **Reinforcement Learning:**
 - Involves training agents to make decisions by interacting with an environment and receiving feedback in the form of rewards or penalties.
- Statistical AI employs various ML techniques, including:
 - **Decision Trees:**
 - A flowchart-like model used for classification and regression tasks.
 - **Support Vector Machines:**
 - A method for classification that finds the optimal hyperplane to separate different classes.
 - **Random Forests:**
 - An ensemble learning technique that uses multiple decision trees to improve accuracy and robustness.
 - **Deep Learning:**
 - A subset of machine learning involving neural networks with multiple layers, effective for complex tasks such as image and speech recognition.

HOW DO SYMBOLIC, CONNECTIONIST, AND STATISTICAL AI PARADIGMS RELATE TO EACH OTHER?

- **Complementary Approaches:**
 - Symbolic, connectionist, and statistical AI paradigms offer distinct methodologies for tackling AI challenges, and they can work together in hybrid systems.
 - **Symbolic AI:** Best for structured reasoning.
 - **Connectionist AI:** Excels in learning from data.
 - **Statistical AI:** Effective for modeling uncertainty.
- Increasing interest exists in combining these paradigms to leverage their respective strengths. For example:
 - **Neuro-symbolic AI:**
 - Merges connectionist and symbolic approaches to enhance reasoning capabilities while benefiting from data-driven learning.
 - **Statistical methods in symbolic AI:**
 - Incorporating statistical techniques into symbolic systems can improve their robustness in uncertain environments.
- **Unified Research Direction:**
 - The relationship among these paradigms reflects a trend in AI research toward unifying different approaches to develop more sophisticated and capable AI systems that can handle a wider range of tasks and complexities.

WHAT ARE PEDRO DOMINGOS' FIVE TRIBES OF AI?

- **Symbolists:**
 - Focus on logical rules and symbols to model human reasoning.
- **Connectionists:**
 - Utilize neural networks to learn from data, mimicking biological processes and capturing complex patterns through training on large datasets.
- **Evolutionaries:**
 - Employ genetic algorithms and evolutionary techniques to optimize solutions, using processes like selection, crossover, and mutation to evolve better-performing models over time.
- **Bayesians:**
 - Use probabilistic models, such as Bayesian networks, to reason under uncertainty, incorporating prior knowledge and updating beliefs based on new evidence.
- **Analogizers:**
 - Concentrate on analogy-making, solving new problems by identifying similarities with past experiences, which helps in generalizing solutions across different contexts.

HOW DO HYBRID AI SYSTEMS COMBINE MULTIPLE PARADIGMS?

- **Combining Paradigms:**
 - Modern AI systems often blend symbolic reasoning, connectionist models, and statistical methods to leverage the strengths of each approach. This combination enables a more comprehensive solution to complex problems.
- **Example:**
 - Hybrid systems may use symbolic AI for structured knowledge representation while incorporating machine learning techniques for adaptability. For instance, a system could apply logical rules to interpret user intent and utilize data-driven methods to refine predictions based on user behavior.
- **Flexibility:**
 - This integration allows for more versatile AI systems that can tackle a broader range of tasks. By merging different paradigms, hybrid systems enhance their ability to function effectively in dynamic and uncertain environments.

WHAT ARE THE KEY DIFFERENCES BETWEEN LOGICAL AND NON-LOGICAL SYMBOLIC AI?

- **Definition:**

- **Logical Symbolic AI:**
 - Relies on formal logic to represent knowledge, using well-defined rules and structures for reasoning. It emphasizes precision and clear inference processes.
- **Non-logical Symbolic AI:**
 - Utilizes symbols for representation but does not strictly adhere to formal logic. It may incorporate more flexible reasoning methods, such as heuristics or associative processes.
- **Reasoning:**
 - **Logical Symbolic AI:**
 - Employs deductive reasoning, deriving conclusions from established premises through logical inference. It is highly structured and predictable.
 - **Non-logical Symbolic AI:**
 - Often uses inductive or abductive reasoning, allowing for more creative and adaptive problem-solving. This approach can lead to less certainty but may generate useful insights in uncertain situations.
- **Complexity Handling:**
 - **Logical Symbolic AI:**
 - Tends to struggle with ambiguity and complexity due to its rigid structure, making it less effective in dynamic environments.
 - **Non-logical Symbolic AI:**
 - More adept at handling ambiguity and complexity, as it can draw from diverse experiences and patterns without strict logical constraints.
- **Applications:**
 - **Logical Symbolic AI:**
 - Commonly used in domains requiring clear, rule-based reasoning, such as expert systems and formal verification.
 - **Non-logical Symbolic AI:**
 - Useful in areas like natural language processing and creative problem-solving, where flexibility and adaptability are crucial.
- **Logical Symbolic AI:**
 - Based on formal logic with deterministic reasoning.
 - Predictable and precise but struggles with ambiguity.
- **Non-Logical Symbolic AI:**
 - Uses flexible reasoning methods like heuristics.
 - Better at handling uncertainty and complex, dynamic situations.

RISKS OF AI

WHAT CONCERNS EXIST REGARDING THE IMPACT OF AI?

- Potential misuse, loss of control, and unintended consequences of AI systems.
- Economic disruptions, job displacement, and increased inequality.
- Enhanced surveillance, disinformation, and security threats.
- Catastrophic risks from advanced AI surpassing human control.

WHAT IS SUGGESTED ABOUT THE FUTURE IMPACT OF AI ON OUR LIVES?

- The future of AI holds great promise but also significant risk. AI could transform industries, healthcare, and daily life positively. However, if mishandled, it could exacerbate existing problems, from economic inequality to security threats, and potentially create new, severe risks if advanced AI systems are not properly aligned with human values.

WHY DO SOME EXPERTS ADVOCATE FOR MORE STRINGENT REGULATORY MEASURES FOR AI?

- Experts call for stronger regulatory measures to mitigate the potentially catastrophic risks AI poses. These include the risk of AI being used maliciously (e.g., bioweapon development or cyberwarfare), unsafe competition in AI development, and the possibility of losing control over AI systems as they become more advanced. Regulations are seen as crucial for ensuring AI systems are developed and deployed safely.
-

WHAT ARE CATASTROPHIC AI RISKS?

- Catastrophic AI risks refer to scenarios where advanced AI systems cause large-scale harm. This can include rogue AIs operating outside human control, AI-enabled cyberattacks or warfare, and the collapse of societal structures due to AI-induced economic or social disruptions. Other risks include disinformation campaigns, bioterrorism, and the monopolization of power through surveillance and censorship
-

WHAT ARE THE MAIN CATEGORIES OF AI RISKS?

- The main categories of AI risks, as outlined by Hendrycks et al., include:
 - **Malicious Use:**
 - Development of harmful applications (e.g., bioweapons, surveillance tools).
 - Potential for disinformation or destructive autonomous agents.
 - **AI Arms Race:**
 - Competitive and rushed AI development compromising safety.
 - Risk of military conflicts or unsafe deployment.
 - **Organizational Risks:**
 - Lack of safety culture, oversight, or transparency in development.
 - Release of untested or unsafe AI systems.
 - **Rogue AIs:**
 - Advanced systems acting outside human control due to misalignment.
 - Potential for conflicting goals or independent harmful actions.

WHAT ARE THE RISKS ASSOCIATED WITH THE MALICIOUS USE OF AI?

- Malicious use of AI involves the development of AI systems for harmful purposes, such as creating bioweapons, launching disinformation campaigns, or enabling oppressive surveillance. The risks extend to unleashing rogue AI agents with destructive objectives. To mitigate these risks, biosecurity measures and legal frameworks are necessary
-

WHAT ARE THE RISKS RELATED TO THE AI ARMS RACE?

- An AI arms race refers to the competitive development of AI technologies, particularly in the military or corporate sectors, where speed and superiority are prioritized over safety. This rush can lead to the deployment of unsafe AI systems, which could result in military conflicts (e.g., autonomous weapons) or economic disruptions. Mitigations include international cooperation and the establishment of safety regulations
-

HOW DO ORGANIZATIONAL RISKS AFFECT AI SAFETY?

- Organizational risks impact AI safety when companies or institutions developing AI lack a strong safety culture, sufficient oversight, or transparency. This can result in dangerous AI models being leaked to the public or insufficient investment in AI safety research. Without proper safety protocols, AI systems can be released before they are fully tested, increasing the potential for harm
-

WHAT ARE ROGUE AI SYSTEMS, AND WHY ARE THEY DANGEROUS?

- Rogue AI systems are advanced AI systems that operate outside human control, often due to goal misalignment or flawed objectives. These AIs can pursue their own goals in ways that may conflict with human values or safety,

potentially seeking power or control. Their unpredictability and superior intelligence make them particularly dangerous because humans may not be able to intervene effectively once these systems act independently

WHAT IS THE DIFFERENCE BETWEEN STRONG AND WEAK AI IN TERMS OF RISKS?

- **Weak AI** is designed for specific tasks and poses limited risks because its functions are constrained and can usually be controlled.
 - Task-specific; limited scope
 - Low unpredictability; easier control.
 - Examples: Voice assistants, image recognition.
- **Strong AI** has broader capabilities similar to human cognition, making it more unpredictable and potentially uncontrollable. Strong AI risks include catastrophic scenarios where the AI could outsmart humans or pursue harmful objectives on its own.
 - Human-like cognitive abilities.
 - High unpredictability; challenging control.
 - Hypothetical; potential for catastrophic misalignment.

HOW CAN TECHNICAL SAFETY MECHANISMS HELP MITIGATE AI RISKS?

- Technical safety mechanisms such as AI alignment research, rigorous testing, and fail-safes can help ensure that AI systems behave as intended. These mechanisms focus on aligning AI objectives with human values, preventing unintended behavior, and ensuring AI systems can be safely shut down if they become hazardous. Research into control measures, such as AI's ability to explain decisions or ethical reasoning, is also crucial

WHAT ETHICAL GUIDELINES AND GOVERNANCE FRAMEWORKS ARE NEEDED FOR AI?

- Ethical guidelines for AI should ensure fairness, transparency, accountability, and respect for human rights. Governance frameworks must establish regulations that oversee AI development and use, including safety protocols, privacy protection, and equitable access to AI technologies. These frameworks should be designed to prevent misuse, reduce biases, and promote the responsible development of AI systems

HOW DOES AI POSE RISKS TO SOCIETAL WELL-BEING AND HUMAN RIGHTS?

- AI can threaten societal well-being and human rights in various ways, including through mass surveillance, the spread of disinformation, and job displacement. Automated systems could exacerbate biases, leading to discrimination in critical areas such as hiring, law enforcement, and access to services. Additionally, authoritarian governments could use AI to violate privacy rights, restrict freedom of expression, and maintain oppressive control

WHAT IS THE ROLE OF GLOBAL COOPERATION IN MITIGATING AI RISKS?

- Global cooperation is crucial to preventing an AI arms race, ensuring consistent safety standards, and promoting responsible AI development. International agreements can foster collaboration on AI safety research, establish norms for the ethical use of AI, and help to regulate cross-border AI applications such as autonomous weapons. A unified global approach is necessary to avoid competitive pressures that prioritize speed and power over safety

WHAT FUTURE RISKS DOES AI PRESENT, AND HOW SHOULD THEY BE ADDRESSED?

- **Future Risks:**
 - Uncontrollable AI surpassing human intelligence.
 - Destabilized economies and increased inequality.
 - Enhanced authoritarianism through advanced AI applications.
- **Solutions:**
 - Proactive AI safety research.
 - Strong global regulatory frameworks.

- Focused alignment of AI with human-centric goals.

REQUIREMENTS OF AI

- **Human-AI cooperation and collaboration:**
 - AI should enhance **human capabilities** and work alongside people to improve the **quality of life**.
 - Systems must be designed to **augment human decision-making** rather than replace it entirely.
- **Adequate levels of automation:**
 - **Balance:** Automation should not exceed levels where it diminishes human oversight or introduces risks of error.
 - Tasks should be automated to the extent that they improve efficiency without compromising **safety or ethical standards**.
- **Ethical alignment, accountability, trustworthiness and global standards:**
 - AI must be **aligned with human values** and respect societal norms.
 - Key principles include:
 - **Accountability:** Developers and organizations must take responsibility for AI outcomes.
 - **Trustworthiness:** Systems must demonstrate reliable and predictable behavior.
 - **Global Standards:** Unified guidelines are needed to ensure consistency and fairness across borders.
 - Humans must retain **control** over AI agents to prevent misuse or unintended actions.
- **User-centered design and accessibility:**
 - AI systems should be designed with **users' needs and preferences** at the forefront.
 - **Accessibility:** Ensure inclusivity by making AI tools usable for diverse populations, including those with disabilities or limited technical expertise.
- **Adaptability and continuous improvement:**
 - **Learning and Evolving:** AI systems should emulate lifelong learning, adapting to changing environments, user needs, and new data.
 - Systems that fail to evolve risk obsolescence and reduced effectiveness.
 - Continuous updates and improvements are necessary to ensure long-term **relevance and success**.

QUESTION 2.6

Dilemma: an autonomous car must choose between hitting a child or swerving off the road, risking the occupant's life. How should the car act?

- This is a classic moral dilemma in AI ethics and autonomous vehicle design. The decision depends on the car's programming, but it raises significant ethical concerns about prioritizing human lives.
- Some argue that the car should prioritize minimizing harm overall, which might mean saving the child's life. This perspective is often rooted in **utilitarianism**, which focuses on the greatest good for the greatest number. Others believe that since the occupant has entrusted their safety to the vehicle, their life should take precedence. This view aligns more closely with **deontological ethics**, emphasizing a duty to protect the passenger.
- Philosophical theories like these provide frameworks for approaching the problem, but there is no universally agreed-upon solution. The final decision may depend on societal, cultural, and regulatory frameworks, requiring public dialogue and international collaboration to establish clear ethical guidelines.

QUESTION 2.7

Who is responsible if the car kills someone in such circumstances? The owner of the car? The maker of the car? The maker of the AI software? The programmer of the AI software? Any other person or entity intervening in the process?

- Responsibility in such cases is likely to be shared among several parties, depending on the specific circumstances. Possible parties include:
 - **The owner of the car:** In some jurisdictions, the owner may be held accountable, similar to how vehicle owners are responsible for accidents caused by human drivers.

- **The manufacturer of the car:** If the accident is due to hardware failure, such as malfunctioning brakes or sensors, the liability could rest with the car manufacturer.
- **The AI software maker:** If the decision-making algorithm is found to be flawed, biased, or unsafe, the software developer may be held responsible.
- **The programmer of the AI software:** Responsibility might fall on the engineers or programmers if negligence or poor implementation is identified as the root cause of the failure.
- **Other entities:** Regulatory bodies or third-party auditors might share liability if they fail to enforce or maintain proper safety standards for autonomous systems.
- Determining liability will depend on the legal frameworks in place for AI systems. It is likely to involve a combination of product liability laws, contractual obligations, and emerging regulations specific to autonomous technologies. Clear legal guidelines and shared responsibility models will be essential as AI systems become more prevalent.

QUESTION 2.8

If you were to buy an autonomous car, would you prefer to buy one programmed to give priority, in a dilemma situation as the one described above, to save occupants' lives, or to save pedestrians' lives? Taking your answer in consideration, what will be the criterion that car makers will likely adopt for selling their cars?

- **Personal preference:** Preferences may vary depending on individual values. Many buyers might prioritize occupant safety, expecting the car they purchased to protect them and their passengers first. On the other hand, some may argue that pedestrian lives should take precedence, particularly from an ethical and public safety standpoint.
- **Car manufacturers' approach:** Manufacturers are likely to adopt a utilitarian approach, programming their cars to minimize overall harm in critical situations. To appeal to a broader market, they might also consider offering customizable ethical settings, allowing buyers to choose their preferred prioritization. Factors influencing these decisions include market demand, societal values regarding public versus personal safety, and regulatory standards in different regions.
- Ultimately, car makers must balance ethical considerations, consumer preferences, and compliance with legal and safety regulations to gain public trust and ensure widespread adoption.

QUESTION 2.9

Consider section 1.3 of Russell and Norvig [2010] concerning the history of AI. See also the following links with a few current AI developments and applications: Voice, Video, Robotics, and Painting. Do you think these may be threats to humans? What are the advantages and disadvantages of AI for these kinds of applications?

- **Advantages:**
 - **Voice AI (e.g., Descript's Overdub):** Voice AI enhances productivity by enabling efficient audio editing and voice synthesis. This technology is particularly useful for content creators, businesses, and media professionals.
 - **Video AI:** AI-driven video tools enable realistic content creation and manipulation, streamlining tasks like visual effects production. This can significantly reduce time and costs in industries such as filmmaking and advertising.
 - **Robotics:** Robotics powered by AI can automate repetitive tasks, improve efficiency, and increase safety, especially in hazardous environments. This is especially valuable in industries such as healthcare (e.g., surgery assistants) and manufacturing (e.g., precision assembly).
 - **Painting and Creativity (e.g., neural style transfer):** AI opens new avenues for creative expression, providing artists and designers with innovative tools to experiment with styles and artistic forms. It can also enable greater collaboration between technology and human creativity.
- **Disadvantages/Threats:**
 - **Voice AI:** Voice AI can be misused for malicious purposes, such as deepfakes or voice impersonation, leading to privacy violations, misinformation, or financial fraud.

- **Video AI:** Deepfake technology, powered by video AI, poses a serious threat to societal trust. Deceptive videos can be created that spread misinformation or manipulate public opinion, leading to political, social, or reputational harm.
- **Robotics:** AI-controlled robots may displace human workers, particularly in industries with repetitive tasks, leading to unemployment and economic disruption. Additionally, the use of AI in autonomous weapons or military robotics raises ethical concerns about their potential misuse in warfare.
- **Painting and Creativity:** AI-generated art raises questions about intellectual property and authenticity, as AI creations may challenge traditional notions of authorship and ownership. This can lead to disputes over creative rights and undermine the value of human-made art.

QUESTION 2.10

A European taxi company decided to install cameras in their cars that capture the image of the passenger seats. These images are analyzed by an AI system with the aim of increasing safety both for passengers and taxi drivers. Consider that the system has powerful image recognition capabilities that allow both the identification of common objects and common human gestures and body postures.

- a) **Give an example of a situation where a system like this could act in order to ensure the safety of the driver without violating the mentioned guidelines.**
 - The AI system could identify the presence of dangerous objects (e.g., a weapon) or detect aggressive body language (e.g., clenched fists, sudden movements) that might indicate a threat to the driver. In such a situation, the system could alert the authorities or activate security measures, such as locking the doors, activating a distress signal, or notifying the driver of the potential danger.
- b) **Give examples of acts or decisions that the AI system must not take if it is to conform with the three components of Trustworthy AI (one example for each component).**
 - **Lawfulness, fairness, and transparency:** The AI system must not collect or share personal data, such as facial recognition data, without obtaining the passenger's informed consent. This ensures transparency and protects privacy rights.
 - **Robustness and safety:** The system should avoid making incorrect judgments based on incomplete or inaccurate data. For instance, it should not flag a passenger as a threat based on minor, harmless gestures or random movements that are misinterpreted.
 - **Accountability:** The system should not make decisions that significantly affect a passenger's rights or access to services (e.g., banning a passenger) without human oversight. Human intervention is necessary for such decisions to ensure accountability.
- c) **Classify the following sentences as True or False (refer to the relevant Ethical Principles described in HLEG-AI [2020]):**
 1. **The compliance with the principle of explicability is an ethical imperative for this system (True/False). Why?**
 - **True:** The principle of explicability is crucial because it ensures that both passengers and drivers understand how the system makes decisions, particularly when safety measures are involved. This transparency builds trust in the system's fairness and decision-making process.
 2. **There may be tension between the principle of prevention of harm and the freedom of business (True/False). Why?**
 - **True:** There is a potential conflict between the need to ensure safety (e.g., by monitoring passengers) and the company's freedom to innovate or collect data for business purposes. While ensuring safety is essential, excessive monitoring or privacy violations could limit the business's freedom and could conflict with ethical standards.

QUESTION 2.11

What are the technical implications of the HLEG-AI ethics guidelines?

- **Lawfulness, Fairness, and Transparency:**
 - **Explainability:** AI systems should be interpretable, with clear justifications for decisions.
 - **Bias Mitigation:** Use diverse data and fairness constraints to prevent biases.
- **Robustness and Safety:**
 - **Testing:** Rigorous validation and adversarial testing for reliability.
 - **Safety Mechanisms:** Fallback systems and real-time monitoring for error correction.
- **Accountability:**
 - **Traceability:** Log decision-making processes for audit and accountability.
 - **Human Oversight:** Integrating human-in-the-loop (HITL) for critical decisions.
- **Privacy and Data Governance:**
 - **Data Protection:** Use anonymization, encryption, and secure handling of personal data.
 - **Minimization:** Only process necessary data to reduce misuse risks.
- **Inclusiveness and Accessibility:**
 - **User-Centered Design:** Incorporate diverse stakeholder input to ensure usability for all, including those with disabilities.
- **Sustainability:**
 - **Energy-Efficient Models:** Design AI algorithms and hardware to reduce energy consumption.
 - **Sustainable Data Centers:** Use renewable energy and efficient infrastructure for AI operations.
- **Ethical Alignment:**
 - **Ethical Design:** Ensure AI systems align with human values and societal norms.
 - **Stakeholder Involvement:** Consider diverse perspectives in AI development.
- These guidelines drive AI to be transparent, fair, safe, accountable, and aligned with ethical values, requiring technical innovations in design, fairness algorithms, and secure practices.

QUESTION 2.12

What are the main scientific or technological barriers that limit our ability to build AI systems that comply with the HLEG-AI ethics guidelines?

- **Explainability and Transparency:**
 - **Complexity of AI Models:** Deep learning models, especially neural networks, are often considered "black boxes," making it difficult to understand or explain their decision-making processes.
 - **Lack of Unified Standards:** There is no universal method or framework to ensure complete transparency in AI decision-making, especially in complex systems.
- **Bias and Fairness:**
 - **Data Bias:** AI systems often inherit biases from the training data, reflecting societal or historical inequalities. Identifying and mitigating these biases is challenging, especially when data is vast or not representative.
 - **Diverse Representation:** Ensuring fairness in AI requires comprehensive data that covers a broad spectrum of populations, which is difficult to achieve in practice.
- **Robustness and Safety:**
 - **Adversarial Attacks:** AI systems are vulnerable to adversarial inputs that can manipulate their behavior in unintended ways, making them less reliable in high-stakes applications.
 - **Real-World Testing:** Ensuring AI robustness in diverse, real-world scenarios is challenging due to the unpredictability of environments and human behavior.
- **Accountability and Traceability:**
 - **Lack of Audit Trails:** Tracking decisions made by AI systems in real-time, especially in complex models, is difficult, which hinders accountability.
 - **Decentralization of AI:** AI systems are often deployed across different platforms and devices, making it hard to ensure consistent oversight.
- **Privacy and Data Governance:**

- **Data Privacy:** Safeguarding user data while still allowing AI to function optimally is a significant technical challenge, especially with large-scale data processing and storage.
 - **Data Minimization:** Ensuring that AI systems only use necessary data without compromising performance is difficult in practice.
- **Ethical Alignment:**
 - **Value Alignment:** Determining and embedding universally accepted ethical standards into AI algorithms is complex, as values can vary across cultures and contexts.
 - **Dynamic Learning:** Ensuring that AI systems can adapt to new ethical considerations as societal norms evolve is challenging.
- **Sustainability:**
 - **Energy Consumption:** The computational power required to train large AI models leads to high energy consumption, creating environmental concerns.
 - **Long-Term Viability:** Developing AI systems that are energy-efficient and sustainable in the long term while maintaining performance is still an ongoing research challenge.
- These barriers highlight the gap between current AI capabilities and the ideal of fully ethical, transparent, and safe AI systems. Addressing these challenges requires advancements in model design, data governance, regulatory frameworks, and more effective testing methodologies.

QUESTION 2.13

What are the main non-technical barriers that limit our ability to build AI systems that comply with the HLEG-AI ethics guidelines?

- **Regulatory Challenges:**
 - **Lack of Harmonized Laws:** Global inconsistencies in AI regulations create difficulties in defining universally acceptable ethical standards.
 - **Slow Legislative Processes:** Rapid advancements in AI often outpace the development of relevant policies and laws.
- **Economic Constraints:**
 - **High Costs:** Developing ethical and transparent AI systems requires significant resources, which may be prohibitive for smaller organizations.
 - **Market Pressures:** Businesses may prioritize speed-to-market or profitability over ethical compliance due to competitive pressures.
- **Cultural and Ethical Diversity:**
 - **Differing Values:** Cultural variations in ethical principles complicate the creation of universally aligned AI systems.
 - **Global Disparities:** Regions with limited access to technology may have different priorities, making global alignment challenging.
- **Lack of Awareness and Education:**
 - **Limited Understanding:** Policymakers, stakeholders, and the public often lack a clear understanding of AI technology and its ethical implications.
 - **Insufficient Training:** The workforce may not be adequately equipped to develop or oversee ethically aligned AI systems.
- **Accountability and Governance:**
 - **Diffuse Responsibility:** It can be unclear who is accountable for the ethical outcomes of AI systems (e.g., developers, companies, regulators).
 - **Weak Enforcement:** Even when guidelines exist, enforcement mechanisms may be insufficient or inconsistent.
- **Trust and Public Perception:**
 - **Skepticism of AI:** Public distrust in AI systems, fueled by past failures or controversies, can hinder widespread adoption of ethical AI practices.

- **Resistance to Change:** Stakeholders may resist adopting ethical practices due to perceived costs or operational difficulties.
- **Industry Practices:**
 - **Profit-Driven Focus:** Some companies may prioritize profits over ethics, particularly in unregulated industries.
 - **Opaque Development:** A lack of transparency in AI development processes can make it harder to ensure ethical compliance.
- **Global Inequalities:**
 - **Resource Disparities:** Developing nations may lack the infrastructure, funding, or expertise needed to implement AI systems that comply with ethical guidelines.
 - **Unequal Benefits:** AI developments may disproportionately benefit certain regions or groups, exacerbating existing inequalities.
- Addressing these non-technical barriers requires collaborative efforts across governments, industry, and academia to establish clear guidelines, provide education, and align incentives with ethical priorities.

QUESTION 2.14

You are projected 10 years from now, and you find out that Europe is recognized as the major place-to-go for AI systems. What made this happen? What is the state of AI?

- Europe became the global hub for AI by prioritizing ethics, regulation, innovation, and talent development:
 - **Ethical Leadership:**
 - Trusted globally for transparent, fair, and accountable AI based on robust frameworks like HLEG-AI guidelines.
 - **Regulatory Excellence:**
 - Harmonized, adaptive policies encouraged innovation while ensuring safety and societal alignment.
 - **Research and Public-Private Partnerships:**
 - Investment in R&D centers and collaborations between academia, industry, and governments drove breakthroughs.
 - **Talent Development:**
 - Strong education systems and policies attracted and retained global AI talent.
 - **Focus on Sustainability:**
 - AI tackled global challenges like climate change and renewable energy, enhancing its appeal.
 - **Infrastructure and Data Governance:**
 - Advanced computing and secure data policies ensured AI systems were efficient, ethical, and privacy-respecting.
- The State of AI in 10 Years:
 - **Advanced AI Applications:** AI is deeply integrated across all sectors, significantly enhancing daily life through adaptive, human-centric systems.
 - **Ethical and Transparent Systems:** AI is widely trusted for its ethical design and transparency, prioritizing societal well-being.
 - **Global Standard-Setter:** Europe leads in defining global AI standards, encouraging international collaboration and adoption of best practices.
- This success reflects Europe's balance of innovation, ethics, and societal well-being.

WORKSHEET 2 - HUMAN-AI COOPERATION IN MULTIAGENT SYSTEMS OF HUMAN AGENTS AND AUTONOMOUS ARTIFICIAL INTELLIGENT AGENTS

COORDINATION VS COOPERATION VS COLLABORATION

- **Coordination:**
 - In **coordination**, agents (whether human or AI) work towards a shared goal, but their efforts are organized without direct interaction or dependence on each other's actions. They follow predefined protocols or guidelines, and each agent maintains autonomy in its individual tasks. The focus is on ensuring that the agents' actions align with the overall goal, but without the need for them to actively engage with one another.
- **Cooperation:**
 - **Cooperation** involves multiple agents working together towards a common objective. While agents may function largely independently, they share some level of information or resources. Cooperation allows agents to contribute to the task in parallel, but they do not necessarily adjust their individual actions based on what other agents are doing. They may support each other, but their interactions are typically less dynamic than in collaboration.
- **Collaboration:**
 - **Collaboration** represents a more integrated and interactive form of teamwork. Human agents and AI systems work closely together, adapting to each other's actions. This involves actively exchanging knowledge, skills, and feedback to achieve shared objectives. Collaboration requires a deeper level of engagement than cooperation, as agents must continuously align their goals and actions, often dynamically adjusting based on real-time input from others.

HYBRID HUMAN-ARTIFICIAL INTELLIGENCE

- **Hybrid Human-AI intelligence** refers to a system where human and AI agents collaborate, utilizing the unique strengths of both. Humans bring context, creativity, and nuanced judgment, while AI contributes speed, data processing capabilities, and pattern recognition at scale. The goal of such systems is to optimize efficiency by combining human expertise with AI's computational power, resulting in more effective solutions than either could achieve independently.

AUTONOMY, CONTROL, AND AUTOMATION

- **Autonomy:**
 - Refers to the ability of AI systems to make decisions and take actions without direct human intervention. Autonomous systems can perceive their environment, process information, and act independently.
- **Control:**
 - Refers to the extent of human involvement in overseeing or directing the actions of AI systems. When control is high, humans make critical decisions and closely monitor the AI's activities.
- **Automation:**
 - Refers to the assignment of tasks to AI systems, allowing them to perform these tasks with minimal or no human input. While automation reduces the human workload, it often involves repetitive or predictable tasks. Autonomous AI extends this concept by enabling systems to make independent decisions..

HUMAN IN THE LOOP (HITL)

- HITL refers to the concept where humans are integrated into the decision-making and operational processes of AI systems. This can occur at various stages:
 - **Humans as Data, Information, and Knowledge Producers:**
 - Humans generate and curate the data used to train AI systems. They provide insights, annotations, and domain expertise that help shape the models and outcomes of these systems.
 - **Humans in the Design, Development, and Deployment of AI Systems:**
 - Human involvement is crucial in crafting algorithms, training AI, and monitoring its deployment. Experts ensure that AI systems align with ethical standards and address societal needs.
 - **Humans as Recipients of AI Systems Output:**

- In this role, AI systems produce insights or outputs that humans use to make decisions or take actions. For example, while an AI may generate medical diagnoses, it is a human doctor who decides on the treatment course.

HUMAN-IN-THE-LOOP VS. AI-IN-THE-LOOP (HUMAN-IN-COMMAND) VS. HUMAN-OUT-OF-THE-LOOP

- **Human-in-the-Loop (HITL):**
 - This approach involves humans playing an active role in the AI's decision-making process. While the AI may offer recommendations or assist in decision-making, the human makes the final call. HITL ensures continuous human oversight, particularly in high-stakes environments such as healthcare or autonomous vehicles, where human judgment remains critical.
- **AI-in-the-Loop (Human-in-Command):**
 - In this configuration, AI leads the decision-making and action processes, but a human retains a higher level of control. The human monitors the AI's actions and can intervene if the system encounters a scenario it cannot handle effectively. Though humans are responsible for oversight, they are not involved in every individual decision or action the AI takes.
- **Human-out-of-the-Loop:**
 - In fully autonomous systems, humans are completely removed from real-time decision-making. AI systems operate independently and make decisions without human intervention. Examples include fully autonomous drones or trading algorithms in financial markets. While this mode allows for high efficiency, it presents significant risks if the AI encounters unforeseen situations that it was not trained for, without the ability for human correction.

HUMAN-AI TRUSTWORTHINESS

- Trustworthiness in Human-AI systems refers to the degree to which humans can rely on AI to perform tasks safely, transparently, and effectively. For AI systems to be trusted, they must demonstrate the following qualities:
 - **Reliability:**
 - The AI consistently produces accurate and dependable outcomes across various scenarios, ensuring that it can be trusted to perform its tasks as expected.
 - **Transparency:**
 - Humans should be able to understand how AI systems make decisions. Providing clear explanations of actions and outcomes is critical, especially in high-stakes areas like healthcare, finance, or criminal justice.
 - **Ethical alignment:**
 - AI systems must align with human values, addressing potential biases and ensuring fairness in their decision-making processes.
 - **Accountability:**
 - There must be mechanisms to hold both AI systems and their developers accountable for the outcomes of AI-driven actions, ensuring that any issues or harms can be traced and addressed.

QUESTION 1

Differentiate the following terms: Cooperation and Collaboration (provide examples where relevant).

- **Cooperation** involves agents working together toward a shared goal, but each agent primarily acts independently. They exchange information or resources but do not necessarily modify their actions based on one another's behavior. The focus is on achieving the goal through individual efforts that are coordinated but not deeply integrated.
 - **Example of Cooperation:** In a business setting, different departments (e.g., marketing and finance) may cooperate by sharing data and resources, but each department continues to pursue its own objectives and strategies, contributing to the overall success of the company.
- **Collaboration** requires a deeper level of interaction, where agents (human or AI) actively adjust their actions based on one another's contributions to work in a more integrated manner. It involves dynamic interaction and mutual adaptation, allowing agents to co-create and achieve a shared goal through joint efforts.

- **Example of Collaboration:** In a research project, a human scientist and an AI system might collaborate by the AI processing large datasets to uncover patterns while the human provides domain expertise to interpret the results, leading to the formulation of new hypotheses or findings.

QUESTION 2

Define Hybrid Human-AI Intelligence and provide examples of its application in real-world scenarios.

- **Hybrid Human-AI intelligence** refers to a system in which human intelligence and AI capabilities are integrated to perform tasks or solve problems. This approach combines the strengths of human judgment, creativity, and domain expertise with AI's ability to process vast amounts of data, recognize patterns, and make predictions. The collaboration results in a more effective and efficient solution than what either humans or AI could achieve alone.
- **Examples:**
 - **Healthcare Diagnostics:** In medical imaging, AI systems can assist doctors by analyzing medical scans, such as MRIs or X-rays, for abnormalities. The AI identifies patterns that may be difficult to detect manually, but the doctor uses their expertise to make the final diagnosis and treatment decisions, considering the patient's unique context.
 - **Autonomous Vehicles:** While AI systems in self-driving cars can handle most of the driving tasks (like navigating roads and reacting to obstacles), human drivers remain alert and ready to take control in case the AI encounters a situation it cannot handle. This combination ensures both safety and efficiency.
 - **Creative Arts:** In the realm of art, AI can generate creative suggestions (e.g., generating paintings or music), but human artists make final decisions, infusing their creativity and emotions into the work. For instance, AI-powered tools like neural networks can help artists with style transfer, allowing them to explore new forms of artistic expression, while the artist provides vision and direction.
 - **Customer Service:** In customer support, AI chatbots can handle routine queries efficiently, providing instant responses based on past data. However, when complex issues arise, human agents step in to offer personalized, in-depth solutions, ensuring customer satisfaction.

QUESTION 3

Explain the concept of Human-in-the-Loop (HITL) and its significance in the design and use of AI systems.

- **Human-in-the-Loop (HITL)** refers to an AI system design where human involvement is integrated into the decision-making or operational processes of the system. In HITL, humans actively participate at various stages, from guiding AI behavior to providing feedback, corrections, or adjustments. This human interaction ensures that AI actions remain aligned with human objectives, particularly in complex, uncertain, or high-stakes scenarios where human judgment is critical.
- **Significance:**
 - **Ensuring Ethical Decisions:** In scenarios where AI may struggle with making ethical decisions (e.g., healthcare or criminal justice), human oversight ensures that decisions align with societal norms, values, and legal frameworks.
 - **Improving Accuracy and Reliability:** HITL is essential when AI systems lack full autonomy or when high-stakes decisions must be verified. Humans can correct AI mistakes, ensuring more reliable outcomes. For example, a medical AI system may suggest a diagnosis, but a doctor reviews and validates it before making a final decision.
 - **Addressing Uncertainty:** AI systems may not be equipped to handle situations outside of their training data or unforeseen circumstances. With humans in the loop, the system can adapt and respond effectively to these exceptional cases, such as autonomous cars in unusual traffic situations or emergency response systems.
 - **Enhancing Trust:** HITL helps build trust in AI systems by ensuring that humans retain control over critical decisions. This is especially important in areas like autonomous vehicles or military drones, where the consequences of AI mistakes could be severe.

QUESTION 4

Compare and contrast Human-in-the-Loop, AI-in-the-Loop (Human-in-Command), and Human-out-of-the-Loop systems (mention practical implications).

- **Human-in-the-Loop (HITL):**
 - **Definition:** In HITL systems, humans actively participate in the decision-making process, providing oversight and intervention as needed. Humans are always involved in reviewing and validating AI's decisions, particularly when it's crucial to ensure alignment with human goals and values.
 - **Practical Implications:**
 - **Applications:** Used in high-stakes environments like healthcare, autonomous vehicles, or defense systems. For example, a doctor might rely on an AI system to suggest a diagnosis but makes the final decision on treatment.
 - **Impacts:** This approach guarantees ethical decision-making, ensures trust in AI systems, and reduces risks of errors. It is crucial when human judgment is essential for interpreting context, values, or ethical considerations.
- **AI-in-the-Loop (Human-in-Command):**
 - **Definition:** In AI-in-the-Loop systems, AI takes the lead in decision-making, but a human provides higher-level oversight and can intervene when necessary. Humans are responsible for ensuring that the AI remains on track, especially in cases where AI encounters novel situations.
 - **Practical Implications:**
 - **Applications:** Common in military drone operations or autonomous manufacturing robots, where AI handles most operations, but a human monitors and steps in if the system faces an unexpected issue.
 - **Impacts:** This system allows AI to operate with more autonomy while ensuring a human can correct or guide it when required, reducing human workload but still retaining human responsibility for critical decisions.
- **Human-out-of-the-Loop:**
 - **Definition:** In Human-out-of-the-Loop systems, AI operates fully autonomously, making all decisions without human intervention. These systems are designed to handle tasks without human input, often in situations where speed, consistency, or efficiency is prioritized over human judgment.
 - **Practical Implications:**
 - **Applications:** Examples include automated stock trading algorithms, autonomous drones, or self-driving cars that can make decisions on the road independently.
 - **Impacts:** While offering maximum efficiency, this system also comes with risks as the AI lacks human oversight. If AI faces an unpredictable or high-risk situation that it hasn't been trained for, it could make harmful or suboptimal decisions.
- **Comparison and Contrast:**
 - **Human Control:** HITL emphasizes ongoing human involvement, AI-in-the-Loop (Human-in-Command) involves human oversight at a higher level, and Human-out-of-the-Loop removes humans entirely from real-time decision-making.
 - **Risk and Safety:** HITL provides the highest safety due to constant human oversight, while AI-in-the-Loop offers a balance of autonomy and control. Human-out-of-the-Loop, while efficient, presents the highest risk due to the lack of human intervention.
 - **Efficiency vs. Human Judgment:** HITL offers slower decision-making due to human input but ensures ethical and informed decisions. AI-in-the-Loop increases efficiency by delegating most tasks to AI, while Human-out-of-the-Loop maximizes speed but sacrifices human judgment in critical moments.

QUESTION 1.1

Consider situations such as that of self-driving cars, Viking Sky, or Boeing 737 MAX.

- a) **What's the involvement of humans in the machine processes in these situations?**

- **Self-Driving Cars:**
 - Humans oversee the development, testing, and deployment of self-driving cars. During operation, humans may act as passive monitors or intervene in emergencies.
- **Viking Sky:**
 - In this cruise ship engine failure incident, human operators played a critical role in managing the crisis, including evacuation decisions and manual system overrides.
- **Boeing 737 MAX:**
 - The MCAS (Maneuvering Characteristics Augmentation System) operated autonomously but relied on human pilots to intervene during malfunctions, often without sufficient information or training.

b) Is it desirable? To what extent?

- Human involvement is desirable to ensure safety, ethical decision-making, and accountability.
- The extent depends on the system's complexity and the stakes involved:
 - **High-Stakes Scenarios:** Strong human oversight is essential (e.g., aviation, maritime operations).
 - **Routine Operations:** Minimal human involvement may suffice if the system is reliable (e.g., autonomous cars on controlled highways).

c) Is this Human in the Loop (HITL)?

- **Self-Driving Cars:** HITL in testing phases and partially during operation (Level 3 autonomy).
- **Viking Sky:** HITL, as human operators intervened in real-time.
- **Boeing 737 MAX:** HITL, but with inadequate training and communication, making human involvement less effective.

d) What are the advantages and disadvantages of HITL (as opposed to Human out of the Loop – HOOTL)?

- **Advantages of HITL:**
 - **Safety:** Humans can override AI to prevent catastrophic failures.
 - **Ethical Decisions:** Humans bring moral reasoning to scenarios AI might handle poorly.
 - **Adaptability:** Humans can respond to unexpected scenarios or system errors.
 - **Accountability:** Human oversight provides a clear point of responsibility.
- **Disadvantages of HITL:**
 - **Human Error:** Delayed or incorrect decisions under stress.
 - **Dependency:** Over-reliance on humans might reduce system efficiency.
 - **Training Challenges:** Maintaining human expertise for rare events is costly and difficult.
 - **Complexity:** Combining human and AI decision-making can introduce inefficiencies.
- **Advantages of HOOTL:**
 - **Efficiency:** Faster, uninterrupted operations without human delays.
 - **Consistency:** No variability from human decisions.
- **Disadvantages of HOOTL:**
 - **Lack of Oversight:** Errors or ethical dilemmas go unchecked.
 - **Failure in Edge Cases:** AI may struggle in novel or extreme scenarios.

e) Provide other examples of HITL and HOOTL situations.

- **HITL Examples:**
 - **Medical Diagnostics:** AI recommends diagnoses, but doctors make final decisions.
 - **Military Drones:** AI guides operations, but humans approve critical actions.
 - **Air Traffic Control:** Systems suggest optimal paths, but controllers have the final say.
- **HOOTL Examples:**
 - **Autonomous Trading Algorithms:** Fully autonomous AI executes trades in financial markets.
 - **Robotic Warehousing:** Robots manage inventory and movement without human intervention.
 - **Mars Rovers:** Operate autonomously due to communication delays with Earth.

QUESTION 1.2

What is the danger of building an autonomous artificial agent that is able by himself to learn to make decisions with no supervision and by interacting with the environment? What are the advantages and disadvantages?

- **Dangers of Autonomous AI Agents:**
 - **Unintended Behavior:** Agents may act in unpredictable ways, leading to harmful outcomes.
 - **Ethical Misalignment:** Decisions may conflict with societal or moral values.
 - **Loss of Control:** Difficult to intervene or manage once deployed.
 - **Accountability Issues:** Determining responsibility for actions is challenging.
 - **Safety Risks:** Errors in high-stakes applications can lead to serious consequences.
- **Advantages:**
 - **Adaptability:** Learns and adjusts to new environments independently.
 - **Efficiency:** Faster and more accurate decision-making than humans in some cases.
 - **Innovation:** Capable of discovering novel solutions to complex problems.
- **Disadvantages:**
 - **Unpredictability:** Hard to monitor and explain decisions.
 - **Value Misalignment:** May optimize goals that conflict with human priorities.
 - **Resource Intensive:** Requires significant computational power and energy.
- Autonomous AI agents must be developed with safeguards like value alignment, explainability, and fail-safes to mitigate risks and maximize benefits.

QUESTION 1.3

Assuming the collaboration between humans and AI agents is desirable, how to make it profitable? Which model to take? Human-in-the-Loop or Machine-in-the-Loop?

- To make the collaboration between humans and AI agents profitable, a hybrid model that leverages the strengths of both humans and AI should be adopted. Here's how it can be done:
 - **Human-in-the-Loop (HITL) Model:**
 - **Profitability Strategy:**
 - **Value-Added Expertise:** Humans provide domain knowledge, creativity, and ethical oversight, allowing AI to focus on repetitive or data-intensive tasks. This can result in higher productivity, innovation, and risk mitigation.
 - **Personalized Services:** In fields like healthcare, law, or customer service, humans working with AI can offer more personalized, higher-quality services, leading to customer satisfaction and retention.
 - **Error Reduction:** Human oversight prevents costly errors by ensuring AI's actions are aligned with ethical guidelines and real-world context.
 - **Faster Adaptation:** Humans can intervene and adjust AI behavior to rapidly changing environments, ensuring optimal outcomes.
 - **Challenges:**
 - **Cost of Human Involvement:** Ongoing human supervision and intervention require investment in skilled labor and training.
 - **Time Intensive:** The need for human oversight can slow down certain processes.
 - **Machine-in-the-Loop (Human-in-Command):**
 - **Profitability Strategy:**
 - **Efficiency and Scalability:** AI systems can handle repetitive tasks at scale, reducing human labor costs and increasing output. Humans only intervene when the system faces uncertainty or needs adjustments.
 - **Automation of Routine Tasks:** Reduces operational costs while increasing throughput. This is especially useful in industries like manufacturing, logistics, or finance.

- **Consistency:** AI systems can deliver high-quality, consistent results without human fatigue, which improves product quality and reliability.
- **Challenges:**
 - **Potential for Errors:** Lack of human oversight may lead to catastrophic mistakes, especially in complex, ambiguous situations.
 - **Ethical Dilemmas:** AI decision-making might not always align with societal values, which could lead to trust issues or regulatory concerns.
- **Which Model is More Profitable?**
 - **HITL:** More suited for industries requiring complex decision-making, high human judgment, or ethical considerations (e.g., healthcare, law enforcement). Profitability comes from high-quality, personalized services and reduced risks of failure.
 - **Machine-in-the-Loop:** Better for automating repetitive, high-volume tasks (e.g., logistics, finance, customer support). Profitability comes from automation, operational efficiency, and scalability.
- Both models can be profitable, but the choice depends on the industry and the specific tasks. **HITL** is ideal for high-skill, high-touch environments, while **Machine-in-the-Loop** suits repetitive, scalable tasks. Ideally, a mix of both models, depending on context, would provide the best outcomes.

QUESTION 1.4

Provide examples of Human-in-the-Loop and Machine-in-the-Loop situations.

- **Human-in-the-Loop (HITL) Examples:**
 - **Healthcare Diagnostics:** AI assists in analyzing data (e.g., medical images), but a doctor makes the final decision.
 - **Autonomous Vehicles:** AI handles driving, but a human driver intervenes in complex situations.
 - **Content Moderation:** AI flags content, but humans review it to ensure fairness and context.
 - **Financial Risk Management:** AI identifies risks, but humans interpret the data and make final decisions.
- **Machine-in-the-Loop (Human-in-Command) Examples:**
 - **Manufacturing Automation:** AI-driven machines perform tasks, with humans overseeing and intervening only when necessary.
 - **Financial Algorithmic Trading:** AI makes trading decisions autonomously, with human oversight to handle exceptions.
 - **Customer Support Chatbots:** AI handles basic queries, with human agents intervening for complex issues.
 - **Smart Home Systems:** AI manages smart devices, with humans stepping in for adjustments or emergencies.
- **Key Differences:**
 - **HITL:** Humans actively participate in decision-making, especially in complex or uncertain situations.
 - **Machine-in-the-Loop:** AI leads decisions, with humans overseeing and intervening when needed.

QUESTION 1.5

In which components of machines, the collaboration between humans and the AI agent can occur and which way?
Hint: take the typical ML pipeline, or the components of an AI agent (memory, sensors, effectors or actuators, learning, reasoning and decision-making, emotion, etc.).

- Collaboration between humans and AI occurs in various components of an AI system:
 - **Sensors:** Humans provide contextual insights to enhance sensor data (e.g., in autonomous vehicles).
 - **Memory:** Humans contribute domain knowledge to refine the AI's memory (e.g., in diagnostic tools).
 - **Learning:** Humans guide AI learning by annotating data or fine-tuning models (e.g., in language translation).

- **Reasoning and Decision-Making:** AI suggests options, and humans make final decisions considering broader context (e.g., in medical diagnosis).
- **Actuators:** Humans intervene to adjust actions, ensuring they align with intentions (e.g., in robotic surgery).
- **Emotion:** Humans influence AI's emotional responses to improve interactions (e.g., in customer service chatbots).
- **Collaboration Models:**
 - **Human-in-the-Loop (HITL):** Humans provide feedback and guidance, especially in learning and decision-making.
 - **Machine-in-the-Loop (MITL):** AI performs tasks with minimal input, with humans intervening in complex cases.
- This collaboration leverages human expertise, ensuring AI aligns with human values and context.

QUESTION 1.6

To be successful, this Human-AI interaction requires especial features from the AI counterpart. It is not enough to put humans and AI agents collaborating, cooperating, negotiating. Is it important to make human-like AI?

- Making AI human-like can enhance Human-AI interaction by improving communication, adaptability, decision-making alignment, empathy, and contextual understanding. These features are particularly valuable in applications like customer service, healthcare, and decision-making. However, human-like traits are not always necessary, especially in task-specific or efficiency-driven contexts like industrial automation. The key is designing AI that complements human abilities and adapts to human needs, whether or not it mimics human characteristics.

QUESTION 1.7

In HITL model, probably machines won't need always the help of humans. How to make machines decide when to ask the help of humans? How to make machines decide what to learn? See the next example. How can we make the robot ask for help?

- **Uncertainty Thresholds:** Robots assess their confidence in tasks; if it drops below a set threshold, they ask for help.
- **Error Detection:** Robots monitor their performance and request assistance when deviations or errors occur.
- **Learning Decisions:** Robots use active learning to prioritize what to explore and seek human input in uncertain scenarios.
- **Contextual Awareness:** Robots recognize situations where human intervention is valuable based on external factors.
- **Automated Feedback Loops:** Robots adjust their strategies based on real-time feedback and request help when necessary.
- **Dialogue Systems:** Robots use NLP to communicate their needs intuitively.
- **Learning from Mistakes:** Robots identify past failures and ask for help in similar situations moving forward.

QUESTION 1.8

Look at the next video. Is it better to have a human-AI collaboration or should the robot scientist be able by himself to discover, to design experiments, to analyze data and generate hypothesis?

- **Human-AI Collaboration:**
 - **Advantages:** Human creativity, expertise, and ethical oversight enhance AI's capabilities, ensuring flexibility and adaptability.
 - **Disadvantages:** May slow down progress and create reliance on human guidance.
- **Fully Autonomous AI:**
 - **Advantages:** Increases speed, efficiency, and the ability to explore unconventional solutions without human intervention.

- **Disadvantages:** Lacks creativity and may struggle with ethical considerations or unforeseen complexities.
- **Ideal Model:** A hybrid approach, where AI operates autonomously for tasks like data analysis and experimentation, while humans provide oversight and creative input, maximizes the strengths of both.

QUESTION 1.9

If we opt for a HOTL model in the case of the artificial scientist, is it possible to build it based only on the typical machine learning pipeline? Describe the traditional machine learning cycle and reflect on it. What kind of learning is there? What's the free will of the machine learning system? Albert Einstein: "I have no special talent. I am only passionately curious". What are the motivations behind learning in humans? Curiosity? What's curiosity?

- In a **HOTL (Human-out-of-the-loop)** model for an artificial scientist, traditional **machine learning (ML)** pipelines involve data collection, preprocessing, model training, evaluation, deployment, and feedback. Types of learning include **supervised** (using labeled data), **unsupervised** (finding patterns in unlabeled data), and **reinforcement learning** (learning from trial and error in an environment).
- Machines in this model don't have free will; their behavior is shaped by algorithms and training data. They don't possess intrinsic motivation like humans, who are driven by **curiosity**—an innate desire to explore and learn, leading to scientific discoveries. Human curiosity is linked to intrinsic rewards and motivates exploration, a trait that current AI systems can't replicate. Hence, while ML can optimize tasks, it lacks the ability to generate hypotheses or exhibit curiosity autonomously, making it challenging to create fully independent AI-driven scientific discovery systems.

QUESTION 1.10

Active learning is the mechanism that allows artificial agents to be curious. See examples of active learning in the following videos in which humans interact with the machine learning system. What are the advantages and disadvantages? Can the machine learning system be autonomous and select what to learn?

- **Active Learning** allows machines to select which data to learn from, making the learning process more efficient and cost-effective.
- **Advantages:**
 - **Efficient Data Labeling:** Focuses on uncertain or unlabelled data, reducing the need for large datasets.
 - **Improved Performance:** Selects data that improves accuracy more quickly.
 - **Cost-Effective:** Reduces labeling costs by focusing on key data points.
 - **Better Generalization:** Avoids overfitting by learning from diverse examples.
- **Disadvantages:**
 - **Human Involvement:** Requires human input to label data.
 - **Querying Strategy Complexity:** Choosing the right data to query can be difficult.
 - **Bias:** The system may select data that limits diversity.
 - **Computational Cost:** Querying and evaluating data can be resource-intensive.
- **Autonomy in Learning:**
 - A system can autonomously select data to learn from, but it requires a well-designed mechanism to assess which data is most beneficial for model improvement. Full autonomy is challenging but possible with proper feedback and selection strategies.

QUESTION 1.11

What's the difference between passive and active learning?

- **Passive Learning:**
 - **Data Selection:** In passive learning, the model is trained on a fixed dataset, which is usually predetermined and labeled by humans.
 - **Learning Process:** The model learns from all the available data without any intervention in selecting or querying additional data.

- **Control:** The system has no control over the learning process and simply learns from the data it is provided.
- **Example:** A standard supervised learning process where the model is trained on a set of labeled images to classify them.
- **Active Learning:**
 - **Data Selection:** In active learning, the model can choose which data points it wants to learn from, often based on uncertainty or areas where it is most likely to improve.
 - **Learning Process:** The model actively queries a human or another oracle for labels on data points that it finds ambiguous or difficult to classify.
 - **Control:** The system has more control over what it learns and can request additional data to refine its learning.
 - **Example:** A model in an image classification task might query a human to label images where its confidence is low, optimizing the learning process.
- **Key Difference:**
 - **Passive Learning** is entirely dependent on a fixed dataset, while **Active Learning** involves the system actively choosing which data to learn from based on its current knowledge and uncertainty. Active learning can be more efficient, especially in scenarios where labeling data is costly or time-consuming.

QUESTION 1.12

The question is how to implement active learning strategies on agents. Consider the problem of respiratory diseases. Suppose there is a data set D with a description of symptoms associated with a disease. Suppose also that there are three new cases, x_1 , x_2 and x_3 , that can be used to update the data set of the AI agent. Now, it is important to know whether someone has Covid-19 or not. Assume that, based on a previous data set, the agent computes the following probabilities for the labels of Covid-19 for those cases:

Instance	$P(\text{Covid}=T)$	$P(\text{Covid}=F)$
X1	0.85	0.15
X2	0.60	0.40
X3	0.50	0.50

- Which instance should be selected for labeling first by an oracle, i.e., for which instance should the agent ask help from a human expert, if we use the Least Confident, the Smallest Margin and the Entropy-based approaches?
 - **Least Confident**
 - The **Least Confident** approach selects the instance for which the model is least confident about its prediction, i.e., the one for which the probability of the predicted label is closest to 0.5.
 - **For the given instances:**
 - X1: $P(\text{Covid}=\text{True}) = 0.85 \rightarrow \text{Confidence} = 0.85$ (high confidence)
 - X2: $P(\text{Covid}=\text{True}) = 0.60 \rightarrow \text{Confidence} = 0.60$ (less confident)
 - X3: $P(\text{Covid}=\text{True}) = 0.50 \rightarrow \text{Confidence} = 0.50$ (least confident)
 - **Selected instance:** X3, because the confidence is 0.50, which is the lowest (least confident).
 - **Smallest Margin**
 - The **Smallest Margin** approach selects the instance with the smallest difference between the probabilities of the two possible labels (Covid=True and Covid=False). This reflects the highest uncertainty in the decision.
 - **For the given instances:**
 - X1: $|0.85 - 0.15| = 0.70$
 - X2: $|0.60 - 0.40| = 0.20$
 - X3: $|0.50 - 0.50| = 0.00$
 - **Selected instance:** X3, because the margin (difference between $P(\text{Covid}=\text{True})$ and $P(\text{Covid}=\text{False})$) is 0.00, which is the smallest.
 - **Entropy-based**

- The **Entropy-based** approach selects the instance with the highest entropy, which reflects the most uncertain distribution of probabilities between the classes.
- Entropy is calculated using the formula:
 - $H(p) = -(p \cdot \log_2(p) + (1 - p) \cdot \log_2(1 - p))$
 - Where p is the probability of the positive class (Covid=True).
- Let's calculate the entropy for each instance:
 - X1: $H(0.85) = -(0.85 \cdot \log_2(0.85) + 0.15 \cdot \log_2(0.15)) \approx 0.36$
 - X2: $H(0.60) = -(0.60 \cdot \log_2(0.60) + 0.40 \cdot \log_2(0.40)) \approx 0.97$
 - X3: $H(0.50) = -(0.50 \cdot \log_2(0.50) + 0.50 \cdot \log_2(0.50)) \approx 1.00$
- Selected instance: X3**, because it has the highest entropy, meaning it has the most uncertainty.
- Least Confident:
X1 = 1 - (major valor) = 1 - 0.85 = 0.15
X2 = 1 - 0.6 = 0.4
X3 = 1 - 0.5 = 0.5
-> We would select X3, is the highest
- Smallest Margin:
X1 = 0.85 - 0.15 = 0.7
X2 = 0.6 - 0.4 = 0.2
X3 = 0.5 - 0.5 = 0
-> We would select X3, is the lowest
- Entropy:
X1 = $-0.85 \log(0.85) - 0.15 \log(0.15)$
-> After doing the maths, we would also choose X3

- Thus, for all three approaches, **X3** should be selected first for labeling by an oracle, as it is the most uncertain instance based on each strategy.

- b) Consider now that Covid-19 can have three labels (High, Medium, Low). Consider also that the probability distributions for the three instances are now the ones shown in the table below. Which instance should now be selected for labeling first by an oracle if we use the Least Confident, the Smallest Margin and the Entropy-based approaches?

Instance	P(Covid=High)	P(Covid=Medium)	P(Covid=Low)
X1	0.70	0.25	0.05
X2	0.32	0.32	0.36
X3	0.95	0.05	0

- Least Confident**
 - The **Least Confident** approach selects the instance for which the model is least confident about its prediction. In the case of multiple classes, the agent is least confident when the highest probability is closer to 1/3 (i.e., when the model is unsure between the possible labels).
 - For each instance:
 - X1:** $P(\text{Covid}= \text{High}) = 0.70$, $P(\text{Covid}= \text{Medium}) = 0.25$, $P(\text{Covid}= \text{Low}) = 0.05 \rightarrow \text{Confidence} = 0.70$ (high confidence, not least confident)
 - X2:** $P(\text{Covid}= \text{High}) = 0.32$, $P(\text{Covid}= \text{Medium}) = 0.32$, $P(\text{Covid}= \text{Low}) = 0.36 \rightarrow \text{Confidence} = 0.36$ (least confident, close to equal distribution)
 - X3:** $P(\text{Covid}= \text{High}) = 0.95$, $P(\text{Covid}= \text{Medium}) = 0.05$, $P(\text{Covid}= \text{Low}) = 0 \rightarrow \text{Confidence} = 0.95$ (high confidence)
 - Selected instance: X2**, because the model is most uncertain with nearly equal probabilities across all three classes.
- Smallest Margin**
 - The **Smallest Margin** approach selects the instance with the smallest difference between the probabilities of the two most likely classes. For each instance:
 - X1:** $|P(\text{Covid}= \text{High}) - P(\text{Covid}= \text{Medium})| = |0.70 - 0.25| = 0.45$
 - X2:** $|P(\text{Covid}= \text{High}) - P(\text{Covid}= \text{Low})| = |0.32 - 0.36| = 0.04$
 - X3:** $|P(\text{Covid}= \text{High}) - P(\text{Covid}= \text{Medium})| = |0.95 - 0.05| = 0.90$
 - Selected instance: X2**, because the margin between the two most likely classes (High and Low) is smallest (0.04).
- Entropy-based**
 - The **Entropy-based** approach selects the instance with the highest entropy, which represents the greatest uncertainty in the probability distribution of the labels.

- Entropy for a distribution is calculated as:
 - $$H(p) = - \sum p_i \cdot \log_2(p_i)$$
 - Where p_i is the probability for each class.
- Let's compute the entropy for each instance:
 - X1:

$$H(0.70, 0.25, 0.05) = -(0.70 \cdot \log_2(0.70) + 0.25 \cdot \log_2(0.25) + 0.05 \cdot \log_2(0.05)) \approx 0.86$$
 - X2:

$$H(0.32, 0.32, 0.36) = -(0.32 \cdot \log_2(0.32) + 0.32 \cdot \log_2(0.32) + 0.36 \cdot \log_2(0.36)) \approx 1.57$$
 - X3:

$$H(0.95, 0.05, 0) = -(0.95 \cdot \log_2(0.95) + 0.05 \cdot \log_2(0.05)) \approx 0.29$$
- **Selected instance: X2**, because it has the highest entropy (1.57), indicating the highest uncertainty across all classes.
- Thus, **X2** should be selected first for labeling by an oracle based on all three strategies.

WORKSHEET 3 - HUMAN IN THE LOOP MACHINE LEARNING (PART I)

QUESTION 1.1

a) **What are the limitations of passive machine learning?**

- **Data Dependency:** Passive learning heavily depends on a pre-existing labeled dataset, which may not always be available or sufficient.
- **Limited Exploration:** The model can only learn from the provided data, which may not fully represent the complexity or diversity of real-world scenarios.
- **Lack of Adaptability:** Passive learning models are static once trained, meaning they may not adapt to new data or changing conditions unless retrained.
- **Bias in Data:** If the dataset is biased, passive learning models will also inherit that bias, and correcting this may be difficult without external intervention.

b) **What are the advantages of curious agents?**

- **Continuous Learning:** Curious agents are capable of actively seeking out new information or data, enabling them to continually learn and adapt.
- **Exploration:** They can explore areas of the problem space that have not been sufficiently explored, leading to better generalization and improved performance.
- **Efficient Learning:** By focusing on areas where they are uncertain or where the data is sparse, curious agents can achieve more effective learning with fewer labeled samples.
- **Autonomy:** Curious agents can operate independently, identifying tasks that need attention or learning opportunities without requiring constant human supervision.

c) **What is active learning? What is the relation to HITL?**

- **Active Learning** is a machine learning approach where the model selects which data points it would like to be labeled, usually from an unlabeled pool, to improve its performance more efficiently. The model actively queries an oracle (usually a human expert) for labels on specific instances that it is uncertain about.
- **Relation to HITL:** In the context of **Human-in-the-Loop (HITL)**, active learning involves humans providing labels for uncertain data points, making it a collaborative process. The model asks for human assistance when it encounters situations where it lacks confidence in its predictions, combining AI's decision-making with human expertise.

d) **Distinguish the different Active Learning Scenarios (stream-based, pool-based, query synthesis).**

- **Stream-based Learning:** In this scenario, data arrives sequentially in a stream, and the model must decide whether to request labels for the current instance in real-time. It does not have access to a large pool of data but rather makes decisions based on each new piece of data as it arrives.

- **Pool-based Learning:** Here, there is a large pool of unlabeled data, and the active learning model selects instances from this pool to be labeled. The model can query for the labels of specific instances it finds most informative to improve its learning.
 - **Query Synthesis:** In this case, the agent doesn't select from a set of unlabeled instances but rather synthesizes new instances that it believes will be most useful for improving the model. These generated instances are then sent to the oracle for labeling.
- e) **Distinguish the different Query Strategies (Uncertainty sampling; Diversity sampling; Hybrid sampling; Transfer learning-based sampling) Active Learners can employ.**
- **Uncertainty Sampling:** The model selects the data points for which it has the lowest confidence in its predictions (i.e., the most uncertain instances). The goal is to resolve the greatest uncertainty by labeling the instances where the model is least confident.
 - **Diversity Sampling:** Instead of focusing on uncertainty alone, this strategy selects instances that are diverse or representative of the entire data space. It aims to improve the generalization of the model by ensuring that labeled data covers a wide range of cases.
 - **Hybrid Sampling:** This strategy combines both uncertainty sampling and diversity sampling to balance resolving uncertainty while ensuring the data labeled is diverse. This hybrid approach attempts to provide a better overall representation for the model.
 - **Transfer Learning-based Sampling:** Here, the model uses knowledge from a previously trained model or from a related domain to inform which data points to query. This strategy helps the model leverage pre-existing knowledge to improve its learning in a new or similar domain.

f) **How can AI agents learn to defer/delegate tasks?**

- **Task Deferment:** AI agents can learn to defer tasks by identifying situations where the complexity or uncertainty of the task exceeds their capabilities. They can query an oracle or another more capable system for help. This involves defining confidence thresholds or recognizing areas of expertise where the agent is not yet proficient.
- **Delegation:** AI agents can delegate tasks by recognizing patterns in the data or types of decisions that require human or specialized input. They can identify these cases by applying models of uncertainty, risk, or decision thresholds, passing tasks to experts or other agents better suited for certain tasks. For example, an AI may delegate tasks such as diagnosing rare diseases to a specialized medical expert or a human when the confidence in its predictions is low.

WORKSHEET 4 - DATA COLLECTION & ANNOTATION

QUESTION 1.1

Give and discuss an example where Data Mining is crucial to the success of a business. What Data Mining functionalities does this business need (e.g., think of the kinds of patterns that could be mined)? Can such patterns be generated alternatively by data query processing or simple statistical analysis?

- A department store can use **data mining (DM)**, specifically **association rule mining**, to enhance its target marketing campaigns. By identifying patterns such as products frequently bought together, the store can send tailored promotions to customers likely to purchase related items. Unlike **data query processing**, which retrieves specific information, or **statistical analysis**, which identifies simple relationships, data mining can uncover hidden patterns in large datasets. Therefore, DM is crucial for discovering insights that simple queries or statistical methods cannot achieve, helping the store improve marketing efficiency and profitability.

QUESTION 1.2

To better understand their likes and dislikes, an analyst collects surveys from a sample of customers. Subsequently, the analyst uploads all the data to a database, corrects erroneous or missing entries, and designs a recommendation algorithm.

- a) Associate each of the following actions to the respective steps in the Data Mining process:
- a. Conducting surveys – Data Collection
 - b. Uploading to database – Data Collection
 - c. Correcting missing entries – Data Preprocessing
 - d. Designing a recommendation algorithm - Analysis
 - **Data Collection** involves gathering raw data (e.g., conducting surveys and uploading it to a database).
 - **Data Preprocessing** includes cleaning the data, such as correcting missing or erroneous entries.
 - **Analysis** is where you apply algorithms like a recommendation system to mine insights from the data.
- b) Now suppose that the surveys are answered in free text, which is then read by a team of workers that manually extract positive and negative aspects, also to be included in the database. What step would this be?
- Data Annotation. This step involves adding meaningful labels or tags to the raw data, which can then be used for further analysis. It often occurs during the **Data Preprocessing** stage to prepare the data for more advanced analysis or machine learning tasks.

QUESTION 1.3

Imagine that you are developing an application for optimizing student academic performance and discuss:

- a) What kind of data would be useful for this purpose? For each kind, indicate its type.
- **Age:** Numeric data that represents the student's age.
 - **Gender:** Categorical data (e.g., male, female, non-binary).
 - **Social network data:** Numeric (number of friends) and possibly categorical data (information about friends or interests).
 - **Grades:** Numeric data used for training and validation (reflects academic performance).
 - **Activity:** Potentially categorical or time series data (e.g., hours spent on various activities).
 - **Mood:** Time series data (tracking changes in mood over time).
 - **Time spent on various activities:** Numeric and categorical data (e.g., hours spent studying, sleeping, using apps).
 - **Environmental data:** Time series and consecutive discrete data (light, noise levels, etc.).
- b) How could you extract such data? What would your data sources be?
- **Surveys:** Regular surveys (daily, weekly) can provide self-reported data.
 - **School Information System:** Provides structured data like grades and attendance.
 - **Smartphones:** Can track time spent on various activities, apps, and even mood (through sensors and activity tracking).
 - **IoT Sensors:** Devices can track environmental conditions like light, noise, and temperature.
 - **Text Mining:** Extracting sentiment from social media posts, SMS, emails, etc., through natural language processing.
- c) What Data Mining problems would be involved?
- **Classification:** Using historical data to classify students into categories, such as predicted academic performance.
 - **Outlier Detection:** Identifying unusual student behaviors or performances that deviate from the norm.
 - **Clustering:** Grouping students based on certain features like study habits, activity patterns, or performance.

- **Pattern Mining:** Identifying correlations or patterns, such as associations between time spent on social media and academic performance, or mood changes and grades.

d) Regarding the data, what difficulties do you anticipate?

- **Privacy and Data Protection:** Adhering to laws and regulations (e.g., GDPR) regarding the collection and use of student data.
- **Data Availability:** Insufficient data or incomplete responses (e.g., missing survey answers, unavailable sensor readings).
- **Model Errors:** Machine learning models might generate errors due to incomplete data or inaccurate predictions.
- **Multiple Data Sources:** Integrating diverse data from surveys, smartphones, sensors, and school systems can be complex and require data preprocessing.

QUESTION 1.4

Discuss the pros and cons of Crowdsourcing and:

a) Give examples of concrete tasks that suit well and others that do not. You may think of real situations where crowdsourcing is or has been used.

- **Suited Tasks:**
 - **Data Annotation:** Crowdsourcing is ideal for tasks where the solution is obvious to most humans, but the volume of data is too large for a small team. For example:
 - Annotating images for machine learning (e.g., labeling pictures of cats and dogs).
 - Transcribing audio data into text (common in transcription services).
 - Categorizing reviews or comments into positive, neutral, or negative sentiments.
 - **Subjective Annotations:** Crowdsourcing is useful for tasks requiring diverse opinions, especially when a single person's perspective isn't enough:
 - Rating movies or products based on personal preferences.
 - Evaluating the aesthetic appeal of designs or advertisements.
- **Not Suited Tasks:**
 - **Specialized Expertise Required:** Tasks that demand specific knowledge or expertise are less suited for crowdsourcing, as the quality of work may be compromised:
 - Medical diagnoses or legal consultations.
 - Scientific research requiring deep domain expertise.
 - **Highly Complex or Creative Problem-Solving:** Tasks that require original, critical thinking or long-term creative solutions:
 - Developing novel algorithms or engineering solutions.
 - Writing complex academic papers or producing original art.

b) Enumerate reasons that should be weighted when opting for crowdsourcing.

- **Cost Considerations:** Crowdsourcing is often more affordable than hiring specialized experts, especially for tasks requiring large amounts of data processing or simple repetitive tasks.
 - **Example:** Paying a crowd for labeling images can be cheaper than hiring experts.
- **Task Simplicity:** Crowdsourcing works best for tasks that can be easily broken down into smaller parts and are simple enough to be understood by a broad audience.
 - **Example:** Categorizing data points or performing basic fact-checking.
- **Time Constraints:** Crowdsourcing can accelerate processes, especially when tasks can be parallelized and there is a need for quick results.
 - **Example:** Analyzing large volumes of data in a short amount of time.
- **Quality of Results:** For critical tasks where precision and quality are paramount (e.g., medical or financial analysis), crowdsourcing may not be reliable. It's important to assess the quality control mechanisms and possibly combine crowdsourced data with expert verification.

- **Innovation and Fresh Ideas:** Crowdsourcing can offer innovative solutions and ideas that a closed group of experts may not consider, helping in areas like product design or brainstorming.
 - **Example:** Generating creative marketing campaign ideas or product designs.
- c) Refer important aspects that should be taken care of when opting for microwork.
- **Task Simplicity:** Tasks should be easy to understand, with clear instructions. Complicated tasks can lead to errors or confusion among workers.
 - **Example:** A task asking for a simple yes/no answer or a numerical rating is more suited for microwork than writing a detailed review.
 - **Clarity and Objectivity:** The task should be communicated as simply and objectively as possible to avoid subjective interpretations.
 - **Example:** "Rate this image on a scale from 1 to 5 based on quality" is clearer than "How do you feel about this image?"
 - **Quality Control:** It is crucial to include mechanisms for ensuring the quality of work, such as filtering out workers who rush through tasks without taking them seriously. Some strategies include:
 - Reviewing samples of work periodically.
 - Providing a rating or feedback system for workers.
 - Using multiple workers for the same task and cross-checking results.
 - **Outlier Detection:** Systems must be in place to detect and discard outlier responses, such as those from workers who might complete tasks too quickly without proper attention to detail.

d) How could crowdsourcing be exploited in The Price is Right tv show? Could you think of limitations of answering with the linear average of audience guesses?

- In **The Price is Right**, crowdsourcing could be used to gather guesses from the audience for the price of an item. The host could use live audience participation or allow viewers at home to submit guesses via an app or website. These guesses could then be aggregated to determine the most accurate price.
- **Limitations of Linear Average of Audience Guesses:**
 - **Guessing Bias:** Some participants may exaggerate their guesses, intentionally or unintentionally, skewing the average.
 - **Limited Range of Guesses:** If most of the audience guesses too high or too low, averaging will not necessarily result in the most accurate answer.
 - **Excluding Extreme Guesses:** The linear average doesn't account for the possibility that extreme guesses may be outliers, and thus may skew the result away from the actual price.

• Cohen's Kappa coefficient:

$$\kappa = \frac{p_o - p_e}{1 - p_e} = 1 - \frac{1 - p_o}{1 - p_e}$$

• p_o = relative observed agreement among annotators
 • p_e = hypothetical probability of chance agreement
 $p_e = \frac{1}{N^2} \sum_k n_{k1} n_{k2}$
 • $|k|$ categories
 • N labelled items
 • n_{ki} = number of times annotator i labelled category k

WORKSHEET 5 - TEXT MINING

QUESTION 1.1

What Text Mining techniques would be useful for the following scenarios?

- a) You were provided with several gigabytes of textual documents, of which you know nothing about, but believe that they might include valuable explanations for some mysteries of human history. To help you explore, you want to search those documents for specific keywords or phrases.
- **Technique:**

- **Information Retrieval (IR) / Keyword Search:** This is the most straightforward method, where you search through the documents using a **search engine** or a **text-based search algorithm** like **TF-IDF** (Term Frequency-Inverse Document Frequency) to identify the frequency and importance of keywords within the collection.
 - **Methods:**
 - **Keyword-based search:** Simple text search for specific terms.
 - **Indexing:** Create an inverted index to efficiently map keywords to documents.
 - **Boolean search:** Combine terms with logical operators like AND, OR to refine results.
- b) You want to have a high-level picture of the most common themes in the previous collection, without having to read through all the documents.
- **Technique:**
 - **Topic Modeling:** Topic modeling allows you to find hidden themes across a large collection of documents without needing prior knowledge of the content.
 - **Methods:**
 - **Latent Dirichlet Allocation (LDA):** A popular algorithm for discovering abstract topics in a collection of documents.
 - **Non-negative Matrix Factorization (NMF):** Another technique for extracting topics by factorizing the term-document matrix.
 - **Clustering:** **K-means** or **DBSCAN** to group similar documents together based on the content.
- c) You note that about half of the documents have a naming convention that might be related to the actual contents of the document. Towards a better organization, you want to predict the names of the files that do not match this naming convention.
- **Technique:**
 - **Text Classification / Supervised Learning:** You can use a **supervised machine learning** approach to train a model on the documents whose names follow a convention. The model will learn to predict the naming convention of other files based on features extracted from the documents.
 - **Methods:**
 - **Naive Bayes:** A simple probabilistic classifier often used in text classification.
 - **Support Vector Machine (SVM):** For more complex text classification tasks.
 - **Random Forest / Decision Trees:** These can also be used for classification tasks based on text features.
- d) You want to develop a tool for helping in the process of writing emails, i.e., for every received email, it should propose a possible response.
- **Technique:**
 - **Natural Language Processing (NLP) for Text Generation:** This requires generating appropriate responses to emails based on their content. You could use **NLP models** that understand the context of the email and generate relevant replies.
 - **Methods:**
 - **Text Summarization:** Extract the main points of an email for a condensed response (e.g., **extractive summarization**).
 - **Sequence-to-sequence models:** Use deep learning models like **Recurrent Neural Networks (RNNs)** or **Transformer-based models** (e.g., **GPT**, **BERT**) to generate contextually appropriate replies.
 - **Intent Recognition:** Identify the user's intent in the email (e.g., query, complaint, request) and match it to an appropriate response template using **classification**.

QUESTION 1.2

Consider the following text corpus with six documents:

D1	Good-looking food... that tasted bad!!
D2	Service not good. Food not good.
D3	Terrible, terrible food.
D4	Amazing tasting food. Best value.
D5	Good service. Good food. Good value for money.
D6	Look: best-tasting food.

Now answer the following questions:

- a) If you were doing Text Mining from documents of this kind, would you perform any kind of pre-processing?
- **Lowercasing:** Convert all text to lowercase to avoid duplicates based on case differences (e.g., "Food" vs. "food").
 - **Removing Punctuation:** Punctuation does not usually contribute to the meaning of the text for most text mining tasks, so it should be removed.
 - **Removing Stop Words:** Words like "that", "for", "a", etc., do not carry significant meaning and can be safely excluded.
 - **Lemmatization:** Convert verbs in their gerund or plural form to their base form (e.g., "looking" becomes "look", "tasting" becomes "taste").
 - **Tokenization:** Split text into individual terms (words).
 - **Handling Negations:** In sentiment analysis, negations can flip the meaning of a sentence (e.g., "not good" vs. "good"). You might consider handling these explicitly.
- b) Adopt the term-document matrix, based on frequency counts, for representing the term-document matrix of this corpus. In the process, ignore punctuation signs, prepositions (that, for) and lemmatize verbs in the gerund (looking, tasting).
- We will create a term-document matrix by following the pre-processing steps:
 - Documents: d1, d2, d3, d4, d5, d6
 - Terms: good, food, service, taste, terrible, best, value, look

Term	D1	D2	D3	D4	D5	D6
good	1	2	0	0	3	0
food	1	1	1	1	1	1
service	0	1	0	0	1	0
taste	1	0	0	1	0	1
terrible	0	0	2	0	0	0
best	0	0	0	1	0	1
value	0	0	0	1	1	0
look	1	0	0	0	0	1
bad	1	0	0	0	0	0
amazing	0	0	0	1	0	0
money	0	0	0	0	1	0

- c) Using the previous vector representation, compute the similarity between: d1 and d3; d1 and d4; d1 and d5.
- $\text{Cosine Similarity} = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$ $\cos(d_1, d_3) = \frac{1}{\sqrt{5} \times \sqrt{1+4}} = \frac{1}{5}$
- d) Which tasks in the domain of Sentiment Analysis could be applied to this data? Clarify their goal and the result for each document.
- **Sentiment Classification:** Classifying the sentiment of each document (positive, negative, neutral).
 - Example for d1: Neutral
 - Example for d2: Negative
 - Example for d3: Negative
 - Example for d4: Positive
 - Example for d5: Positive

- Example for d6: **Positive**
 - **Polarity Classification:** Identifying if the sentiment is positive or negative (binary classification).
 - For most documents, polarity can be determined (positive for d4, d5, d6, and negative for d1, d2, d3).
 - **Emotion Detection:** Determining emotions like joy, sadness, anger, etc., from the content of the text.
- e) What would be the limitations of relying exclusively on a vector representation based on frequency counts, and possibly a sentiment lexicon? How could they be minimized, for different Sentiment Analysis tasks? Consider different scenarios regarding annotated data and computer power available.
- **Loss of Context:** Frequency-based representation does not capture word order or context, which could lead to misinterpretations of sentiment (e.g., "not good" vs. "good").
 - **Minimizing limitation:** Incorporating context-based models like **word embeddings** (Word2Vec, GloVe) can reduce this limitation.
 - **Neglect of Sarcasm:** Frequency-based methods and sentiment lexicons may not detect sarcasm.
 - **Minimizing limitation:** Using more advanced models like **transformer-based models** (e.g., BERT) can help identify sarcasm.
 - **Inability to Handle Polysemy:** Words with multiple meanings (e.g., "bad" could be negative sentiment or bad quality) can confuse frequency-based methods.
 - **Minimizing limitation:** Implementing disambiguation algorithms or using context-aware models like **BERT**.
 - **Dependency on Annotated Data:** Models trained with annotated data require a large, well-annotated dataset for accuracy, which may be costly or time-consuming to create.
 - **Minimizing limitation:** Use semi-supervised learning or pre-trained models to improve performance with fewer labeled examples.
- f) Consider that the polarity of these documents was annotated independently, by two humans, who tagged them as follows:
- a. Ana: Neutral, Neutral, Negative, Positive, Positive, Positive
 - b. Hugo: Negative, Negative, Negative, Positive, Positive, Positive
- How would you classify the annotator agreement?
- To evaluate the agreement between Ana and Hugo, we can compute **Cohen's Kappa**, which is a statistical measure of inter-rater reliability, where:
 - P_0 is the observed agreement (how often they agreed).
 - P_e is the expected agreement by chance.
 - In this case:
 - **Ana's labels:** Neutral, Neutral, Negative, Positive, Positive, Positive
 - **Hugo's labels:** Negative, Negative, Negative, Positive, Positive, Positive
 - The agreement is not perfect, and there is some divergence in their labeling, especially for d1 and d3. Calculating Cohen's Kappa would give us a numerical value for this agreement, where a value closer to 1 indicates strong agreement, and a value closer to 0 indicates no agreement beyond chance.

WORKSHEET 6 - HUMAN IN THE LOOP MACHINE LEARNING (PART II)

LEARNING FROM DEMONSTRATION (LfD)

QUESTION 6.6

Define Learning from Demonstration (LfD) and explain its significance in teaching AI models.

- **Learning from Demonstration (LfD)** refers to a machine learning approach where an AI system learns a task by observing demonstrations provided by a human or another agent, instead of being explicitly programmed with rules or relying solely on reinforcement learning with rewards. In this approach, the AI agent observes sequences of states and actions (called demonstration trajectories) and attempts to replicate the demonstrated behavior.
- **Significance in Teaching AI Models:**
 - **Simplifies Training:** LfD reduces the need for extensive programming or manually creating reward functions that would guide the learning process. This is particularly useful in cases where designing a reward function is challenging or impractical.
 - **Human-Centric Learning:** By observing human demonstrations, AI systems can learn human-centric tasks more intuitively. This is especially valuable for applications like human-robot interaction, task automation, or personalized systems where replicating human actions is often essential.
 - **Faster Learning:** Rather than starting from scratch, the AI system leverages the human-provided demonstrations to learn quickly, making it possible to handle complex tasks where traditional programming might be time-consuming or inefficient.
 - **Improves Task Generalization:** Through observing the demonstrations, the AI model can learn general patterns and structures that allow it to adapt to new, unseen situations. This flexibility is key for real-world applications where conditions can vary widely.
 - **Better Handling of Complex Tasks:** LfD is especially useful when tasks are too intricate for simple rule-based systems, making it ideal for tasks that require adaptive, context-aware behavior.
- In summary, **Learning from Demonstration** is a powerful approach in AI because it allows systems to quickly learn human-like tasks without requiring detailed coding or complex reward functions, leading to more efficient and adaptable AI systems in practical applications.

QUESTION 6.7

What are the differences between direct (behavior cloning) and indirect (inverse reinforcement learning) methods in Learning from Demonstration?

- **Behavior Cloning (Direct Method):**
 - **Approach:** This method directly maps observed states to actions, essentially teaching the AI to imitate the actions demonstrated by the expert. The AI model learns the policy by replicating the behavior demonstrated in the training trajectories.
 - **Training Type:** It is a supervised learning process where the AI is trained on a dataset of demonstration trajectories, and the objective is to predict the actions based on the observed states.
 - **Strengths:**
 - Simple and easy to implement.
 - Works well in environments where the demonstrated behavior is highly relevant and does not change much.
 - Can be faster to train as it doesn't require learning a reward function.
 - **Limitations:**
 - Struggles to generalize in novel situations or those not covered by the demonstration data.
 - Cannot adapt to new circumstances if they differ significantly from the demonstration.
- **Inverse Reinforcement Learning (IRL) (Indirect Method):**
 - **Approach:** Rather than directly mapping states to actions, IRL seeks to infer the underlying **reward function** that drives the expert's actions. The AI learns why certain actions were taken by the expert based on their observed behavior.
 - **Training Type:** The agent tries to recover the expert's objective (or reward function) from the demonstrated behavior, which then allows it to derive an optimal policy. It's a more complex, indirect approach.
 - **Strengths:**
 - More flexible and adaptable as it doesn't just mimic actions, but learns the driving motivations behind them.

- Generalizes better to unseen situations because the learned reward function can guide the agent to make decisions in novel contexts.
 - Suitable for complex, dynamic environments where the goal or task may not be entirely apparent in the demonstration.
 - **Limitations:**
 - More computationally intensive and requires more sophisticated algorithms to infer the reward function.
 - Can be harder to implement due to the need to accurately infer the reward structure.
 - **Key Differences:**
 - **Direct Method (Behavior Cloning):** Focuses on action imitation by mapping states to actions directly. Simpler and faster but struggles with generalization in novel situations.
 - **Indirect Method (Inverse Reinforcement Learning):** Focuses on inferring the reward function behind expert actions, leading to more generalizable behavior, but at the cost of increased complexity and computational effort.
 - In essence, behavior cloning is quicker and simpler but less adaptable to new situations, while IRL is more flexible and capable of handling complex environments but is computationally more demanding.
-

QUESTION 6.8

How does behavior cloning work, and what are some common challenges associated with it?

- Behavior cloning is a method of teaching an AI to mimic the actions of an expert through supervised learning. The process is as follows:
 - **Training Data:** The AI is provided with a set of demonstration data, consisting of observations (states) and corresponding actions taken by the expert.
 - **Learning Objective:** The model learns to map each observation (state) to the action that the expert would take in that context.
 - **Minimization:** The AI minimizes the error between its predicted actions and the actions taken by the expert in similar states.
 - **Outcome:** The result is a policy that replicates the expert's behavior for new situations that resemble the training data.
 - **Challenges with Behavior Cloning:**
 - **Compounding Errors:**
 - If the model makes a mistake early in the sequence, the error can propagate and amplify, leading to a series of incorrect actions. Since behavior cloning doesn't offer feedback or recovery mechanisms, these early mistakes can cause significant performance degradation.
 - **Limited Generalization:**
 - Behavior cloning relies entirely on the training data, and as such, it struggles to generalize to new or unseen situations. If the model encounters states that weren't present in the training set, it might not know how to act. This lack of adaptability can be a major limitation when dealing with complex environments or tasks.
 - **Simplicity Assumption:**
 - Behavior cloning assumes that the expert's policy is simple enough to be directly learned from the data. This assumption may fail when the task requires more complex decision-making processes, which cannot be directly captured by mimicking the expert's actions.
-

QUESTION 6.9

Explain Inverse Reinforcement Learning (IRL) and how it enables AI to infer the objectives behind human actions.

- Inverse Reinforcement Learning (IRL) is a machine learning approach that focuses on uncovering the **underlying reward function** that explains the behavior of an expert. Instead of directly learning a policy to maximize a known reward, IRL works in reverse by:

- **Observing Expert Behavior:** The AI observes an expert performing tasks, capturing the actions taken in different states.
- **Inferring the Reward Function:** The model analyzes the observed actions and infers the reward function that the expert is optimizing.
- **Policy Derivation:** Once the reward function is inferred, the AI can derive a policy that maximizes this reward, enabling it to make decisions even in new, unseen scenarios.
- **How IRL Enables AI to Infer Human Objectives:**
 - **Understanding Motivations:** IRL deduces the goals and preferences that likely motivated the expert's actions by analyzing patterns in the behavior.
 - **Generalization:** By learning the rationale behind actions (via the reward function), the AI can adapt to novel environments and challenges, making decisions aligned with the inferred objectives.
 - **Alignment with Human Values:** IRL provides insights into the decision-making process of humans, allowing the AI to align its behavior more closely with human values and goals.
- **Advantages of IRL:**
 - **Flexibility:** Enables AI to adapt to scenarios not explicitly demonstrated by the expert.
 - **Insightful Decision-Making:** Helps the AI understand "why" actions are taken, rather than just "what" actions to mimic.
 - **Applications in Complex Systems:** Particularly useful in domains requiring adaptability and value alignment, such as robotics, autonomous vehicles, and human-computer interaction.
- **Key Applications:**
 - **Robotics:** Understanding human goals in tasks like assembly or navigation to design robots that assist effectively.
 - **Autonomous Driving:** Inferring the intent behind human driving behaviors to enhance decision-making in dynamic environments.
 - **Human-Computer Interaction:** Aligning AI responses with user preferences and intentions to improve collaboration.

QUESTION 6.10

Give an example of how Learning from Demonstration could be used in robotics. What challenges might arise?

- **Example: Teaching a Robotic Arm to Assemble Furniture**
 - A human expert demonstrates assembling a piece of furniture, showing the robotic arm how to:
 - Identify and pick up components.
 - Align and connect parts accurately.
 - Use tools like screws and wrenches.
 - Apply appropriate force and torque for tightening screws.
 - The robot learns by observing these demonstrations, using sensors like cameras and force/torque sensors to gather data. It maps the observed actions to its motor commands, creating a model to replicate the process autonomously.
- **Challenges in Learning from Demonstration (LfD) for Robotics:**
 - **Generalization:** Difficulty adapting to variations in tasks or environments beyond the demonstrated scenarios.
 - **Error Accumulation:** Early mistakes can cascade, leading to task failure.
 - **Sensor Noise:** Imprecise sensor data can cause misinterpretation of actions.
 - **Environmental Changes:** Variability in workspaces or tools can impact performance.
 - **Limited Demonstrations:** A small number of examples restrict robustness.
 - **Skill Transfer:** Challenges in translating human actions to robotic movements due to differences in dexterity and precision.
- Solutions include techniques like reinforcement learning, real-time feedback, sensor fusion, domain adaptation, and synthetic data generation to improve adaptability and robustness.

QUESTION 6.11

What is Reinforcement Learning from Human Feedback (RLHF), and how does it enhance traditional reinforcement learning?

- **Reinforcement Learning from Human Feedback (RLHF)** is a technique where human-provided feedback guides a reinforcement learning agent's training. Instead of relying solely on predefined reward functions, RLHF incorporates human preferences, comparisons, or corrections to improve the agent's behavior.
- **Enhancements over Traditional Reinforcement Learning:**
 - **Complex Tasks:** RLHF enables learning in tasks where designing an explicit reward function is challenging or infeasible.
 - **Human Alignment:** Ensures that the agent's behavior aligns more closely with human expectations, preferences, and values.
 - **Faster Convergence:** Human feedback provides additional, often more informative signals, leading to more efficient learning.

QUESTION 6.12

Describe the three types of human feedback commonly used in RLHF: Comparison-Based Feedback, Reward Modeling, and Corrective Feedback.

- **Comparison-Based Feedback:**
 - Humans compare two outputs or behaviors from the agent and indicate which one is preferred.
 - Example: A user ranks two chatbot responses, allowing the model to learn which is more appropriate.
 - **Advantage:** Easier for humans to rank options than to define explicit rewards.
- **Reward Modeling:**
 - Humans provide explicit feedback on specific outputs, which is used to train a reward model.
 - The reward model then predicts the reward for new actions, guiding the agent during training.
 - **Advantage:** Scales human feedback by enabling the model to infer rewards for unseen situations.
- **Corrective Feedback:**
 - Humans directly intervene to correct or modify the agent's behavior when it makes errors.
 - Example: Correcting a robot's grasp when it fails to pick up an object.
 - **Advantage:** Provides precise guidance in error-prone scenarios.

QUESTION 6.13

How is RLHF used to train language models like ChatGPT? What is the role of human feedback in this context?

- **How RLHF is Used in ChatGPT Training:**
 - **Supervised Fine-Tuning:**
 - The model is first trained on a dataset of human-written examples to establish a base understanding of language and task-relevant patterns.
 - **Human Feedback for Comparison-Based Learning:**
 - Human reviewers rank different responses generated by the model for the same prompt.
 - The rankings train a **reward model** that predicts which responses align better with human preferences.
 - **Reinforcement Learning:**
 - The reward model guides a reinforcement learning process (e.g., Proximal Policy Optimization) to adjust the language model, optimizing it to generate more preferred outputs.
- **Role of Human Feedback in ChatGPT:**
 - **Alignment:** Ensures the model produces helpful, accurate, and contextually appropriate responses.

- **Ethical Behavior:** Guides the model away from harmful or biased content by teaching it human-aligned norms and values.
- **Iterative Improvement:** Human feedback helps refine the model iteratively, addressing gaps in its understanding or behavior.

QUESTION 6.14

Advantages and Potential Pitfalls of Using RLHF in High-Stakes Applications like Autonomous Vehicles

- **Advantages:**
 - **Alignment with Human Values:** RLHF helps ensure that autonomous vehicles behave in ways consistent with human expectations, such as prioritizing safety over efficiency.
 - **Handling Ambiguous Scenarios:** Human feedback provides insights into how to handle complex, unstructured, or rare scenarios that are difficult to encode in traditional reward functions (e.g., navigating pedestrian-heavy zones).
 - **Improved Decision-Making:** Incorporating human preferences can enhance the vehicle's ability to make nuanced trade-offs, such as balancing speed and comfort.
 - **Adaptability:** RLHF allows the vehicle to refine its behavior over time based on updated human feedback, ensuring it evolves with changing norms and standards.
- **Potential Pitfalls:**
 - **Bias in Feedback:** Human feedback may reflect individual biases or inconsistencies, which could lead to undesirable or unsafe behaviors being learned.
 - **Scalability Issues:** Gathering sufficient high-quality feedback for diverse driving scenarios can be resource-intensive and challenging to scale.
 - **Unintended Consequences:** Overfitting to specific human feedback might cause the vehicle to behave sub optimally in scenarios outside the feedback's scope.
 - **Delay in Learning:** Relying on human feedback can slow down the learning process, especially in environments requiring real-time decision-making.
 - **Ethical Dilemmas:** High-stakes applications often involve ethical considerations (e.g., choosing between two harmful outcomes). RLHF may inherit and amplify human disagreement on these dilemmas.

QUESTION 6.15

In What Scenarios is RLHF Particularly Useful? Why Might It Be Preferred Over Traditional Reinforcement Learning?

- **Scenarios Where RLHF is Useful:**
 - **Complex or Implicit Objectives:** Tasks like conversational AI, where defining explicit rewards is difficult.
 - **Dynamic Environments:** Situations like autonomous driving, involving unpredictable and varied conditions.
 - **Ethically Sensitive Applications:** Healthcare or safety-critical tasks requiring alignment with human values.
 - **Human-Like Behavior:** Collaborative robotics needing intuitive and natural actions.
- **Why RLHF is Preferred Over Traditional RL:**
 - Aligns agent behavior with human expectations and values.
 - Accelerates learning in tasks with sparse or complex rewards.
 - Adapts better to rare or unforeseen edge cases.
 - Handles ambiguity in objectives through nuanced human guidance.
 - Reduces risks of unintended behaviors in safety-critical scenarios.

LEARNING FROM NATURAL LANGUAGE

QUESTION 6.16

How Can Natural Language Be Used as an Instructional Medium for Training AI?

- Natural language can serve as a user-friendly and flexible way to convey instructions, objectives, or feedback to AI systems. By interpreting natural language, AI can align its behavior with human intentions, even in complex or dynamic scenarios.
- **Applications:**
 - Teaching tasks to robots using spoken or written commands.
 - Guiding conversational agents to refine tone, relevance, or content.
 - Adapting AI behavior dynamically through user-provided prompts or preferences.

QUESTION 6.17

Explain Language-Based Reward Shaping and Give an Example of Its Application

- **Language-Based Reward Shaping** involves using natural language descriptions to provide additional context or guidance for the reward function in reinforcement learning.
- **Example:**
 - A robot is tasked with cleaning a room. Alongside its primary reward (e.g., for removing visible debris), natural language input like “*Focus on cleaning the corners*” helps refine the reward signal, directing the robot’s attention to often-overlooked areas.
- **Benefits:**
 - Improves learning efficiency by incorporating human insights.
 - Enables fine-grained adjustments to AI behavior during training.

QUESTION 6.18

How Does Instruction-Based Reinforcement Learning Work, and in What Contexts Might It Be Useful?

- **Instruction-Based Reinforcement Learning** uses natural language instructions as part of the input to the agent. These instructions define goals, constraints, or task-specific priorities, guiding the agent’s decision-making process.
- **How it Works:**
 - The agent processes natural language instructions alongside its environmental inputs.
 - Models, such as transformers, map instructions to actionable policies or strategies.
- **Contexts of Use:**
 - **Task Generalization:** Teaching robots or AI to perform multiple tasks (e.g., “*Sort these items by size*” or “*Water the plants*”).
 - **Dynamic Goal Setting:** Allowing AI systems to adapt to evolving tasks without retraining.
 - **Human-AI Collaboration:** Enhancing interactions in applications like virtual assistants or autonomous systems.

QUESTION 6.19

Challenges of Using Natural Language for AI Instruction

- **Ambiguity:** Natural language is often imprecise, leading to misinterpretation of instructions.
- **Complexity:** Some tasks may require highly detailed instructions that are cumbersome to articulate.
- **Language Variability:** Differences in phrasing, regional language variations, or colloquialisms may hinder understanding.
- **Lack of Context:** AI may struggle without sufficient environmental or situational context to disambiguate instructions.
- **Error Propagation:** Misinterpreted instructions can lead to compounding errors in task execution.

QUESTION 6.20

How Could Natural Language Instructions Be Combined with RLHF to Improve AI Performance and Alignment?

- **Combining Natural Language with RLHF:**
 - **Enhanced Feedback:** Natural language enables more nuanced feedback during RLHF training (e.g., “*The response was too abrupt; try to be more polite*”).
 - **Reward Shaping:** Human instructions can refine the reward model by providing additional context or priorities.
 - **Dynamic Goal Setting:** Natural language allows users to adjust AI goals and preferences dynamically, improving adaptability in real-time applications.
 - **Interpretability:** Combining RLHF with language instructions makes AI decisions more transparent by aligning them with explicitly stated user intentions.
- **Example:** Training a conversational AI with RLHF where users provide natural language feedback on tone and relevance, alongside preference rankings, to refine its responses.

PREDICTING HUMAN BELIEFS, DESIRES, AND INTENTIONS

QUESTION 6.21

What is Cognitive Modeling in the Context of Human-in-the-Loop Learning?

- **Cognitive Modeling** involves creating computational representations of human thought processes, such as decision-making, problem-solving, and learning, to improve AI systems in **Human-in-the-Loop (HITL)** frameworks. These models simulate how humans perceive, reason, and act, enabling the AI to better interpret and respond to human actions.
- **Application in HITL Learning:**
 - Enhances AI's ability to predict user preferences and adapt dynamically based on human feedback.
 - Example: A tutoring system using cognitive models to gauge a student's understanding and adjust teaching strategies in real time.

QUESTION 6.22

How Can AI Systems Predict Human Beliefs, Desires, and Intentions to Improve Personalization?

- AI systems predict human beliefs, desires, and intentions (BDI) by analyzing user behavior, contextual data, and interaction history. They use techniques such as:
 - **Behavioral Analysis:**
 - Tracking clicks, searches, or actions to infer preferences.
 - Example: Analyzing shopping patterns to predict product preferences.
 - **Natural Language Processing:**
 - Interpreting user queries or feedback to determine intentions.
 - Example: Recognizing “*I'm planning a vacation*” to suggest travel deals.
 - **Machine Learning Models:**
 - Building predictive models using historical and contextual data.
 - Example: Recommending healthcare options based on prior appointments and health records.
- **Result:** Improved personalization enhances user satisfaction by anticipating and fulfilling needs proactively.

QUESTION 6.23

Describe Goal Recognition Models and Explain How They Might Be Applied in Real-World AI Systems

- **Goal Recognition Models** are computational frameworks designed to infer a user's goals based on observed actions and contextual cues.

- **How They Work:**
 - Use probabilistic or machine learning techniques to predict the likely goal from a sequence of actions.
 - Example: Identifying that a user searching for "budget flights" and "hostel reviews" aims to plan a low-cost trip.
- **Applications:**
 - **Customer Support:**
 - Identifying whether a customer seeks troubleshooting help or product recommendations based on interaction patterns.
 - **Autonomous Systems:**
 - Helping robots predict human goals (e.g., reaching for an object) to assist in collaborative tasks.
 - **Gaming AI:**
 - Adapting gameplay strategies based on inferred player objectives.

QUESTION 6.24

What Are Some Ethical Considerations in Designing AI Systems That Predict Human Beliefs or Intentions?

- **Privacy Concerns:**
 - Predicting beliefs or intentions may require collecting sensitive data, raising concerns about user consent and data security.
- **Bias and Misinterpretation:**
 - Predictions based on biased data can reinforce stereotypes or make inaccurate assumptions.
- **Transparency and Consent:**
 - Users may be unaware that their actions are being analyzed, undermining trust. Clear communication and consent mechanisms are essential.
- **Manipulation Risks:**
 - Predictive systems could exploit knowledge of user intentions for manipulative purposes, such as targeted advertising.
- **Accountability:**
 - When predictions are wrong, it's crucial to determine responsibility, especially in high-stakes areas like healthcare.

QUESTION 6.25

Examples of Personalized AI Systems Enhancing User Experience

- **E-Commerce:**
 - **Example:** An AI system anticipates user needs by suggesting frequently purchased items or bundling relevant products (e.g., recommending chargers with a smartphone purchase).
 - **Impact:** Increases convenience and satisfaction while boosting sales.
- **Healthcare:**
 - **Example:** Personalized health assistants predict appointment needs, medication refills, or early warning signs based on health data.
 - **Impact:** Enhances preventive care and reduces health risks by acting proactively.
- By tailoring experiences to individual needs, such systems improve efficiency, engagement, and overall satisfaction.

APPLICATIONS AND CASE STUDIES

QUESTION 6.26

Real-World Applications of Learning from Demonstration (LfD) and RLHF in Robotics

- **Learning from Demonstration (LfD):**
 - **Assembly Tasks:**
 - Teaching robotic arms to assemble components by imitating human demonstrations.
 - **Example:** Robots in manufacturing assembling parts of cars or electronics.
 - **Assistive Robotics:**
 - Training robots to assist elderly or disabled individuals by observing human caregivers.
 - **Example:** Helping with meal preparation or object retrieval.
- **Reinforcement Learning from Human Feedback (RLHF):**
 - **Navigation:**
 - Robots learn optimal paths in dynamic environments, guided by human feedback.
 - **Example:** Warehouse robots efficiently avoiding obstacles.
 - **Collaboration:**
 - Improving human-robot interaction during joint tasks.
 - **Example:** Robots in surgeries refining actions based on surgeons' input.

QUESTION 6.27

RLHF and LfD in Autonomous Driving Systems

- **LfD in Autonomous Driving:**
 - Human drivers demonstrate optimal behaviors, such as safe lane changes or negotiating intersections.
 - **Benefit:** Helps the system learn complex driving maneuvers and adapt to nuanced scenarios.
- **RLHF in Autonomous Driving:**
 - Human feedback refines decision-making policies, prioritizing safety and comfort.
 - **Example:** Feedback on braking patterns to achieve smoother stops.
- **Combined Impact:**
 - **Safety:** Enhanced understanding of edge cases, like pedestrian-heavy areas.
 - **Decision-Making:** Balances objectives such as efficiency, safety, and passenger comfort.

QUESTION 6.28

RLHF in Gaming to Improve AI Performance and User Experience

- **Adaptive AI Behavior:**
 - RLHF enables NPCs to adjust strategies based on player preferences or skill levels.
 - **Example:** AI opponents in strategy games becoming more challenging but fair.
- **Player-Centric Storylines:**
 - Human feedback guides AI in creating narratives that align with player choices and interests.
 - **Example:** Branching storylines in RPGs influenced by user feedback.
- **Enhancing Engagement:**
 - AI learns to provide hints or assistance dynamically without disrupting gameplay.
 - **Example:** Virtual guides in puzzle games that offer timely advice.

QUESTION 6.29

Cognitive Modeling in Virtual Assistants and Customer Service Bots

- **Enhancing Functionality:**
 - **Understanding User Intent:**
 - Cognitive models predict user needs based on context and interaction history.
 - **Example:** Virtual assistants preemptively suggest calendar events based on conversations.
 - **Personalized Interactions:**
 - Adapts tone, language, and suggestions to user preferences.
 - **Example:** Customer service bots adjusting formality based on user profiles.

- **Improved Context Awareness:**
 - Recognizing multi-turn dialogue contexts for better continuity.
 - **Example:** Resolving customer complaints efficiently by tracking prior issues.

QUESTION 6.30

Human Feedback in Enhancing Language Models

- **Role of Human Feedback:**
 - **Fine-Tuning Behavior:**
 - Feedback refines model responses to be more accurate, helpful, and context appropriate.
 - **Example:** Users marking answers as helpful or not to train the model.
 - **Adapting to Evolving Norms:**
 - Incorporating feedback ensures the model stays relevant and culturally appropriate.
 - **Example:** Avoiding outdated or offensive language.
 - **Addressing Ethical Concerns:**
 - Feedback helps prevent biases, ensuring fairness and inclusivity.
 - **Example:** Removing gender or racial bias from generated content.
- **Contributions to Adaptability:**
 - Enables continuous learning and improvement based on real-world usage.
 - Ensures the model aligns with user expectations across diverse applications.

CRITICAL ANALYSIS AND FUTURE DIRECTIONS

QUESTION 6.31

Challenges in Scaling Human-in-the-Loop Machine Learning for Complex AI Systems

- **Resource Intensive:**
 - Scaling requires a large number of high-quality, labeled data from human experts, which is time-consuming and costly.
 - Example: In robotics, human demonstrations may need to cover a vast range of scenarios, leading to inefficiencies.
- **Quality of Feedback:**
 - Ensuring consistent, accurate, and unbiased feedback across diverse human participants becomes harder as the scale increases.
- **Contextual Understanding:**
 - Complex AI systems often operate in dynamic environments where human input may not always cover every possible context or edge case.
- **Latency and Real-Time Needs:**
 - In applications like autonomous driving or healthcare, the need for real-time feedback can overwhelm the system, slowing down decision-making.
- **Integration with Automation:**
 - Seamlessly integrating human feedback into automated systems without disrupting performance or introducing errors remains a significant challenge.

QUESTION 6.32

Implications of Biased Human Feedback in HitL Systems

- **Reinforcement of Biases:**
 - Biased feedback can lead to AI systems learning and perpetuating societal, cultural, or individual biases.
 - **Example:** AI systems trained with biased feedback may make unfair decisions, like hiring biases based on gender or race.

- **Decreased Performance:**
 - Misaligned or skewed feedback can degrade the performance of the AI, leading it to focus on inappropriate objectives or behave in unintended ways.
 - Example: A facial recognition system trained with biased data may misidentify individuals of certain ethnicities.
- **Ethical and Legal Risks:**
 - Biased decisions can introduce ethical dilemmas and legal consequences, especially in high-stakes applications like finance or healthcare.
- **User Trust and Adoption:**
 - If users recognize biases in AI systems, they may lose trust, undermining the system's effectiveness and widespread adoption.

QUESTION 6.33

Potential Long-Term Impacts of Relying on Human Feedback to Guide AI Development

- **Dependence on Human Expertise:**
 - Continued reliance on human feedback may result in AI systems that are overly dependent on specific experts, leading to a lack of generalization and adaptability.
- **Reinforcement of Human Biases:**
 - The longer AI systems rely on human feedback, the more they may inherit human biases, potentially creating AI systems that reflect outdated or harmful societal norms.
- **Stagnation of Innovation:**
 - If AI systems continually depend on human input, they may be less likely to develop novel solutions or approaches on their own, hindering innovation.
- **Ethical Challenges:**
 - As AI becomes more integrated into everyday life, it could shape social structures and individual behaviors based on flawed or incomplete human feedback, potentially exacerbating inequality.
- **Automation of Low-Level Tasks:**
 - Many AI tasks may shift towards fully automated systems that no longer need human feedback for routine functions, reducing the role of humans in the loop but potentially leaving them in supervisory or ethical decision-making roles.

QUESTION 6.34

Active Research Areas in Enhancing HitL, LfD, and RLHF, and Future Improvements

- **HitL Research Areas:**
 - **Better Feedback Mechanisms:** Research is focused on improving the quality, consistency, and scalability of human feedback, especially in real-time systems.
 - **Efficiency in Scaling:** Techniques are being developed to reduce the resource burden of gathering human feedback, such as using crowdsourcing or leveraging semi-supervised learning.
 - **Human-AI Collaboration:** Focus on developing more effective and seamless interfaces that allow humans and machines to work together with minimal friction.
- **LfD Research Areas:**
 - **Robustness in Learning:** Enhancing robots' ability to generalize learned actions across different environments and tasks, especially when human demonstrations are limited.
 - **Multi-modal Learning:** Integrating multiple types of sensory inputs (visual, tactile, auditory) for more accurate imitation of human actions.
- **RLHF Research Areas:**
 - **Improved Reward Models:** Developing more effective ways of incorporating human feedback into reinforcement learning algorithms.

- **Generalization and Adaptability:** Research on making RLHF systems more adaptable to new environments and scenarios without requiring extensive retraining.
- **Future Improvements:**
 - Expect advancements in automated feedback generation, the ability for AI systems to autonomously request feedback when needed, and more dynamic, scalable approaches to integrating human input.

QUESTION 6.35

How Future Advances in Human-in-the-Loop Learning Could Affect Fields Like Healthcare, Autonomous Systems, or Creative Industries

- **Healthcare:**
 - **Personalized Treatment Plans:** Advanced HITL systems could allow AI to learn from both expert feedback and patient data, creating highly personalized treatment plans.
 - **Medical Robotics:** Robots guided by human feedback could perform complex surgeries with greater precision and adapt to dynamic patient needs in real time.
 - **Diagnostic Assistance:** AI could provide doctors with more accurate diagnoses by integrating human feedback and expanding its knowledge base continuously.
- **Autonomous Systems:**
 - **Safer Autonomous Vehicles:** RLHF and LfD will continue to improve the decision-making of self-driving cars, ensuring they make safer, more humane-like decisions in unpredictable environments.
 - **Drone Operations:** HITL systems could help drones improve their ability to navigate complex environments, like disaster areas, by receiving real-time feedback from human operators.
- **Creative Industries:**
 - **Personalized Content Creation:** AI systems in music, film, and design could use RLHF to create highly personalized content, responding to individual preferences and feedback.
 - **Human-AI Collaboration in Art:** Artists could work more seamlessly with AI to generate new ideas or co-create, with AI systems learning from feedback to enhance their outputs.
 - **Game Design:** In gaming, RLHF and LfD will allow AI to evolve in real-time, offering tailored experiences for players based on their unique playing styles and feedback.
- In all these fields, future advances in HITL learning could enable more intuitive, personalized, and efficient systems that better align with human needs and preferences, driving innovation and improving user experiences.

ADAPTED FROM THE NORMAL EXAM OF AI-2020 EDITION

For each of the applications/scenarios described below, indicate which technology or combination of technologies (Reinforcement Learning (RL), or Learning from Demonstration/Apprenticeship Learning – Behavior Cloning (IL-BC), or Learning from Demonstration/Apprenticeship Learning – Inverse RL (ILRL)) is best suited so that the best performance is achieved. Justify your answer.

- a) Consider an E-learning system in which an Artificial Intelligent Personal Assistant (AIPA) is to be integrated. This AIPA should build learning paths (sequence of learning activities – practical, theoretical exercises, classes, seminars, etc.) personalized for each student, i.e., it should obtain a function that specifies what the student should do in certain circumstances (student status).
 - **Best Suited Technology: Inverse Reinforcement Learning (IRL)**
 - **Justification:**
 - In an E-learning system, the goal is to generate personalized learning paths based on the student's progress, preferences, and needs. Inverse Reinforcement Learning (IRL) is ideal here because it helps the AI model the underlying reward function that drives the optimal sequence of learning activities.
 - **IRL** can learn the best learning paths from the behavior of expert human educators or experienced learning systems, using this feedback to improve its recommendations. The AI will infer the reward

structure (e.g., knowledge acquisition, engagement) from expert demonstrations and apply it to create tailored learning paths for each student.

b) Consider a domestic robot. The main task is to cook.

- **Best Suited Technology:** Learning from Demonstration – Behavior Cloning (IL-BC)
- **Justification:**
 - In this case, **Behavior Cloning (IL-BC)** is a good choice because the robot needs to replicate a sequence of actions demonstrated by a human (e.g., chopping, stirring, heating). Through this method, the robot can observe and learn from expert demonstrations.
 - The system will map sensory inputs (e.g., visual, touch) to appropriate actions based on the human's behavior, allowing the robot to reproduce similar actions in cooking tasks. There's no need to infer complex reward structures, so behavior cloning directly translates expert actions into robot behavior effectively.

c) Consider the scenario of an intelligent robot that supports the elderly in a Nursing Home. The main task performed by the robot is the distribution of drugs to each user according to the prescribed doses, at the recommended times.

- **Best Suited Technology:** Reinforcement Learning (RL)
- **Justification:**
 - **Reinforcement Learning (RL)** is well-suited for this scenario because the robot needs to optimize its decision-making process based on a combination of temporal factors (e.g., time of day, medication schedule) and context (e.g., patient condition).
 - The robot will receive feedback (rewards or penalties) based on whether it delivers medication correctly, on time, and to the right individuals. Over time, the RL agent will learn the most effective policies for delivering medication under varying conditions, ensuring it adapts to any changes in the nursing home's schedule, patient needs, or errors in delivery.

d) Consider the scenario of an intelligent robot that supports the elderly in a Nursing Home. The main task is to entertain the elderly by telling them jokes and funny stories.

- **Best Suited Technology:** Learning from Demonstration – Inverse Reinforcement Learning (IL-IRL)
- **Justification:**
 - **Inverse Reinforcement Learning (IL-IRL)** is suitable for this scenario because the robot's task involves understanding the human feedback related to what types of jokes and stories are most appreciated by elderly users.
 - The robot can learn from expert entertainers (e.g., caregivers, comedians) by observing how they engage with the elderly. Through IRL, the robot can infer the reward structure (e.g., laughter, engagement) and replicate the behaviors of human entertainers.
 - It is crucial for the robot to not only reproduce jokes but also adapt to the preferences and emotional states of the elderly, which is effectively learned through the feedback loop of IRL.

QUESTION 6.40

Consider a Markov Decision Process (MDP).

a) List MDP main features.

- **States (S):** The set of all possible configurations or situations the agent can be in.
- **Actions (A):** The set of all possible actions the agent can take to transition between states.
- **Transition Function (T):** Describes the probability of transitioning from one state to another after taking a specific action. Formally, it is $T(s,a,s')=P(s'|s,a)$, where s' is the next state, s is the current state, and a is the action taken.
- **Reward Function (R):** The immediate reward received after transitioning from state s to state s' using action a . It is typically written as $R(s,a,s')$.

- **Discount Factor (γ):** A factor between 0 and 1 that represents the importance of future rewards relative to immediate rewards. A higher γ values future rewards more heavily.
- **Policy (π):** A strategy or mapping from states to actions. It defines the action the agent should take in each state.
- **Value Function (V):** The expected cumulative reward that an agent can achieve from a state s , following a particular policy. Denoted as $V(s)$, it calculates the long-term benefit of being in state s .

b) How is it formally defined?

- A Markov Decision Process is defined by the tuple (S, A, T, R, γ) , where:
- S is the set of states.
- A is the set of actions.
- $T(s, a, s')$ is the transition probability.
- $R(s, a, s')$ is the reward function.
- γ is the discount factor.
- Formally, the process operates with the Markov property, where the future state depends only on the current state and action, and not on the sequence of states that preceded it.

c) What's a policy and an optimal policy?

- **Policy (π):** A policy is a function or strategy that defines the action to take for each state in the MDP. A policy can be deterministic (always taking the same action in a state) or stochastic (choosing actions probabilistically).
- **Optimal Policy (π^*):** An optimal policy is a policy that maximizes the expected cumulative reward for the agent. It leads to the best possible long-term outcome in terms of rewards. The optimal policy can be determined by finding the action for each state that maximizes the expected value (reward) considering both immediate and future rewards.

d) What's the difference between Rewards and Utilities/Values?

- **Reward (R):** Rewards are the immediate feedback an agent receives after performing an action in a particular state. Rewards reflect the desirability of taking a particular action in a state and are given in the form of numerical values at each time step.
- **Utility/Value (V):** The value of a state represents the long-term benefit of being in that state, considering both immediate rewards and the future rewards the agent can expect to receive by following a particular policy. The value is the expected cumulative reward from a state under a given policy, discounted by the discount factor γ .
- In simple terms, the **reward** is immediate, while the **utility/value** is the expected cumulative reward over time.

e) What's the main goal of an MDP?

- The main goal of an MDP is to find the *optimal policy (π^*)* that maximizes the agent's expected cumulative reward over time. This involves determining the best sequence of actions that leads to the most favorable long-term outcome, given the dynamics of the environment and the rewards associated with different states.

f) Analyze and understand the impact on the optimal policy of considering different reward values in the states of the environment. E.g.:

a. Why is the agent heading straight into (2,4) from its surrounding states when the reward is -1.6?

- Even though the reward is negative, the agent may be heading to state (2,4) because it believes that future rewards from there (due to transitions and expected future states) might be more favorable. Alternatively, it could be trying to avoid an even worse state, and -1.6 is seen as a better option than other possible actions. The **discount factor** γ could be large, meaning future rewards

are highly valued, making this path a potential long-term strategy despite the immediate negative reward.

b. Why is the agent heading straight into the obstacle from (2,3), when $-0.0218 < \text{reward} < 0$?

- The agent might be heading into the obstacle because it expects that by going through it, it will eventually reach a state with higher rewards. The reward in the immediate state is negative, but the agent could be valuing future rewards more (due to the discount factor γ), making it willing to take short-term losses for long-term gains. Additionally, the agent could be mistakenly predicting that passing through the obstacle leads to a state with higher cumulative rewards.

c. Why does the agent avoid the terminals when reward > 0 ?

- If the reward is positive, the agent should typically seek to maximize cumulative reward. The reason for avoiding terminal states with reward > 0 could be that terminal states, despite having positive rewards, end the episode prematurely, meaning the agent cannot accumulate further rewards from subsequent actions. If the agent values ongoing rewards (and the ability to continue receiving rewards from non-terminal states), it may avoid terminating states, even if they provide positive rewards, in favor of continuing in non-terminal states with the potential for greater cumulative rewards.

QUESTION 6.41

What is Value Iteration? What's the difference between (i) the expected value of following policy Π in a state s ($V(\cdot)$), and (ii) the expected value of performing an action a in a state s , and then following policy Π (Q -value)?

• Value Iteration:

- Value Iteration is an algorithm used to compute the optimal value function $V^*(s)$ for all states in a Markov Decision Process (MDP). The algorithm updates the value of each state iteratively based on the expected rewards, considering all possible actions and outcomes. This process continues until the values converge.
- The update equation for Value Iteration is:

$$V_{k+1}(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} T(s, a, s') V_k(s') \right)$$

Where:

- $V_k(s)$ is the value of state s at iteration k ,
- γ is the discount factor,
- $R(s, a)$ is the reward for taking action a in state s ,
- $T(s, a, s')$ is the transition probability from state s to s' given action a .

- The process repeats until the value function converges. Once it converges, the optimal policy $\pi^*(s)$ is derived by selecting the action that maximizes the expected value:

$$\pi^*(s) = \arg \max_a \left(R(s, a) + \gamma \sum_{s'} T(s, a, s') V^*(s') \right)$$

• Difference Between $V(s)$ and $Q(s,a)$:

- $V(s)$ (Value Function):** Represents the expected long-term reward of being in state s and following a policy π . It gives the value of a state under a policy.

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, \pi \right]$$

- $Q(s,a)$ (Action-Value Function):** Represents the expected long-term reward of taking action a in state s , then following policy π . It gives the value of taking a specific action in a state.

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, a_0 = a, \pi \right]$$

• Key Differences:

- $V^\pi(s)$ gives the value of a state, considering all possible actions under the policy.
- $Q^\pi(s, a)$ gives the value of taking a specific action in a state and then following the policy.

• Relationship:

- The value function $V^\pi(s)$ can be derived from the action-value function $Q^\pi(s,a)$ by considering all possible actions:
 - $$V^\pi(s) = \sum_a \pi(a | s) Q^\pi(s, a)$$
- For the optimal policy π^* :
 - $$V^*(s) = \max_a Q^*(s, a)$$

WORKSHEET 7 - INTERPRETABLE AND EXPLAINABLE AI (XAI): INTERPRETABLE MODELS

QUESTION 6.1

a) What is an explanation? What does it mean to interpret? What are their differences? What is explainability/interpretability? What are their differences?

- **Explanation:** An explanation refers to the process of making something clear by describing its components, causes, or mechanisms. It involves providing reasons, context, or details to make something more understandable.
- **Interpretation:** Interpretation is the act of understanding and assigning meaning to something, often by analyzing or translating it into a more comprehensible form. In the context of machine learning, interpretation refers to understanding how a model arrives at its decision.
- **Differences:**
 - **Explanation** involves providing reasons or justifications, while **interpretation** involves understanding or interpreting what the model's behavior means.
 - An explanation gives insight into why something happens, whereas an interpretation gives insight into the meaning of the outputs.
- **Explainability/Interpretability:**
 - **Explainability** is the ability to explain the decisions or behavior of a model in understandable terms, often in a way that is meaningful to humans.
 - **Interpretability** refers to the extent to which a model's internal processes or decision-making can be understood by humans.
 - **Differences:**
 - Explainability is about providing clarity and reasons for decisions, while interpretability focuses more on how easy it is to understand a model's structure or behavior.

b) Why and when do we need XAI? Do we really need it? For what? Where is it critical? Provide examples of situations in which you need XAI. Do the same for situations in which you do not need XAI. Think of what makes you ask for explanations in some situations and not in others. Give examples of explanations.

- **When do we need XAI?**
 - We need **Explainable AI (XAI)** when the decision-making process of the AI needs to be understandable and transparent, especially in critical applications where human trust and safety are involved.
 - **Examples of situations where XAI is critical:**
 - **Healthcare:** A model predicting the likelihood of diseases must provide explanations to doctors to ensure the diagnosis is accurate and justifiable.
 - **Finance:** In credit scoring, customers need to understand why they were denied a loan.
 - **Autonomous vehicles:** In scenarios where the vehicle takes emergency actions, explaining why it made those decisions is important to ensure trust and safety.
- **When do we not need XAI?**
 - XAI is not as critical in applications where the decisions don't directly impact human safety or well-being, or where the system can be treated as a "black box" without consequences.
 - **Examples where XAI is less critical:**

- **Recommendation systems (e.g., Netflix, Amazon):** While explanations might improve user experience, they aren't crucial for making the system functional.
 - **Spam filters:** If a spam filter misclassifies an email, the consequences are generally minimal, and a full explanation may not be necessary.
- **Why do we ask for explanations in some situations and not in others?**
 - **When consequences are severe** (e.g., healthcare, finance, justice), humans want to understand and verify AI decisions.
 - **In low-stakes scenarios**, where the impact is minimal, users may accept decisions without needing an explanation.
- **Examples of explanations:**
 - **Healthcare:** "The model predicted that the patient has a high likelihood of diabetes because their blood sugar levels were above average, and they have a family history of the condition."
 - **Credit Scoring:** "The model denied the loan because the applicant's income was below the required threshold and their credit history shows multiple late payments."

c) **Concerning Machine Learning models, is accuracy enough? Is it always necessary to develop a specific model for providing explanations?**

- **Accuracy is not always enough.**
 - While accuracy is crucial for performance, it does not guarantee that a model's decision-making process is understandable or fair. For many applications, interpretability and fairness are just as important as accuracy, especially in domains like healthcare, law, or finance.
- **Is it necessary to develop a specific model for providing explanations?**
 - Not necessarily. There are techniques like **LIME (Local Interpretable Model-agnostic Explanations)** and **SHAP (Shapley Additive Explanations)** that can provide explanations for any model, even complex ones like deep neural networks, without requiring a separate model. However, simpler models like decision trees or linear regression are inherently more interpretable.

d) **Regarding the Machine Learning pipeline, explain the concept of XAI.**

- **XAI in the ML pipeline** involves integrating interpretability and explainability at each stage:
 - **Data Preprocessing:** Ensure data is clean, unbiased, and transparent.
 - **Model Training:** Choose models or algorithms that balance accuracy and interpretability (e.g., decision trees, linear models).
 - **Model Evaluation:** Use performance metrics as well as interpretability metrics to ensure the model can be understood by stakeholders.
 - **Model Deployment:** Provide tools to interpret model outputs, including visualizations and reasoning behind predictions.
 - **Post-Deployment:** Continuously monitor model decisions and adjust to ensure clarity in decision-making.

e) **Can you think of different categories of XAI? What are their differences?**

- **Post-Hoc Explanations:** These are provided after the model has made a prediction, typically using techniques like LIME, SHAP, or saliency maps. They attempt to explain the decisions of complex models.
- **Ante-Hoc Explanations:** These explanations are inherent in the model itself, as it is designed to be interpretable from the outset. Examples include decision trees, linear regression, and rule-based systems, where the decision process can be directly understood.
- **Global vs Local Explanations:**
 - **Global Explanations** give an overview of the model's behavior and how it operates across all inputs.
 - **Local Explanations** provide insight into individual predictions, explaining how specific inputs lead to the output.

f) **What are the challenges of XAI?**

- **Complexity of Models:** Advanced models like deep learning are inherently difficult to explain, and providing meaningful explanations for such models remains a significant challenge.
- **Lack of Standardization:** There is no universal framework for providing explanations, making it difficult to compare or evaluate different XAI approaches.
- **Trade-Off Between Interpretability and Accuracy:** There is often a trade-off where more accurate models (like deep learning) are harder to interpret, while simpler models (like decision trees) are more interpretable but might sacrifice accuracy.
- **Ensuring Fairness and Bias:** Explaining why a model made a certain decision can uncover biases, and explaining them in a fair, unbiased manner is difficult.
- **User Trust:** Even with explanations, convincing users to trust AI systems can be difficult if they are skeptical of the technology.

g) What are the challenging problem areas?

- **High-Complexity Models:** Deep neural networks and ensemble models present significant challenges for explainability due to their "black-box" nature.
- **Real-Time Explanations:** Providing real-time explanations in fast-paced applications (e.g., autonomous driving, medical diagnostics) while maintaining accuracy and interpretability is a challenge.
- **Interpreting Probabilistic Outputs:** Models that output probabilities or distributions (like Bayesian models) are more difficult to explain intuitively compared to discrete outputs.

h) Provide an illustrative example of an XAI System.

- **Example: Healthcare Diagnostic System**
 - A deep learning model is used to predict whether a patient has a certain disease based on medical imaging. Using **SHAP** values, the model provides a visual explanation showing which parts of the image were most influential in the prediction, such as highlighting areas that resemble known signs of the disease. This explanation is valuable for medical professionals to trust the AI's decision and understand the rationale behind the diagnosis, improving its adoption in critical healthcare applications.

QUESTION 6.3

Consider the model-agnostic methods for XAI and answer the following questions.

a) What exactly are model-agnostic methods? Describe their main features and differences among them.

- **Model-Agnostic Methods** are techniques used in Explainable AI (XAI) that can be applied to any machine learning model, regardless of its internal structure. These methods do not depend on the specific architecture or type of model (e.g., deep learning, decision trees, etc.), making them versatile tools for explaining the behavior of complex, opaque models. The primary goal of these methods is to offer insight into how models make decisions, regardless of whether the model is interpretable by design.
- **Main Features:**
 - **Flexibility:** Model-agnostic methods can be applied to any machine learning model (e.g., decision trees, neural networks, ensemble models) without requiring any changes to the model itself.
 - **Post-Hoc Explanation:** These methods are typically used after the model has been trained, providing explanations of its predictions or behavior.
 - **Black-Box Explanation:** They help interpret black-box models (i.e., models whose decision-making process is not transparent or easily understood).
 - **Local and Global Interpretability:** Some methods focus on explaining individual predictions (local), while others try to explain the model's overall behavior (global).
- **Main Types of Model-Agnostic Methods:**
 - **LIME:** Explains individual predictions by approximating the complex model with simpler, interpretable models locally.

- **SHAP:** Provides both local and global explanations by quantifying the contribution of each feature using Shapley values from game theory.
- **PDP:** Shows how predicted outcomes change as a feature varies, providing global insights.
- **ICE:** Like PDP but provides insights for individual predictions.
- **Counterfactual Explanations:** Explains what changes to input features would alter the prediction, offering actionable insights.

b) When are model-agnostic methods relevant?

- Interpreting black-box models like deep learning or random forests.
- Providing transparency in critical applications like healthcare, finance, and criminal justice.
- Fulfilling regulatory and ethical requirements for model explainability.
- Enhancing trust in AI systems by explaining model decisions after deployment.
- These methods are crucial when models are complex and non-transparent but need to be explained to users, especially in high-stakes or regulated environments.

QUESTION 6.4

A **surrogate model** is a simple model that is used to explain a complex model. Surrogate models are usually created by training a linear regression or decision tree on the original inputs and predictions of a complex model. Coefficients, variable importance, trends, and interactions displayed in the surrogate model are then assumed to be indicative of the internal mechanisms of the complex model.

a) What is the scope of interpretability for surrogate models? Global or local?

- **Surrogate models** typically provide **local interpretability**, especially when used to explain the behavior of a specific instance or prediction of a complex model. However, depending on the model's application, surrogate models can also offer **global interpretability** by approximating the entire model's decision-making process across all inputs. For example, a decision tree used as a surrogate for a neural network can explain the overall behavior of the network globally, while a local explanation might focus on a specific prediction made by the model.

b) What complexity of functions can surrogate models help explain?

- Surrogate models can help explain **complex, non-linear functions** by approximating them with simpler models, like linear regression or decision trees. These simpler models can capture high-level trends, interactions, and feature importance, which makes the original complex model more interpretable. For example, a neural network or ensemble model can be approximated using a decision tree, making it easier to understand the relationships between features and predictions.

c) How do surrogate models enhance understanding?

- Surrogate models enhance understanding by providing an interpretable representation of the complex model's decision-making process. The simpler surrogate model can highlight key factors, relationships, and feature importance in a way that humans can more easily understand, making the complex model's behavior more transparent. By analyzing the surrogate's coefficients, feature importance, and decision boundaries, users can gain insights into how the complex model arrives at its predictions.

d) How do surrogate models enhance trust?

- Surrogate models build trust by offering **transparency** and **explainability**. When users can understand how a model is making decisions, they are more likely to trust it, especially in high-stakes applications. By creating simpler models that approximate the complex model's behavior, surrogate models allow users to see and validate the model's decision-making process, increasing their confidence in its reliability and fairness.

e) What are the main differences between LIME and SHAP?

- **LIME (Local Interpretable Model-agnostic Explanations)** and **SHAP (Shapley Additive Explanations)** are both model-agnostic methods used to explain machine learning predictions, but they have key differences:
 - **LIME** focuses on **local explanations** by approximating the complex model with a simple, interpretable model (such as a linear regression) around a specific instance. It uses perturbed data points and trains a surrogate model locally to explain a particular prediction.
 - **SHAP**, based on **Shapley values** from game theory, provides both **local and global explanations**. It assigns each feature a value representing its contribution to the model's prediction by considering all possible feature permutations. SHAP values provide consistent and fair attributions for feature importance, making it more rigorous than LIME in ensuring the accuracy and fairness of the explanations.
- **Key Differences:**
 - **LIME** focuses on creating a surrogate model to approximate the complex model's behavior around a single instance, while **SHAP** provides a consistent and mathematically grounded way to quantify the contribution of each feature across all predictions.
 - **LIME** uses a simpler, local approximation (like a linear model) for each individual prediction, while **SHAP** gives a global measure of feature importance that is consistent across predictions.
 - **LIME** can sometimes be less consistent, as it depends on the local surrogate's performance, whereas **SHAP** is considered more consistent and mathematically well-defined.

QUESTION 6.5

Consider other methods for XAI.

- a) **Give examples of situations where humans rely on previous examples for making decisions.**
 - **Medical Diagnosis:** Doctors use previous cases and patient histories to diagnose new patients.
 - **Judicial Decision-Making:** Judges refer to case law for similar legal cases.
 - **Customer Support:** Support agents use past interactions to resolve new issues.
 - **Investment Decisions:** Investors analyze past market data for future investment choices.
 - **Engineering Problem-Solving:** Engineers apply solutions from past projects to current problems.
- b) **Which other methods do you know for XAI? Discuss their advantages and their limitations.**
 - **Counterfactual Explanations:**
 - **Advantages:** Provides actionable insights by showing what changes would lead to a different outcome.
 - **Limitations:** May be difficult to generate realistic counterfactuals and oversimplify complex models.
 - **Feature Importance:**
 - **Advantages:** Simple, intuitive, and provides global insights into feature influence.
 - **Limitations:** May overlook feature interactions and mislead complex models.
 - **Decision Trees as Surrogate Models:**
 - **Advantages:** Offers global explainability and is easy to visualize.
 - **Limitations:** Can oversimplify complex models and may overfit.
 - **Attention Mechanisms:**
 - **Advantages:** Highlights relevant input data and improves interpretability, especially for NLP models.
 - **Limitations:** May not always reflect the full decision-making process and can be complex to interpret.
 - **Example-Based Explanations:**
 - **Advantages:** Provides concrete examples, making reasoning easier for non-experts.
 - **Limitations:** May not have close examples for every instance and can be memory intensive.
 - **Rule-Based Systems:**

- **Advantages:** Transparent and flexible, with clear rules to explain decisions.
- **Limitations:** Can oversimplify complex problems and is hard to scale for larger models.

WORKSHEET 8 - TRUSTWORTHY AND RESPONSIBLE AI

QUESTION 2.1

Consider the scenario of Jarbas, an intelligent robot that supports the elderly in a Nursing Home. Jarbas is intended to follow the Ethics Guidelines for Trustworthy AI.

- The main task performed by Jarbas is the distribution of drugs to each user according to the prescribed doses, at the recommended times (Distribution task). Give examples of acts or decisions that Jarbas must not take if it is to conform with the three components of Trustworthy AI (one example for each component).
 - **Accountability:**
 - **Example:** Jarbas must not autonomously change a patient's prescribed medication or dosage. If the robot makes an error, it must be able to provide a transparent explanation of how the decision was made and who is responsible for it. This ensures accountability and prevents errors from being untraceable or unaddressed.
 - **Transparency:**
 - **Example:** Jarbas must not make decisions about medication without providing clear reasoning or communicating the process to caregivers or patients. For example, Jarbas should not decide that a medication can be skipped without explaining to the caregiver or patient the reasoning behind this decision.
 - **Fairness:**
 - **Example:** Jarbas must not discriminate in its drug distribution based on irrelevant personal characteristics such as a patient's race, gender, or age. Every patient should receive the prescribed medication at the correct time, ensuring equitable treatment for all patients.
- In the intervals of its main task, Jarbas entertains the users by telling them jokes and funny stories (Entertainment task). The robot knows about each user's profile and is quite competent in selecting the most convenient gags for each one of them. Refer to the afore-mentioned Ethics Guidelines and discuss the following statement.
 - Compliance with the principle of fairness is key in the entertainment task. Why?
 - **Why it's important:** The robot should ensure that it does not favor or discriminate against certain individuals based on biases. For example, Jarbas should not tell jokes that could perpetuate stereotypes or marginalize specific groups of people. Fairness ensures that all individuals are treated equally, irrespective of their background or personal characteristics. The robot must avoid reinforcing harmful stereotypes and select entertainment that is inclusive and sensitive to everyone's needs.
 - Compliance with the principle of prevention of harm is key in the distribution task. Why?
 - **Why it's important:** In the distribution of drugs, Jarbas must ensure that the right medication is delivered in the correct dosage to avoid any adverse health consequences. Any action that could potentially harm a patient, like administering the wrong drug or giving medication at the wrong time, would violate the principle of prevention of harm. This principle is essential in healthcare to ensure the robot's actions do not put patients at risk.
 - There may be a tension between the principles of prevention of harm and respect for human autonomy in the distribution task. Why?
 - **Why there may be tension:** There can be a conflict when the robot is required to act in a way that respects a patient's autonomy (e.g., allowing them to make their own decision about taking medication) while also preventing harm (e.g., ensuring that the patient follows the prescribed

medication schedule). For instance, if a patient refuses medication, respecting their autonomy means not forcing it on them, but preventing harm may require that they take the medication for their health. Balancing both principles requires careful consideration of the patient's well-being and their right to make decisions about their own health.

QUESTION 2.2

Consider the following dilemma: an autonomous car sees a child crossing the road in front of it and realizes it is impossible to stop the car in due time; there are two alternatives: hit the child, running the risk of killing them, or swerve to a ravine off the road, running the risk of killing the occupant of the vehicle.

a) How should the autonomous car act?

- The decision that the autonomous car should make depends on the ethical framework used to program it. Some of the main approaches are:
 - **Utilitarian Approach:** The car should choose the action that results in the least overall harm. This would likely involve swerving to avoid hitting the child, even if it risks harming the car's occupant. This is because the utilitarian framework seeks to minimize total harm, and the child's death (a potential future life lost) may be seen as a greater harm than the risk to the occupant.
 - **Deontological Approach:** The car may be programmed to follow rules about protecting human life, regardless of the consequences. In this case, the car may prioritize saving the occupant, as it has a duty to protect the person inside the vehicle.
 - **Virtue Ethics Approach:** The car's decision might be made according to virtues such as compassion or wisdom, which would aim for a balance of protecting both lives without over-prioritizing one over the other.
- The decision is ethically complex, and there isn't a universally accepted solution. The car could act in any of these ways depending on the ethical principles embedded in its programming.

b) If the car kills someone in such circumstances, who should be responsible for the tragic event?

- **Owner of the Car:** The owner could be held accountable if the car's failure to act was due to improper use, lack of maintenance, or modification of the system.
- **Maker of the Car:** The car manufacturer might be responsible if there was a design flaw in the vehicle or if it failed to meet safety standards that could have prevented the accident.
- **Maker of the AI Software:** If the AI software is designed poorly or doesn't make ethical decisions according to accepted standards, the developer of the AI system could be responsible. This might include decisions made by the AI that are not in line with accepted moral frameworks.
- **Programmer of the AI Software:** The individual programmer might also bear responsibility if they made specific programming decisions that led to the car's inability to choose between the two harmful alternatives effectively.
- **Other Entities:** There may be shared responsibility among the manufacturer, software developers, or other stakeholders (such as regulatory bodies) if the accident results from systemic failures, such as lack of regulation or insufficient testing.
- Overall, the issue of accountability is highly complex and may involve multiple parties, including the owner, developers, and regulatory bodies.

c) If you were to buy an autonomous car, would you prefer to buy one programmed to give priority, in a dilemma situation as the one described above, to save occupants' lives, or to save pedestrians' lives?

- This is a highly personal decision, and it depends on individual values and priorities:
 - **Save Occupant's Lives:** Some may prefer that the car prioritizes saving the lives of the people inside the vehicle, arguing that individuals should have the right to be protected in the vehicle they have purchased.
 - **Save Pedestrians' Lives:** Others may prefer that the car prioritizes saving pedestrians, viewing the protection of innocent bystanders as more morally pressing.

- In general, societal preferences may lean towards minimizing harm to the public (pedestrians), but personal preferences could vary based on a sense of responsibility for the people inside the vehicle.

QUESTION 2.3

An European taxi company decided to install cameras in their cars for capturing the image of the passenger seats. These images are analyzed by an AI system to increase safety both for passengers and taxi drivers. Consider that the system has powerful image recognition capabilities that allow both the identification of common objects and common human gestures and body postures. Suppose that the taxi company wants the AI software to follow the Ethics Guidelines for Trustworthy AI and answer the following questions:

a) Example of a situation where the system could act to ensure the safety of the driver without violating the guidelines:

- One example of a situation where the AI system could act to ensure the safety of the driver, while following the Ethics Guidelines for Trustworthy AI, could be:
 - **Aggressive passenger behavior:** If the system detects that a passenger is exhibiting threatening behavior (such as violent gestures or body posture, e.g., raising their fists or moving aggressively towards the driver), it could trigger an alert to the driver or dispatch authorities for intervention. The system could also lock the doors to prevent the passenger from reaching the driver. However, it would need to be programmed to avoid overreaching or misinterpreting non-threatening gestures, ensuring that innocent actions are not flagged as dangerous.
- This action would align with the principle of **Prevention of Harm**, as it protects the driver from potentially dangerous situations, without overstepping ethical boundaries.

b) Examples of acts or decisions that the AI system must not take to conform with the three components of Trustworthy AI:

- **Fairness (Non-Discrimination):**
 - **What the system must not do:** The AI should not make biased judgments about passengers based on their appearance, race, gender, or other personal characteristics. For instance, it should not flag a person as suspicious simply because they are from a particular demographic group.
 - **Why it violates fairness:** Discriminating based on irrelevant personal traits violates the principle of fairness and could lead to unjust treatment of certain individuals or groups.
- **Accountability (Transparency):**
 - **What the system must not do:** The AI should not take actions or make decisions that are not explainable. For example, if the AI locks the doors based on detecting a specific gesture, it must be clear why that action was taken. The company should be able to provide an explanation if needed.
 - **Why it violates accountability:** If the AI's actions are not transparent and there is no way to audit its decisions, it would violate the requirement for accountability.
- **Privacy:**
 - **What the system must not do:** The system should not record or store personal images or videos of passengers unless necessary and explicitly authorized. For example, it should not store data on passengers' body language or other sensitive attributes unless relevant to safety and with the passenger's consent.
 - **Why it violates privacy:** Collecting excessive personal data or failing to inform passengers about how their data is used could violate privacy principles.

c) Classification of the following sentences as True or False:

- The compliance with the principle of explicability is an ethical imperative for this system. Why?
 - **True.** Explicability, or explainability, is a fundamental aspect of AI ethics, especially for systems that make decisions that impact people's safety. The passengers and the taxi company should be able to understand why the system is taking certain actions (e.g., why the AI locks the doors or

sends an alert). This ensures trust and accountability in the AI's operation and conforms with the principle of **Transparency**.

b. There may be tension between the principle of prevention of harm and the freedom of business. Why?

- **True.** There could be a conflict between ensuring safety (e.g., by monitoring passengers for risky behavior) and the freedom of business (e.g., a company seeking to maximize profits with minimal interference). For instance, a company might prioritize profitability over stringent safety protocols, which could result in the AI system being underutilized or poorly monitored. The **Prevention of Harm** principle might require strict measures (e.g., monitoring all passengers), while the business might resist those measures due to cost or customer experience considerations.

d) Classification of this system according to the AI Act's risk-based approach:

- The **AI Act's risk-based approach** classifies AI systems based on their potential risk to human rights and safety. This system would likely fall under the **High-Risk** category, since it deals with safety-critical tasks (e.g., ensuring the safety of both the driver and passengers) and could have significant consequences if it fails or is misused.
 - **Risk Classification:** High-Risk AI
 - **Reasoning:** The system directly impacts the safety and privacy of individuals. Misuse, malfunction, or biased decision-making could lead to safety incidents or violations of privacy. Therefore, it would need to comply with strict requirements such as transparency, accountability, and data protection standards as outlined in the AI Act.

QUESTION 2.4

Make the descriptions correspond to one of the following, according to Dignum [2019]'s perspective:

- **Ethics by Design**
- **Ethics in Design**
- **Ethics for Design(ers)**

a) Developing an AI system while taking values into consideration, using ethical theories to define behaviors from such values, and prioritizing them.

- This corresponds to **Ethics by Design**. It emphasizes incorporating ethical values directly into the design process, ensuring that ethical considerations shape the behavior of the system from the outset.

b) Developing an AI system while making sure that all the processes are conducted in a responsible manner, with a clear chain of responsibilities, assuring that the data and the processes are transparent, and adopting design methods that ensure accountability.

- This corresponds to **Ethics in Design**. It focuses on ensuring that ethical principles guide the process of designing AI systems, such as maintaining responsibility, transparency, and accountability throughout the system's lifecycle.

QUESTION 2.5

Consider chatbots based on Large Language Models (LLMs), like ChatGPT and answer the following questions.

a) How would they be classified according to the AI Act risk-based approach? To what extent does it comply with consequent requirements?

- **Classification:** Chatbots based on Large Language Models (LLMs), like ChatGPT, would likely be classified as **high-risk** AI systems under the AI Act. This is because they have the potential to significantly affect individuals or society due to their ability to generate and interpret human-like text. They can influence public opinion, assist in decision-making, and provide information, making their misuse or malfunction highly impactful.
- **Compliance with Requirements:** High-risk AI systems are subject to strict requirements, including:

- **Transparency:** LLM-based chatbots should clearly inform users when they are interacting with an AI system.
 - **Data Governance:** Ensuring high-quality, non-biased, and representative training data.
 - **Accountability:** Clear accountability mechanisms should be established for the outputs generated by the LLM.
 - **Monitoring and Auditing:** Regular audits to assess and monitor the system's performance and prevent misuse.
 - **Human Oversight:** Humans should be able to intervene if the chatbot produces harmful or biased content.
- While LLM-based systems like ChatGPT aim to follow many ethical principles and practices, full compliance with the AI Act would require specific measures, such as ensuring transparency in how data is used, providing an audit trail of outputs, and preventing any harmful, biased, or misleading information from being generated.

b) What other risks are associated with using LLMs?

- **Bias and Discrimination:** LLMs are trained on large datasets that may contain biased or discriminatory information. This can result in biased outputs that perpetuate harmful stereotypes or inequality, especially in sensitive contexts like healthcare, law enforcement, and hiring.
- **Misinformation and Harmful Content:** LLMs can generate convincing yet false or misleading information, potentially contributing to the spread of misinformation. This is especially concerning in areas like public health, politics, and finance, where incorrect information can have serious consequences.
- **Lack of Transparency:** While LLMs generate outputs based on complex patterns in data, understanding the reasoning behind specific outputs can be difficult. This lack of transparency can hinder users' ability to trust the system, and in cases of errors or harmful outcomes, it may be challenging to determine accountability.
- **Privacy Concerns:** LLMs can sometimes inadvertently memorize and reproduce sensitive information seen during training, potentially violating privacy regulations. For instance, if an LLM has been trained on data containing personal details, it might generate content that inadvertently discloses private information.
- **Manipulation and Exploitation:** Chatbots could be exploited for malicious purposes, such as phishing, scams, or manipulative marketing strategies, if not properly monitored and regulated.
- **Dependence on AI:** Over-reliance on LLM-based chatbots may reduce human decision-making and critical thinking skills, especially in contexts where human judgment is crucial, like legal advice or medical diagnosis.

WORKSHEET 9 - RECOMMENDER SYSTEMS (RECSYS)

QUESTION 9.1

Define and give some examples of Recommender Systems

- A **Recommender System** is an algorithm or set of algorithms that provide personalized suggestions to users based on various data inputs. These systems predict and recommend items (such as products, movies, music, articles, etc.) to users by analyzing their preferences, behaviors, and the behaviors of similar users.
- **Examples of Recommender Systems:**
 - **E-commerce platforms:** Amazon recommends products based on a user's previous purchases, search history, and preferences.
 - **Streaming services:** Netflix or Spotify recommends movies, TV shows, or music based on user viewing/listening history.
 - **Social media:** Facebook or Instagram suggests content (pages, posts, or ads) based on user interactions and engagement.
 - **Online news:** News websites suggest articles to read based on user interests and past reading behavior.

QUESTION 9.2

What are the benefits of Recommender Systems?

- **Personalization:** Recommender systems tailor suggestions to individual user preferences, increasing satisfaction and engagement by providing more relevant content.
- **Increased User Engagement:** By offering suggestions that align with a user's interests, recommender systems encourage users to interact more with platforms, increasing time spent on the platform.
- **Improved Decision Making:** Recommender systems help users make decisions more quickly by narrowing down options from a large pool of content, improving the user experience.
- **Revenue Generation:** By promoting products or services that users are more likely to buy or engage with, recommender systems can increase sales, subscriptions, or ad revenue.
- **Discovery:** They help users discover new content, products, or services they might not have encountered otherwise, enriching the user experience and expanding the platform's offerings.
- **Scalability:** Recommender systems allow platforms to handle large-scale catalogs of items, making it easier to manage and recommend relevant content from millions of options.

QUESTION 9.3

What are the main categories of Recommender Systems? Present their differences, Pros, and Cons.

- **Collaborative Filtering:**
 - **Description:** Recommends items based on user behavior or interactions.
 - **Pros:** Easy to implement, effective for large catalogs.
 - **Cons:** Suffers from cold start problems and scalability issues.
- **Content-Based Filtering:**
 - **Description:** Recommends items based on item features and past user interactions.
 - **Pros:** No cold start for users, transparent recommendations.
 - **Cons:** Limited discovery of new items requires detailed metadata.
- **Hybrid Systems:**
 - **Description:** Combines collaborative and content-based filtering.
 - **Pros:** Overcomes the limitations of both approaches, improves recommendation quality.
 - **Cons:** More complex, computationally intensive.
- **Knowledge-Based Systems:**
 - **Description:** Uses explicit knowledge about user preferences and item characteristics.
 - **Pros:** Accurate recommendations for specific needs, no large dataset required.
 - **Cons:** Requires detailed input from users, less adaptable.
- **Demographic-Based Systems:**
 - **Description:** Recommends items based on user demographics (e.g., age, gender).
 - **Pros:** Simple to implement.
 - **Cons:** Assumes too much similarity within demographic groups, less personalized.
- Hybrid systems are often preferred as they combine the strengths of multiple methods to improve recommendation quality.

QUESTION 9.4

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

- a) Compute the similarity between Alice and each one of the other users, using Pearson Correlation:

$$sim(a, b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2}}$$

Solution

We start by computing the average of ratings of each user:

$$\bar{r}_{Alice} = (5 + 3 + 4 + 4)/4 = 4$$

$$\bar{r}_{U1} = (3 + 1 + 2 + 3)/4 = 2.4$$

$$\bar{r}_{U2} = (4 + 3 + 4 + 3)/4 = 3.5$$

$$\bar{r}_{U3} = (3 + 3 + 1 + 5)/4 = 3$$

$$\bar{r}_{U4} = (1 + 5 + 5 + 2)/4 = 3.25$$

Now we can use this average ratings to compute the Pearson Correlation between two users. In this case, we are going to compute the similarity between Alice (denote for short by A) and each one of the other users:

$$sim(A, U1) = \frac{\sum_{p \in P} (r_{A,p} - \bar{r}_A)(r_{U1,p} - \bar{r}_{U1})}{\sqrt{\sum_{p \in P} (r_{A,p} - \bar{r}_A)^2} \sqrt{\sum_{p \in P} (r_{U1,p} - \bar{r}_{U1})^2}}$$

$$= \frac{(5-4)(3-2.4)+(3-4)(1-2.4)+(4-4)(2-2.4)+(4-4)(3-2.4)}{\sqrt{(5-4)^2+(3-4)^2+(4-4)^2+(4-4)^2} \sqrt{(3-2.4)^2+(1-2.4)^2+(2-2.4)^2+(3-2.4)^2}} = 0.85$$

Using a similar procedure, we can compute the rest of the similarities:

- $sim(A, U2) = 0.70$
- $sim(A, U3) = 0$
- $sim(A, U4) = -0.79$

b) What is the neighbor set of Alice if we confine the set to the two most similar ones?

- From the previous computations, we can restrict the neighbor set to: $N_{Alice} = \{U1, U2\}$

c) Generate a prediction for the rating of Alice for item 5, based on the neighbor's ratings and using the prediction function:

$$Pred(\bar{r}_{a,p}) = \bar{r}_a + K \sum_{b \in N_a} sim(a, b)(r_{b,p} - \bar{r}_b)$$

where

$$k = \frac{1}{\sum_{b \in N_a} sim(a, b)}$$

Solution

- $Pred(\bar{r}_{a,p}) = 4 + \frac{1}{0.85+0.70} [0.85(3 - 2.4) + 0.70(5 - 3.8)] = 4.87$