
Lesson 01.1: AI and Machine Learning



MSML610: Advanced Machine Learning

Lesson 01.1: AI and Machine Learning

Instructor: Dr. GP Saggese - gsaggese@umd.edu

References: - AIMA (Artificial Intelligence: a Modern Approach), Chap 1



1 / 19

2 / 19: ML, AI, and Intelligence

ML, AI, and Intelligence



- Machine Learning is a subset of AI
 - All of it confused with deep learning, large-language models, predictive analytics, ...
- What is artificial intelligence?
- What is intelligence?



2 / 19

- **ML, AI, and Intelligence**
 - **Machine Learning is a subset of AI**
 - * Machine Learning (ML) is a part of Artificial Intelligence (AI). Think of AI as the broad field that aims to create machines capable of performing tasks that would require human intelligence. ML is a specific approach within AI that focuses on teaching machines to learn from data and improve over time without being explicitly programmed.
 - * It's common for people to confuse ML with other terms like deep learning, large-language models, and predictive analytics. Deep learning is a more advanced form of ML that uses neural networks with many layers to analyze complex data. Large-language models, like GPT, are a type of deep learning model designed to understand and generate human language. Predictive analytics involves using data, statistical algorithms, and ML techniques to identify the likelihood of future outcomes based on historical data.
- **What is artificial intelligence?**
 - Artificial Intelligence is the science and engineering of creating intelligent machines that can perform tasks typically requiring human intelligence. These tasks include understanding natural language, recognizing patterns, solving problems, and making decisions. AI systems can be rule-based, where they follow specific instructions, or they can learn from data, as in the case of ML.
- **What is intelligence?**
 - Intelligence, in a general sense, refers to the ability to learn, understand, and apply knowledge to adapt to new situations and solve problems. In humans, intelligence encompasses a range of cognitive abilities, including reasoning, problem-solving, planning,

abstract thinking, and learning from experience. In the context of AI, intelligence is about creating systems that can mimic these human cognitive functions to some extent.

3 / 19: Human Intelligence

Human Intelligence



- We call ourselves “homo sapiens” because intelligence sets us apart from animals
- For thousands of years, we tried to understand how we think
- One of the biggest mysteries
 - Brain is a small mass of matter
 - Our brain can understand nature secrets, e.g., theory of relativity, quantum mechanics, black holes in the universe
 - How can brain understand, predict, and manipulate a world more complicated than itself?



3 / 19

- **Human Intelligence**
 - **We call ourselves “homo sapiens” because intelligence sets us apart from animals**
 - * Humans are known as *homo sapiens*, which translates to “wise man” in Latin. This name highlights the importance of intelligence in distinguishing humans from other species. Unlike animals, humans have developed complex languages, cultures, and technologies, all of which are rooted in our cognitive abilities.
 - **For thousands of years, we tried to understand how we think**
 - * Throughout history, philosophers, scientists, and thinkers have been fascinated by the workings of the human mind. Understanding how we think and process information has been a central question in fields like philosophy, psychology, and neuroscience. This quest has led to various theories and models attempting to explain human cognition.
 - **One of the biggest mysteries**
 - * **Brain is a small mass of matter**
 - The human brain, despite its relatively small size and weight, is an incredibly complex organ. It consists of billions of neurons and connections, allowing it to perform a vast array of functions.
 - * **Our brain can understand nature secrets, e.g., theory of relativity, quantum mechanics, black holes in the universe**
 - The brain’s ability to comprehend and formulate complex scientific theories, such as Einstein’s theory of relativity or the principles of quantum mechanics, is remarkable. These theories help us understand the universe’s fundamental workings, showcasing the brain’s extraordinary capacity for abstract thought

and problem-solving.

* **How can brain understand, predict, and manipulate a world more complicated than itself?**

- This question highlights the paradox of human intelligence: how can a biological organ, limited in size and resources, grasp and manipulate concepts and phenomena that are vastly more complex than itself? This remains one of the most intriguing puzzles in the study of human cognition and intelligence.

4 / 19: Artificial Intelligence

Artificial Intelligence



- The term “Artificial Intelligence” was coined in 1956
- **AI aims to:**
 - Understand human intelligence
 - Create intelligent entities
 - “*What I cannot create, I do not understand*” (Feynman, 1988)
- **AI is a technology**
 - Universal and applicable to any human activity and task
 - Its impact greater than any previous historical event
 - Currently generates trillions of dollars annually in revenue
 - Presents many unresolved problems
 - E.g., major concepts in physics might be established



4 / 19

- **Artificial Intelligence**
 - The term “Artificial Intelligence” was first introduced in 1956 during a conference at Dartmouth College. This marked the beginning of AI as a field of study. The goal was to explore the possibility of creating machines that could simulate aspects of human intelligence.
- **AI aims to:**
 - **Understand human intelligence:** AI research seeks to understand how human intelligence works by trying to replicate it in machines. This involves studying cognitive processes and how they can be modeled computationally.
 - **Create intelligent entities:** The ultimate goal of AI is to create machines that can perform tasks that typically require human intelligence, such as reasoning, learning, and problem-solving.
 - The quote “*What I cannot create, I do not understand*” by Richard Feynman emphasizes the idea that to truly understand something, one must be able to recreate it. This is a guiding principle in AI research.
- **AI is a technology**
 - AI is a versatile technology that can be applied to virtually any human activity or task, from healthcare to finance to entertainment. Its potential applications are vast and varied.
 - The impact of AI is considered to be greater than any previous technological advancement in history, as it has the potential to transform industries and societies on a global scale.
 - AI is already a significant economic force, generating trillions of dollars in revenue each year. This highlights its importance and the rapid pace of its development.
 - Despite its advancements, AI still faces many unresolved challenges. For example, while

AI can solve specific problems, understanding and replicating the full scope of human intelligence remains a complex task. Additionally, AI could potentially contribute to solving major scientific problems, such as those in physics, by providing new insights and approaches.

AI Formal Definition

- AI is defined around **two axes**:
 - Thinking vs. Acting
 - Human vs. Rational (ideal performance)
- Four possible definitions of AI as a machine that can:
 1. Think humanly
 2. Think rationally
 3. Act humanly
 4. Act rationally
- **Q**: Which one do you think is the best definition?
- We will see that building machines that can **"act rationally"** should be ultimate goal of AI



5 / 19

- **AI Formal Definition**
 - AI is defined around **two axes**:
 - * **Thinking vs. Acting**: This axis differentiates between AI systems that focus on cognitive processes (thinking) and those that focus on behavior (acting).
 - * **Human vs. Rational (ideal performance)**: This axis distinguishes between AI systems that aim to mimic human behavior and those that aim for optimal, logical performance, regardless of human-like characteristics.
 - Four possible definitions of AI as a machine that can:
 1. **Think humanly**: This involves creating machines that replicate human thought processes. It focuses on understanding and modeling how humans think.
 2. **Think rationally**: This involves machines that use logic and reasoning to solve problems, aiming for the most logical outcome.
 3. **Act humanly**: This involves machines that behave like humans, often evaluated through the Turing Test, which assesses a machine's ability to exhibit human-like behavior.
 4. **Act rationally**: This involves machines that make decisions and take actions that are logically optimal, aiming for the best possible outcome.
 - **Q**: Which one do you think is the best definition?
 - * This question encourages reflection on the goals of AI development and which approach aligns best with those goals.
 - We will see that building machines that can **"act rationally"** should be the ultimate goal of AI
 - * The emphasis here is on creating AI systems that make the best decisions based on available information, prioritizing logical and optimal outcomes over mimicking

human behavior. This approach is often seen as the most practical and beneficial in real-world applications.

1. AI as Thinking Humanly

- To build machines that think like humans we need to **determine how humans think**
- **Pros**
 - Express precise theory of the human mind as a computer program
- **Cons**
 - Unknown workings of the human mind
 - Anthropocentric definition



6 / 19

- **AI as Thinking Humanly**
 - The idea here is to create machines that can think in the same way humans do. This involves understanding the processes and mechanisms behind human thought. Essentially, it's about mimicking human cognitive processes in machines.
- **Pros**
 - By trying to replicate human thinking in machines, we can develop a precise theory of the human mind. This theory can be expressed as a computer program, which helps in understanding and simulating human cognition. It allows us to create models that can potentially predict human behavior and decision-making.
- **Cons**
 - One major challenge is that we don't fully understand how the human mind works. This lack of understanding makes it difficult to replicate human thought processes accurately in machines. Additionally, this approach is *anthropocentric*, meaning it is centered around human characteristics and may not consider other forms of intelligence or ways of thinking. This can limit the scope and applicability of AI systems.

2. AI as Thinking Rationally



- What are the rules of **correct thinking**?
 - Given correct premises, yield correct conclusions
- **Logic** studies the “laws of thought”
 - Formalize statements about objects and their relations
- **Automatic theorem proving**
 - Programs solve problems in logical notation
 - Run indefinitely if no solution exists (related to the halting problem)



7 / 19

- 2. AI as Thinking Rationally
- What are the rules of **correct thinking**?
 - The idea here is to understand how AI can be designed to think in a way that is logically sound. When we talk about *correct thinking*, we mean that if an AI system is given true information (premises), it should be able to process this information and come up with true conclusions. This is similar to how humans use logic to solve problems and make decisions.
- **Logic** studies the “laws of thought”
 - Logic is a branch of philosophy and mathematics that deals with the principles of valid reasoning. It involves creating formal systems to express statements about objects and their relationships. By using logic, we can ensure that the reasoning process is structured and follows specific rules, which helps in making sure that the conclusions drawn are valid.
- **Automatic theorem proving**
 - This is a field within AI where computer programs are designed to prove mathematical theorems automatically. These programs use logical notation to represent problems and attempt to solve them. However, a challenge arises because some problems might not have a solution, and the program could run indefinitely trying to find one. This is related to the *halting problem*, which is a concept in computer science that deals with determining whether a program will eventually stop running or continue indefinitely.

Thinking Rationally: Cons

1. **Formalizing informal knowledge is difficult**

- Example: “A handshake occurs when two people extend, grip, shake hands, then release.”
- Formal logic representation:

$$\begin{aligned} &\exists x, y \text{ (Person}(x) \wedge \text{Person}(y) \wedge x \neq y \wedge \\ &\quad \text{Hand}(x, h_x) \wedge \text{Hand}(y, h_y) \wedge \\ &\quad \text{MoveToward}(h_x, h_y) \wedge \text{Contact}(h_x, h_y) \wedge \\ &\quad \text{Shake}(h_x, h_y) \wedge \\ &\quad \text{Release}(h_x, h_y)) \end{aligned}$$

2. **Probabilistic nature of knowledge**

- Example in medicine: “Fever, cough, and fatigue could indicate flu, COVID-19, or another illness.”

3. **Scalability challenges**

- Large problems may need heuristics for practical solutions

4. **Intelligence requires more than rational thinking**

- Importance of agent interaction with the world
- Problem of the “body”



8 / 19

- **Formalizing informal knowledge is difficult**

- When we try to convert everyday actions or knowledge into a formal, logical structure, it becomes quite complex. For instance, consider the simple act of a handshake. In our daily lives, we understand this as a friendly gesture between two people. However, when we attempt to describe it using formal logic, it involves defining each step and condition, such as identifying two people, their hands, the movement towards each other, the contact, the shaking motion, and finally the release. This example illustrates how challenging it is to capture the nuances of informal knowledge in a rigid, logical framework.

- **Probabilistic nature of knowledge**

- Many real-world situations involve uncertainty and cannot be described with absolute certainty. For example, in medicine, symptoms like fever, cough, and fatigue can be indicative of multiple illnesses such as the flu, COVID-19, or other conditions. This uncertainty requires a probabilistic approach to reasoning, where we assess the likelihood of various outcomes rather than relying on definitive logic.

- **Scalability challenges**

- As problems grow in size and complexity, finding solutions using purely logical reasoning becomes impractical. In such cases, heuristics—rules of thumb or simplified strategies—are often employed to make the problem manageable. These heuristics help in finding good enough solutions within a reasonable time frame, even if they are not perfect.

- **Intelligence requires more than rational thinking**

- While rational thinking is a crucial component of intelligence, it is not sufficient on its own. Intelligent agents, whether human or artificial, need to interact with the world to gather information and learn. This interaction highlights the importance of having

a “body” or a means to perceive and act within an environment, which is essential for developing a comprehensive understanding and making informed decisions.

3. AI as Acting Humanly

- **Agent** is something that perceives and acts to reach a goal
- **Definition:** AI designs **agents that can act like humans**
- **Turing test**
 - “A computer passes the Turing test if a human cannot tell whether the answers to questions came from a person or a computer”
- Passing the (embodied) Turing test requires:
 1. Natural language processing to communicate
 2. Knowledge representation to store information
 3. Automated reasoning to use stored knowledge and answer questions
 4. Machine learning to detect patterns
 5. Computer vision and speech recognition to perceive objects and understand speech
 6. Robotics to manipulate objects and move



9 / 19

- **Agent** is something that perceives and acts to reach a goal
 - An *agent* in AI is essentially a system or entity that can observe its environment through sensors and act upon it using actuators to achieve specific objectives. This concept is fundamental in AI as it forms the basis for creating systems that can operate autonomously.
- **Definition:** AI designs **agents that can act like humans**
 - The goal of AI is to create agents that can mimic human behavior. This involves designing systems that can perform tasks typically requiring human intelligence, such as understanding language, recognizing objects, and making decisions.
- **Turing test**
 - “A computer passes the Turing test if a human cannot tell whether the answers to questions came from a person or a computer”
 - The Turing test, proposed by Alan Turing, is a measure of a machine’s ability to exhibit intelligent behavior indistinguishable from that of a human. If a machine can engage in a conversation with a human without the human realizing they are talking to a machine, it is said to have passed the test.
- Passing the (embodied) Turing test requires:
 1. **Natural language processing to communicate**
 - This involves the ability of a machine to understand and generate human language, allowing it to interact naturally with people.
 2. **Knowledge representation to store information**
 - AI systems need to store and organize information in a way that allows them to retrieve and use it effectively.
 3. **Automated reasoning to use stored knowledge and answer questions**
 - This is the capability of a system to apply logic and reasoning to the information it

has to solve problems or answer questions.

4. **Machine learning to detect patterns**

- Machine learning enables systems to learn from data, identify patterns, and improve their performance over time without being explicitly programmed.

5. **Computer vision and speech recognition to perceive objects and understand speech**

- These technologies allow machines to interpret visual and auditory information from the world, similar to human senses.

6. **Robotics to manipulate objects and move**

- Robotics involves the design and use of robots that can physically interact with their environment, performing tasks that require movement and manipulation.

10 / 19: Turing Test: Pros and Cons

Turing Test: Pros and Cons

- **Pros**
 - Operational definition of intelligence
 - Sidestep philosophical vagueness
 - “What is consciousness?”
 - “Can a machine think?”
 - ...
- **Cons**
 - **Anthropomorphic** criteria define intelligence in human terms
 - Multiple forms of non-human intelligence exist
 - Intelligence in terms of Turing test is **fooling humans** into thinking it's human
 - E.g., aeronautical engineering is about:
 - Yes: Focus on wind tunnels and aerodynamics
 - No: Designing machines that imitate birds



10 / 19

- **Pros**
 - **Operational definition of intelligence:** The Turing Test provides a clear and practical way to measure machine intelligence. Instead of getting bogged down in complex theories, it offers a straightforward test: if a machine can engage in conversation indistinguishably from a human, it is considered intelligent.
 - **Sidestep philosophical vagueness:** The Turing Test helps avoid deep philosophical questions that are hard to answer, like “What is consciousness?” or “Can a machine think?” By focusing on observable behavior, it bypasses these debates and provides a more tangible goal for AI development.
- **Cons**
 - **Anthropomorphic criteria define intelligence in human terms:** The Turing Test assumes that human-like behavior is the only measure of intelligence. This is limiting because there are many forms of intelligence that don't resemble human thinking, such as animal intelligence or even unique machine intelligence.
 - **Intelligence in terms of Turing test is fooling humans into thinking it's human:** The test measures a machine's ability to mimic human conversation, which might not truly reflect its intelligence. It's more about deception than genuine understanding or problem-solving.
 - **E.g., aeronautical engineering is about:** This analogy highlights that just as aeronautical engineering focuses on principles like aerodynamics rather than mimicking birds, AI should focus on developing unique capabilities rather than just imitating human behavior.

4. AI as Acting Rationally



- **Rational agents:** agents that do the “right thing” given what they know
- Agents that **act rationally** should:
 1. Operate autonomously
 2. Perceive environment
 3. Persist over a prolonged time period
 4. Adapt to change
 5. Create and pursue goals



11 / 19

- **Rational agents:** These are systems or entities designed to make decisions that are considered the “right” or most effective based on the information they have. The idea is that these agents should act in a way that maximizes their chances of success or achieving their goals.
- Agents that **act rationally** should:
 1. **Operate autonomously:** This means the agent should be able to function on its own without needing constant human intervention. It should make decisions independently.
 2. **Perceive environment:** The agent must be able to gather information from its surroundings. This could involve sensors, cameras, or other data-gathering tools that help it understand the context in which it operates.
 3. **Persist over a prolonged time period:** The agent should be able to continue functioning effectively over time, not just in short bursts. This involves maintaining performance and adapting as needed.
 4. **Adapt to change:** As the environment or circumstances change, the agent should be able to adjust its behavior or strategies to remain effective.
 5. **Create and pursue goals:** The agent should have the ability to set objectives and work towards achieving them, making decisions that align with these goals. This involves planning and prioritizing actions to reach desired outcomes.

12 / 19: Acting Rationally as Ultimate Goal of AI

Acting Rationally as Ultimate Goal of AI

- Which definition of AI to use?
 - Acting vs. Thinking
 - Rational vs. Human
- **Acting > Thinking**
 - Acting rationally is broader than just thinking rationally
- **Rational > Human**
 - Rationality can be mathematically defined
 - Human behavior is shaped by evolutionary conditions
- AI focuses on **agents acting rationally**



12 / 19

- *Acting Rationally as Ultimate Goal of AI*: This slide discusses the ultimate aim of Artificial Intelligence (AI), which is to act rationally. It highlights the importance of choosing the right definition of AI to guide its development.
- **Which definition of AI to use?**
 - *Acting vs. Thinking*: The slide contrasts two approaches to defining AI. Acting involves taking actions to achieve goals, while thinking focuses on the process of reasoning. The emphasis here is on the outcome (acting) rather than the process (thinking).
 - *Rational vs. Human*: This point compares rational behavior, which can be defined using mathematical models, with human behavior, which is influenced by evolutionary factors and may not always be rational.
- **Acting > Thinking**: The slide suggests that acting rationally encompasses more than just thinking rationally. It implies that the ability to make decisions and take actions is more comprehensive and practical than merely processing thoughts.
- **Rational > Human**: Rationality is preferred over mimicking human behavior because it can be precisely defined and measured. Human behavior, on the other hand, is complex and often irrational due to its evolutionary background.
- ***AI focuses on agents acting rationally***: The ultimate goal of AI is to create agents that can act rationally. This means they should be able to make decisions and take actions that are optimal or near-optimal in achieving their objectives, based on the information available to them.

13 / 19: Rationally is Not Absolute

Rationally is Not Absolute

- AI wants to build agents that **do the right thing**
 - What is the right thing?
- E.g., you leave the house and a branch strikes you
 - **Q:** Did you act rationally?
 - Probably
- E.g., you cross the street and a car knocks you over
 - **Q:** Did you act rationally?
 - It depends, but probably no
- E.g., moral issues with self-driving car
 - Swerve and hit a pedestrian to avoid a frontal crash that would kill 2 people



13 / 19

- **Rationality is Not Absolute**
 - **AI wants to build agents that do the right thing**
 - * The goal of AI is to create systems that can make decisions that are considered “right” or appropriate. However, determining what the “right thing” is can be complex and context-dependent.
 - **E.g., you leave the house and a branch strikes you**
 - * **Q:** Did you act rationally?
 - * Probably. In this scenario, you couldn’t predict or control the branch falling. Your decision to leave the house was rational based on the information available to you at the time.
 - **E.g., you cross the street and a car knocks you over**
 - * **Q:** Did you act rationally?
 - * It depends, but probably no. If you crossed without checking for traffic, it might be considered irrational. However, if you followed all safety rules and the car was at fault, your actions could still be rational.
 - **E.g., moral issues with self-driving car**
 - * This example highlights the ethical dilemmas AI systems face. A self-driving car might have to choose between swerving to avoid a crash, potentially harming a pedestrian, or staying on course and risking the lives of its passengers. These situations challenge the notion of rationality, as the “right” decision can vary based on ethical perspectives and societal values.

14 / 19: Problems of a Rational Agent

Problems of a Rational Agent

- **Probabilistic environment**
 - A rational agent aims for:
 - The best outcome in a deterministic setup
 - The best expected outcome under uncertainty
- **Best** is determined by the objective function:
 - E.g., cost function, sum of rewards, loss function, utility
- Omniscience vs **no-regrets**
 - Best based on available information
- Sometimes **no provably correct action** exists
 - Yet, an action must be taken
- Even **with perfect information** rationality can't be feasible due to:
 - Cost of acquiring all data (e.g., in medicine)
 - Computational demands
- Perfect good enough vs perfect
 - Acting appropriately ("satisficing")



14 / 19

- **Probabilistic environment**
 - A *rational agent* is an entity that makes decisions to achieve the best possible outcome. In a *deterministic setup*, where everything is predictable, the agent aims for the absolute best result. However, in real-world scenarios, uncertainty is common, so the agent strives for the best *expected* outcome, which means making decisions that are likely to lead to the best results based on probabilities.
- **Best** is determined by the objective function:
 - The term “best” is subjective and depends on the *objective function* used. This could be a *cost function* (minimizing costs), a *sum of rewards* (maximizing benefits), a *loss function* (minimizing errors), or *utility* (maximizing satisfaction or usefulness).
- Omniscience vs **no-regrets**
 - Omniscience implies having complete knowledge, which is rarely possible. Instead, agents aim for a *no-regrets* approach, making the best decisions possible with the information available at the time.
- Sometimes **no provably correct action** exists
 - In many situations, there isn't a clear-cut correct action. Despite this, decisions still need to be made, often based on the best available evidence or predictions.
- Even **with perfect information** rationality can't be feasible due to:
 - Even if all information were available, it might not be practical to use it all. For example, in fields like medicine, gathering all possible data can be too costly or time-consuming. Additionally, processing this data might require more computational power than is feasible.
- Perfect good enough vs perfect
 - In many cases, aiming for “perfect” is unrealistic. Instead, agents often aim for “good

enough,” a concept known as *satisficing*. This means making decisions that are satisfactory and sufficient, rather than optimal, given the constraints.

Machine Learning: Definitions

- How to define machine learning?
- “Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed” (Samuel, 1959)
- **Machine learning** is about building machines to do **useful things** without being **explicitly programmed**
 - E.g. a computer learns to play checkers by playing against itself, memorizing positions that lead to winning
- “A computer program is said to learn from experience E with respect to some task T and some performance measure P , if $P(T)$ improves with experience E ” (Mitchell, 1998)
- E.g.,
 - Computer vision
 - Speech recognition
 - Natural language processing



15 / 19

- **Machine Learning: Definitions**
- **How to define machine learning?**
 - Machine learning is a branch of computer science focused on developing algorithms that allow computers to learn from and make predictions or decisions based on data. This means that instead of being programmed with specific instructions for every possible scenario, the computer can improve its performance over time by learning from examples.
- “Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed” (Samuel, 1959)
 - This classic definition by Arthur Samuel highlights the core idea of machine learning: enabling computers to learn from data and experiences rather than relying solely on pre-defined rules. This approach allows computers to adapt to new situations and improve their performance over time.
- **Machine learning** is about building machines to do **useful things** without being **explicitly programmed**
 - For example, a computer can learn to play checkers by playing against itself repeatedly. Through this process, it identifies and memorizes board positions that lead to winning outcomes, improving its strategy over time.
- “A computer program is said to learn from experience E with respect to some task T and some performance measure P , if $P(T)$ improves with experience E ” (Mitchell, 1998)
 - Tom Mitchell’s definition emphasizes the importance of experience in machine learning.

A program is considered to be learning if its performance on a specific task improves as it gains more experience.

- **Examples of machine learning applications:**

- *Computer vision*: Teaching computers to interpret and understand visual information from the world, such as recognizing objects in images.
- *Speech recognition*: Enabling computers to understand and process human speech, allowing for applications like voice-activated assistants.
- *Natural language processing*: Allowing computers to understand, interpret, and generate human language, which is crucial for applications like chatbots and translation services.

16 / 19: Limits of ML Compared to Human Intelligence

Limits of ML Compared to Human Intelligence

- **AI differs from human intelligence**
 - Machines don't learn like humans (e.g., LLMs)
- **Fragility to input variations**
 - ML models fail with slight input distortions
 - Adversarial attacks cause misclassification by altering one pixel
 - A model trained for a video game may fail if the screen is slightly rotated; humans continue effortlessly
- **Lack of transfer learning**
 - ML systems cannot apply knowledge across domains without retraining
- **Massive data and compute requirements**
 - ML requires enormous datasets and computational resources
 - A teenager learns to drive in hours
 - Self-driving systems need billions of compute hours and extensive data
- **Poor common sense and reasoning**
 - ML lacks built-in world knowledge and intuitive logic



16 / 19

- **AI differs from human intelligence**
 - While both AI and humans can learn, the way they do it is quite different. For example, *Large Language Models (LLMs)*, which are a type of AI, learn patterns from vast amounts of text data, but they don't understand or learn in the same way humans do. Humans can learn from a few examples and apply reasoning, whereas AI relies heavily on data patterns.
- **Fragility to input variations**
 - Machine learning models can be surprisingly fragile. Even small changes in the input data can lead to significant errors. For instance, *adversarial attacks* can trick a model into making mistakes by altering just a single pixel in an image. Similarly, a model trained to play a video game might fail if the screen is slightly rotated, while humans can easily adapt to such changes without any problem.
- **Lack of transfer learning**
 - Unlike humans, who can apply knowledge from one area to another, machine learning systems struggle with this. They typically need to be retrained from scratch when faced with a new domain or task, as they can't naturally transfer their learning across different contexts.
- **Massive data and compute requirements**
 - Machine learning models require a lot of data and computing power to learn effectively. For example, a teenager can learn to drive a car in just a few hours of practice. In contrast, self-driving car systems need to process vast amounts of data and require billions of compute hours to achieve a similar level of proficiency.
- **Poor common sense and reasoning**
 - Machine learning models often lack *common sense* and the ability to reason intuitively.

They don't have an inherent understanding of the world or the ability to apply logic in the way humans do. This means they can make decisions that seem illogical or lack the depth of understanding that comes naturally to humans.

17 / 19: Limits of ML Compared to Human Intelligence

Limits of ML Compared to Human Intelligence

- **Opaque decision-making**
 - Many ML models offer little transparency into decision processes
 - Limits trust, interpretability, and accountability in critical applications
- **Dependence on narrow objectives**
 - ML systems excel at optimizing narrow tasks but fail with ambiguous goals
 - E.g., an algorithm maximizing user engagement may promote harmful content
- **Susceptibility to bias and data quality**
 - Models inherit and amplify biases in training data
- **Lack of embodiment and physical interaction**
 - Human cognition is grounded in physical and sensory experience



17 / 19

- **Opaque decision-making**
 - Machine learning models, especially complex ones like deep neural networks, often act as “black boxes.” This means that while they can make accurate predictions or decisions, it’s not always clear how they arrived at those conclusions. This lack of transparency can be problematic, especially in critical areas like healthcare or finance, where understanding the reasoning behind a decision is crucial for trust and accountability.
- **Dependence on narrow objectives**
 - Machine learning systems are designed to perform specific tasks very well, but they struggle when the goals are not clearly defined. For example, an algorithm designed to increase user engagement might inadvertently promote content that is misleading or harmful because it focuses solely on engagement metrics without understanding the broader context or consequences.
- **Susceptibility to bias and data quality**
 - Machine learning models learn from data, and if that data contains biases, the models will likely reflect and even amplify those biases. This can lead to unfair or discriminatory outcomes, highlighting the importance of using high-quality, representative data and implementing strategies to mitigate bias.
- **Lack of embodiment and physical interaction**
 - Unlike humans, machine learning models do not have physical bodies or sensory experiences. Human intelligence is deeply connected to our physical interactions with the world, which provide context and understanding that machines currently lack. This absence of embodiment limits the ability of machines to fully replicate human-like understanding and reasoning.

18 / 19: The 3 Machine Learning Assumptions

The 3 Machine Learning Assumptions

- In practice, ML involves solving a practical problem by:
 - Gathering a dataset
 - Building a statistical model from the dataset algorithmically
- The **three assumptions** of machine learning
 - A **pattern exists**
 - Pattern cannot be **precisely defined mathematically**
 - **Data is available**
- Which ML assumption is **essential**?
 - A pattern exists
 - If no pattern, try learning, measure effectiveness, conclude it doesn't work
 - Pattern cannot be precisely defined mathematically
 - If solution is direct, ML not recommended, but may still apply
 - Data is available
 - Without data, no progress can be made
 - **Data is crucial**



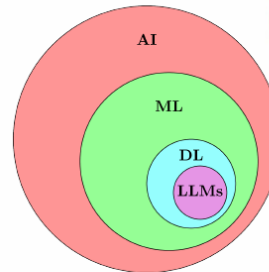
18 / 19

- **In practice, ML involves solving a practical problem by:**
 - **Gathering a dataset:** This is the first step in any machine learning project. You need data to train your model. The data should be relevant to the problem you're trying to solve.
 - **Building a statistical model from the dataset algorithmically:** Once you have your data, you use algorithms to create a model that can make predictions or decisions based on that data.
- **The three assumptions of machine learning:**
 - **A pattern exists:** Machine learning relies on the idea that there is some underlying pattern in the data that can be learned. If no pattern exists, the model won't be able to make accurate predictions.
 - **Pattern cannot be precisely defined mathematically:** If you can define the solution with a simple formula, you don't need machine learning. However, ML can still be useful for complex patterns that are hard to define.
 - **Data is available:** You need data to train your model. Without data, you can't make any progress. This is why data is often considered the most crucial element in machine learning.
- **Which ML assumption is essential?**
 - **A pattern exists:** This is the foundation of machine learning. If there's no pattern, the model won't work. You can try to learn from the data, but if it doesn't work, it might be because there's no pattern.
 - **Pattern cannot be precisely defined mathematically:** If you can solve the problem with a straightforward formula, machine learning might not be necessary. However, it can still be useful for more complex problems.

-
- **Data is available:** Without data, you can't train a model. This makes data availability crucial for any machine learning project.

AI vs ML vs Deep-Learning

- **AI**
 - Machines programmed to reason, learn, and act in a rational way
- **ML**
 - Machines capable of performing tasks without being explicitly programmed
- **AI without ML:**
 - Example: Rule-based systems (e.g., IBM Deep Blue playing chess)
- **Deep Learning (DL)**
 - ML using neural networks with many layers
- **Large Language Models (LLM)**
 - Neural networks trained on massive text datasets and RLHS



19 / 19

- **AI**
 - **Artificial Intelligence (AI)** refers to the broad concept of machines being able to carry out tasks in a way that we would consider “smart.” This involves programming machines to mimic human reasoning, learning, and decision-making processes. AI encompasses a wide range of technologies and applications, from simple rule-based systems to complex neural networks.
- **ML**
 - **Machine Learning (ML)** is a subset of AI focused on the idea that machines can learn from data. Instead of being explicitly programmed to perform a task, machines use algorithms to identify patterns and make decisions based on data. This allows them to improve over time as they are exposed to more data.
- **AI without ML**
 - AI can exist without ML, as seen in rule-based systems like IBM’s Deep Blue, which played chess by following pre-defined rules and strategies rather than learning from data.
- **Deep Learning (DL)**
 - **Deep Learning** is a specialized form of ML that uses neural networks with many layers (hence “deep”) to analyze various levels of data abstraction. This approach is particularly effective for tasks like image and speech recognition.
- **Large Language Models (LLM)**
 - **Large Language Models** are a type of deep learning model specifically designed to understand and generate human language. They are trained on vast amounts of text data and often use techniques like Reinforcement Learning from Human Feedback (RLHF) to refine their outputs, making them capable of tasks like translation, summarization, and conversation.