Ivan Manov

# AI Ethics

## Course Notes

365 DataScience

# Table of Contents

365√DataScience

365√DataScience

# Abstract

AI Ethics is a field that explores how artificial intelligence systems can be developed and deployed responsibly. It centers around core principles like transparency, fairness, accountability, and privacy—ensuring that technology serves the common good without causing harm.

In this course, we examine the ethical implications of AI across its entire lifecycle—from data collection and model training to deployment and regulation. We begin by introducing foundational concepts of AI and its rapid growth, particularly in generative technologies. Then, we dive into real-world ethical challenges such as bias, misinformation, privacy violations, and job displacement. Next, we explore practical strategies for ethical data sourcing, inclusive model development, and mitigating issues like hallucinations or inconsistent outputs. Finally, we analyze global AI regulatory frameworks, highlighting how governments and organizations are addressing the risks and opportunities presented by AI.

These course notes provide a structured summary of the core ideas, tools, and dilemmas introduced in the video lessons. They act as a companion resource to help you navigate the ethical dimension of AI, lesson by lesson.

*Keywords: artificial intelligence, ethics, bias, privacy, fairness, accountability, transparency, regulation, data, generative AI*

# 1. Getting started

## 1.1 Introduction

AI ethics aims to ensure that artificial intelligence is developed and used in ways that are fair, transparent, and respectful of human rights. It helps you understand how to identify ethical risks, make informed choices, and create AI systems that people can trust.

The field addresses a wide range of concerns—from biased algorithms and misinformation to data misuse and the loss of human oversight. It also involves understanding the broader social, legal, and economic consequences of deploying intelligent systems.

AI ethics is no longer a niche concern. It is central to how modern technologies are designed and deployed. From hiring algorithms to chatbots and facial recognition systems, ethical considerations shape how these systems interact with people, influence decision-making, and impact society.

## 1.2 Intro to AI

*What makes human beings unique?*

From solving puzzles to painting masterpieces and engineering rockets, intelligence is the common thread that runs through all our achievements. The Oxford Dictionary defines *intelligence* as *"the ability to acquire and apply knowledge and skills."* That definition captures how our minds naturally absorb and act on information—from everyday decisions to world-changing inventions.

Our *natural intelligence* allows us to learn new skills, solve problems, and adapt to the world around us. It's how we improve day by day, whether it's through understanding music theory or coming up with innovative technology.

AI stems from the idea of building systems that replicate some aspects of human intelligence. This course teaches you what AI is, how it works, and most importantly—how to use it ethically and responsibly.

## 1.3 AI vs data science vs machine learning

- Artificial Intelligence (AI): the science of making machines "smart"—capable of learning, adapting, and performing tasks that typically require human cognition.
- Machine Learning (ML): a subfield of AI focused on using data to make predictions or decisions without being explicitly programmed.
- Data Science encompasses data processing, statistical analysis, and ML to extract insights.

Data scientists often use ML to solve business problems but also rely on non-ML tools like visualization and hypothesis testing.

## 1.4 The AI lifecycle

The 365 AI Lifecycle introduces a structured view of AI development, breaking it into stages:

1. Data collection – gathering raw data (text, images, transactions, etc.).
2. Data preprocessing – cleaning and transforming data for training (e.g., handling missing values, encoding).
3. Model training – feeding data into algorithms to learn patterns.
4. Model evaluation – assessing performance on new data using metrics like accuracy.
5. Model deployment – integrating the model into real-world applications.
6. Monitoring and maintenance – ongoing oversight, retraining with new data.

365√DataScience

These technical steps raise major ethical questions, especially in areas like:

- Data privacy
- Fairness in decision-making
- Transparency of model behavior

To address this, the course groups lifecycle stages into three **e**thical categories:

- Ethical Data Collection = data collection + preprocessing
- Ethical AI Development = training + evaluation
- Ethical AI Deployment = deployment + monitoring



*Figure 1: The 365 AI lifecycle infographic*

365√DataScience

# 2. Introduction to AI and Data Ethics

## 2.1 The rise of Generative AI and its ethical challenges

Generative AI refers to systems that are designed to create new content—text, images, audio, code, or video—based on the data they were trained on. Unlike traditional AI models that classify or predict, generative models produce original outputs by learning from millions of examples.

An example is ChatGPT, launched in November 2022, which demonstrated the power of large language models (LLMs) to simulate human-like conversation. Models like this can summarize text, solve problems, write code, and generate creative content with remarkable fluency.

*How does Generative AI work?*

At the heart of generative AI are deep learning architectures like transformers, which identify patterns in massive datasets. These models are "pre-trained" on broad data (like books, websites, and dialogues), then fine-tuned for specific tasks.

As generative AI tools become more powerful and widely adopted, many *ethical concerns* have emerged:

- *Misinformation*: Generative models can produce highly convincing but false content. For example, fabricated news stories or realistic-looking images can spread disinformation quickly, leading to public panic or manipulation.
- *Bias and discrimination*: If training data contains historical biases—such as overrepresenting certain schools or demographics—models may reproduce these patterns. This can lead to unfair outcomes, like job applicants from underrepresented backgrounds being filtered out.
- *Privacy violations*: Generative models may unintentionally expose sensitive information if trained on private or proprietary data. They can also output personal data fragments without user consent.

365√DataScience

- *Job displacement*: Automation of creative and knowledge work raises concerns about economic inequality. Fields like content writing, graphic design, and customer support may see significant shifts, potentially reducing employment opportunities.

While generative AI unlocks new efficiencies and creativity, unchecked use can lead to serious harm. Building responsible AI means:

- Auditing training data for quality and representativeness
- Implementing safeguards to detect and correct biased or harmful outputs
- Establishing clear accountability and oversight mechanisms
- Respecting data privacy and obtaining proper consent
- Ensuring transparency around how outputs are generated

## 2.2 Why AI Ethics matters more than ever

As artificial intelligence advances at breakneck speed, ethical oversight becomes not just important—but essential. While AI unlocks unprecedented capabilities, its rapid deployment across high-impact sectors like media, education, healthcare, and finance often outpaces the systems meant to regulate it. Misuse or poorly designed AI systems can amplify societal inequalities, compromise democratic processes, and erode public trust.

*"Ethics plays the role of a bicycle brake on an intercontinental airplane."*

This quote by Ulrich Beck highlights how ethical considerations often lag behind technological momentum—but that doesn't mean they're unnecessary. On the contrary, ethical guidelines may be the only safeguards in the absence of fully developed laws.

Several high-profile cases show what can happen when ethics are ignored:

- Deepfake content: AI-generated images and videos—like the viral image of the pope wearing a luxury designer coat—demonstrate how easily false visuals can circulate. These raise concerns about identity manipulation, misinformation, and political influence.

- Misinterpretation by AI: Language models may misread context, leading to bizarre and harmful conclusions. Such errors, while sometimes humorous, underline the risks of context-insensitive automation.

- Job disruption and inequality: AI is projected to displace a significant portion of the workforce, especially among lower-income or lower-education populations. Automation in industries like customer service, transportation, and even creative arts contributes to economic displacement and widening inequality.

- Creative industry impacts: In film and media, AI-generated digital extras and visual effects threaten traditional entry-level roles. This technology reduces costs but also eliminates jobs and alters how human labor is valued in creative fields.

- Voice cloning and consent: Real-world cases, like voice actor Greg Marston discovering his voice was cloned without permission, raise serious questions about consent, ownership, and long-term data use.

- Unethical data collection: Scraping personal data from the internet without consent can lead to serious consequences. The company Clearview AI faced over $30 million in fines across several countries, demonstrating that the absence of ethics can lead to massive financial and legal troubles.

These examples reinforce a central truth: ethical failures are often built into AI systems long before deployment. Issues in data sourcing, model training, and design choices ripple outward into real-world consequences.

To address this, ethics must be considered proactively, not retroactively. That means:

- Using consent-based data practices
- Auditing systems for fairness and bias
- Being transparent about how systems work
- Establishing clear lines of accountability

## 2.3 Ethics vs laws

- Laws are codified rules established by governments. They are enforceable in courts and define what is legally permissible.
- Ethics is broader and more flexible. It informs decision-making in situations where no laws exist or where legal rules lack nuance.

Laws often lag behind technological progress. By the time legislation is drafted and passed, the technology may have already evolved. This makes ethics critical in rapidly moving fields like AI, where waiting for legal guidance could result in unchecked harm.

Ethics doesn't replace law—it supplements it. While laws define the minimum standards of acceptable behavior, ethics challenges individuals and organizations to aim higher, anticipate consequences, and consider broader social impacts.

Responsible AI requires:

- Acting in morally responsible ways, even when no rules exist
- Designing systems that reflect shared human values

## 3. The core principles of AI Ethics

### 3.1 Privacy

Privacy is one of the most essential and widely discussed ethical principles in artificial intelligence. It concerns the protection of individuals' personal information from misuse, unauthorized access, or exploitation by AI systems.

AI models often depend on vast amounts of data—much of which comes from real people. This includes online activity, biometric information, health records, and personal conversations. When these datasets are used without clear consent or proper safeguards, individuals' privacy can be seriously compromised.

## 3.2 Transparency

Transparency in AI refers to the ability to understand how a system works, what data it uses, how it makes decisions, and why it produces certain outcomes. It's about openness and clarity—not only for developers, but also for users, regulators, and those impacted by the system.

In the context of AI ethics, transparency means communicating clearly:

- How data is being handled
- What an AI system is doing
- What potential risks or consequences might arise

This principle ensures people are not left in the dark about how their information is being used or how automated decisions are made.

## 3.3 Accountability

Accountability means being responsible for the outcomes of AI systems—whether those outcomes are positive or harmful. It refers to the ethical obligation of individuals, organizations, and institutions to ensure that AI is developed, deployed, and used responsibly.

This principle is essential because AI systems often involve many actors:

- Developers who build the algorithms
- Engineers who train the models
- Organizations that use the systems

365 DataScience

- Stakeholders who influence implementation decisions

When something goes wrong, it's crucial to identify who is responsible and who will take corrective action.

*Why Is Accountability So Difficult in AI?*

Unlike traditional products, AI systems are complex, dynamic, and data-driven. They may behave in unpredictable ways or be influenced by unseen biases in training data. This makes accountability harder to assign.

For example, when an AI hiring tool discriminates against certain candidates:

- Is the fault with the developers?
- The team that collected the training data?
- Or the company that deployed the system without adequate testing?

Without clearly defined roles and oversight, ethical breaches can go unresolved.

## 3.4 Fairness

Fairness means that AI systems should treat all individuals and groups equitably, avoiding bias and discrimination. It ensures that outcomes are not only accurate but also objective—especially when the AI impacts people's lives in areas such as criminal justice, healthcare, education, or employment.

Fairness is about preventing harm caused by biased data or algorithms and making sure systems work consistently well for everyone.

Fairness must be considered at every stage of the AI development process. Let's walk through the 365 AI Lifecycle:

- Data collection: Collect data that represents diverse groups. For instance, in medical diagnostics, datasets should include people of various ages, genders, and ethnicities to avoid one-size-fits-all outcomes.

- Data preprocessing: Identify and correct imbalances or biases. Techniques like rebalancing datasets or removing problematic features can help prevent biased learning.

- Model training: Conduct fairness inspections. Developers can test how the model performs across different groups and adjust algorithms accordingly (e.g., by adding fairness constraints).

- Model evaluation: Evaluate performance not just for accuracy, but also for fairness. This means checking whether the model performs equally well for different user segments.

- Deployment and monitoring; Continue assessing fairness in the real world. Regular updates, user feedback, and bias audits help detect and fix emerging issues after deployment.

# 4. Ethical data collection

## 4.1 Ethical sourcing and types of data

The AI development process begins with data collection, and the quality of this data plays a critical role in how well a model performs. Poor-quality or irresponsibly sourced data can lead to biased, unsafe, or even unlawful AI systems.

Even though datasets like Common Crawl have powered major models like GPT-3 and Gemini, relying on them can reinforce misinformation and introduce bias from the very beginning of the AI lifecycle.

*What Is Ethical Sourcing?*

Ethical sourcing of data means choosing not only useful data but the right data, with proper attention to:

365√DataScience

- Consent
- Data quality
- Representation
- Legal and social implications

It's not enough to ask whether the data is available. One must also ask:

- Where does the data come from?
- Was it collected with proper consent?
- Does it reflect all groups fairly?

These questions are fundamental to building AI that is trustworthy and responsible.

## 4.2 Proprietary data

- Data collected and owned by an organization (e.g., internal user data, customer records)
- High-quality and controlled, but comes with strong privacy and usage restrictions

## 4.3 Public data

- Freely accessible information (e.g., government statistics, open research data)
- Often lower quality, possibly outdated or biased, and may still include sensitive content

## 4.4 Web-scraped data

- Automatically collected from websites using scraping tools
- Raises legal and ethical questions around consent, terms of use, and privacy—especially when involving user-generated content

Each type of data brings different opportunities and risks. The choice of which type to use affects how the AI will behave, whom it may benefit, and whom it may inadvertently harm.

## 4.5 Dealing with sensitive and protected information

Whether data is proprietary, public, or web-scraped, some of it may be sensitive or legally protected. Handling such data responsibly is a crucial part of ethical AI development. Mishandling protected information—such as personal health records or copyrighted content—can lead to privacy violations, legal consequences, and loss of public trust.

Practical steps to ensure sensitive data is handled ethically and in compliance with relevant laws:

Step 1: Check the Metadata

Metadata is "data about data." It provides useful details such as:

- The data's creator or owner
- Licensing terms
- Any restrictions related to privacy, use, or distribution

Step 2: Verify Licenses

If metadata isn't available or doesn't clarify usage rights, the next step is to check the original source. Most websites or databases have terms of use, licensing details, or user agreements that outline what can and can't be done with the data.

Step 3: Follow Legal Regulations (e.g., GDPR)

Laws like the General Data Protection Regulation (GDPR) lay out specific requirements for handling personal and sensitive data.

365√DataScience

Step 4: Anonymize When Necessary

Anonymization involves removing personally identifiable details (like names, addresses, or IDs) from datasets.

Step 5: Ask for Permission or More Information

If it's unclear whether data is sensitive or protected, the best approach is to ask the source directly. Reaching out to the data provider for clarification or permission is often necessary—especially for proprietary or web-scraped content.

## 4.6 Data bias and fair representation

Bias reflects societal inequalities—and when it's embedded in datasets, it can become part of the AI systems we build. Biased data leads to unfair predictions, harmful outcomes, and discrimination, often against already marginalized groups.

Concrete ways to reduce it in the datasets we use:

1. Ensure diversity in the dataset

Datasets should reflect a range of perspectives, identities, and demographics. For example, when training a facial recognition model, data must include individuals of different:

- Ages
- Ethnicities
- Genders

This helps the system perform equitably across all groups and minimizes the risk of skewed or unfair results.

2. Review and update data regularly

Bias can creep in over time if data becomes outdated or unrepresentative. Frequent updates keep datasets aligned with real-world trends.

365√DataScience

Example: A hiring algorithm should be retrained with updated labor market data to avoid relying on old patterns that may no longer be fair or accurate.

3. Engage affected communities

Involving the people who will be impacted by AI systems—especially during data collection—can dramatically improve fairness.

Example: Voice assistants like Alexa and Siri often struggle with regional accents. By working directly with diverse users, these systems can be trained to understand everyone more accurately.

# 5. Ethical AI development

## 5.1 Ethical challenges in working with labeled data

Labeled data is essential for supervised learning—the process where AI learns by example.

In a labeled dataset, each data point is annotated with the correct output. For example, images are labeled "cat" or "dog," or text comments are marked as "offensive" or "harmless." These labeled pairs are what the model uses to learn how to make predictions.

*Ethical risks:*

- Subjectivity in labeling
- Real-world consequence
- Deliberate manipulation
- Human error and fatigue

Poor-quality or biased labeling doesn't just reduce accuracy—it can result in systems that misjudge, discriminate, or marginalize. That's why it's critical to handle labeled data responsibly and ethically from the start.

365√DataScience

## 5.2 Ethical considerations for unlabeled data

Unlabeled data refers to raw information that lacks predefined tags or categories. This type of data is used in unsupervised learning, where the model detects similarities, groups items, or identifies trends without explicit instructions.

*Ethical risks:*

- Lack of context can lead to bias
- Misinterpretation of patterns

Thoughtful curation, transparency about data sources, and active monitoring are key to using unlabeled data safely and fairly.

## 5.3 Ethical challenges in unsupervised training

Unsupervised training allows AI to learn from data without labeled outcomes. The model identifies patterns, clusters, or relationships by itself—without explicit instruction from developers.

Unsupervised training requires human oversight to avoid ethical disasters. Ethical safeguards include:

- Feeding the AI balanced, representative data
- Closely monitoring its training process
- Documenting decisions and training steps for transparency
- Keeping humans in the loop to detect issues early

365√DataScience

## 5.4 Ethical considerations for supervised fine-tuning

Supervised fine-tuning (SFT) is a post-training phase where a pre-trained model is refined using a carefully curated dataset. While pre-training helps the model understand language structure and general patterns, fine-tuning aligns its behavior with specific goals, use cases, or ethical standards.

*Why it matters ethically?*

Fine-tuning helps AI handle sensitive and complex tasks more responsibly. It ensures that the model's responses reflect desired tone, accuracy, and ethical standards. But this step also presents key ethical challenges—especially if:

- The fine-tuning data is biased or low quality
- The model isn't tested for fairness and appropriateness
- There's no regular review of outputs

If not handled carefully, fine-tuning can reinforce harmful assumptions or lead to discriminatory behavior.

## 5.5 RLHF and ethical AI behavior

Reinforcement Learning from Human Feedback (RLHF) is a method used to train AI systems to behave more ethically and align with human values. Unlike supervised fine-tuning, which relies on labeled examples, RLHF uses human feedback as a guide to shape the model's behavior through rewards and corrections.

While pre-training helps the model understand language structure and general patterns, fine-tuning aligns its behavior with specific goals, use cases, or ethical standards.

RLHF doesn't just improve factual accuracy—it teaches the AI how to:

- Be empathetic and respectful

365√DataScience

- Navigate subjective or sensitive topics

- Respond in ways that feel human-centered

While RLHF offers powerful benefits, it's not without limitations:

- Subjectivity in feedback: Human evaluators may carry their own biases, which could influence what is rewarded or penalized.

- Defining the "right" answer: In complex or nuanced situations, there may not be a single correct response. This makes it hard to set clear ethical standards.

## 5.6 Inclusive and fair AI development practices

Creating inclusive and fair AI systems requires careful calibration, frequent adjustments, and thoughtful oversight. If the system is trained or tested in a narrow way, it may exclude, misinterpret, or disadvantage certain groups.

Practical tips for fair AI development:

1. Use Diverse and Representative Data

2. Avoid One-Size-Fits-All Approaches

3. Test for Edge Cases

4. Keep Humans in the Loop

5. Be Transparent and Accountable

# 6. Ethical AI deployment

## 6.1 Intellectual property and user consent in AI interactions

Intellectual property protects:

365√DataScience

- Ownership rights over AI models and their outputs
- Rights of third-party data sources, such as proprietary research or commercial datasets

User consent ensures that individuals:

- Understand what data is collected
- Know how it will be used
- Can agree or opt out
- Have the option to withdraw consent later

## 6.2 Ethical responsibilities of foundation model developers

Foundation models are large, general-purpose AI systems capable of performing a wide range of tasks—such as writing code, drafting emails, generating images, or composing music. These models are trained on massive datasets and serve as the backbone for many advanced AI applications.

*Key ethical challenges:*

1. Data sourcing and bias
2. Intellectual property risks
3. Opacity and trust

*Developer responsibilities:*

- Prioritizing transparency
- Ensuring user safety and legal compliance
- Developing with inclusion in mind

## 6.3 Common issues in foundation models: Open-source data

Foundation models need enormous amounts of data to function—but most companies don't have enough proprietary data to train them from scratch. That's why many rely on open-source datasets: large, publicly available collections of online content.

These datasets are convenient and cost-effective, but they come with major ethical and practical risks.

*Common problems with open-source data:*

1.  Low-quality or outdated content

    Open datasets often include:

    - Errors
    - Inaccurate facts
    - Outdated or irrelevant information

2.  Bias and representation Issues

    Public data reflects real-world biases. If not addressed, these patterns are absorbed into the AI, which may then reinforce stereotypes or exclude marginalized groups.

3.  Unreliable or harmful content

    Because data is scraped from across the internet, it may include:

    - Hate speech
    - Misinformation
    - Offensive language

    Without proper filtering, these elements can show up in the model's outputs.

*Bias detection methods and fairness metrics*

- Statistical Analysis: Compare how the model treats different demographic groups using traditional metrics.
- Data Visualization: Use charts and graphs to highlight disparities in the model's predictions or outputs.
- Fairness Metrics: These are specific tools to measure equity in decision-making. Examples include:
  - Demographic Parity: Ensures that sensitive attributes (like race or gender) don't influence results.
  - Calibration by Group: Confirms that the model is equally confident in its predictions across all groups.
  - Disparate Impact Ratio: Compares how often different groups receive favorable outcomes (e.g., loan approval rates).

## 6.4 Inconsistency

One of the biggest challenges with foundation models is inconsistency—when the AI gives different answers to the same or similar questions. These inconsistencies can confuse users, spread misinformation, and erode trust in the technology.

*How to reduce inconsistency*

While no system can be perfect, developers can take steps to minimize these issues:

Test for repeatability

- Run multiple variations of the same question
- Identify patterns of inconsistency
- Flag high-risk areas for closer review

Use standardized prompts

- Keep inputs clear, neutral, and consistent
- Reduce ambiguity to guide the AI toward stable outputs

Monitor and document model behavior

- Track how the model responds over time
- Log problematic outputs and address them in future updates

Set expectations

- Inform users that AI may generate different responses
- Provide guidance on when and how to trust outputs
- Suggest double-checking important information with reliable sources

## 6.5. Hallucination

Hallucination occurs when a model produces an answer that is factually incorrect, fabricated, or misleading, while still sounding confident and authoritative. This is one of the most concerning challenges in modern AI, especially with large language models like ChatGPT.

While hallucinations can't be completely eliminated, developers and users can take steps to reduce their frequency and impact:

Improve Training Data

- Use more accurate, high-quality sources—especially in factual domains.

Fine-Tune for Reliability

- Adjust model weights or responses using expert-reviewed content.

Add Disclaimers and Guidance

- Let users know that the AI may generate errors and that outputs should be verified independently.

Reinforce Human Oversight

- Encourage users to double-check important answers, especially when using AI for professional, legal, or health-related advice.

## 6.6. Ongoing monitoring and risk mitigation for deployed AI

Deploying an AI model isn't the finish line—it's the beginning of a continuous responsibility. Once the system is live, developers must monitor it closely to ensure it remains fair, accurate, and ethical over time.

Even well-trained models can drift, behave unpredictably, or encounter unexpected situations in the real world. This is why ongoing monitoring and risk mitigation are essential parts of ethical AI deployment.

*What should be monitored?*

- Fairness and bias: AI systems can evolve in ways that unintentionally start favoring certain groups or ignoring others. Regularly checking fairness metrics helps spot these issues early—before they lead to real-world harm.
- Model performance: Unlike traditional models that give clear right or wrong answers, foundation models produce open-ended responses. This makes evaluation more complex. One effective strategy is to use AI judges—specialized models that assess AI outputs for consistency, factual accuracy, and potential hallucinations.

These tools compare generated content against verified sources and flag unreliable or biased outputs for further review.

- Failure detection and response: No system is perfect. Failures will happen, and organizations must be ready:
    - Set up alert systems for high-risk errors
    - Involve human reviewers in sensitive or impactful decisions
    - Allow users to report issues, giving early visibility into problems

- Regular updates and retraining: Updating the AI with fresh data helps it stay relevant, accurate, and aligned with evolving social norms and regulations.

# 7. Ethical AI for end-users: Individuals

## 7.1 Access to AI technology for businesses of all sizes

AI should be accessible to all businesses—not just those with large budgets and technical teams. Reducing technical and financial barriers is key to enabling widespread use. Fortunately, many tools and platforms now support this goal:

- No-code tools like Google AutoML and Microsoft Power Platform
- Cloud-based services and open-source platforms such as TensorFlow, Hugging Face, and LLaMA
- Programs from tech companies like OpenAI and IBM that help small businesses adopt AI quickly and affordably

Governments are also taking steps.

The UK government has invested £10 million in AI training and support programs. In the U.S., the 2025 announcement of the Stargate project—a private-sector initiative backed by OpenAI, SoftBank, and Oracle—pledged up to $500 billion for AI infrastructure, with 20 data centers planned and over 100,000 jobs expected.

## 7.2 Transparency in AI decision-making process

AI systems are making decisions that affect real lives—who gets hired, who receives a loan, which content is flagged online. But if users don't understand how or why those decisions were made, trust breaks down.

That's why transparency is one of the most important ethical principles in AI. It means people should be able to understand:

365√DataScience

- How an AI system makes decisions
- What kind of data it uses
- Who is responsible for its design and oversight

## 7.3 Ethical use of AI outputs in businesses

When AI outputs are used without proper oversight, businesses risk:

- Sharing misinformation

Example: An AI-generated blog post includes outdated or incorrect facts.

- Making unfair decisions

Example: An AI system recommends one group of customers for discounts while ignoring others without a clear reason.

- Plagiarizing content

Example: AI rephrases material from online sources too closely, risking copyright violations.

- Reinforcing bias

Example: Marketing messages unintentionally exclude certain demographics or promote harmful stereotypes.

These issues aren't always caused by bad intentions. Often, they happen because teams trust the AI too much or assume it doesn't need review.

Best practices for ethical use

To use AI responsibly in a business context, follow these key practices:

- Treat AI as a collaborator—not a decision-maker

Use AI to support human thinking, not replace it. The final judgment should still come from people.

- Always fact-check AI outputs

Don't assume AI-generated content is accurate or unbiased. Review for errors, outdated info, or ethical red flags.

- Credit original sources

If AI helped produce something based on research or existing material, give credit where it's due—just like with any other contributor.

- Watch for bias and exclusion

Be mindful of tone, language, and representation. Ask: *Who might feel left out or misrepresented by this output?*

## 7.4 Responsible AI adoption and risk management for businesses

*A four-step framework for risk-managed AI adoption:*

To adopt AI responsibly, businesses can follow this four-step approach:

- Step 1: Identify the right use case
- Step 2: Select the right tools
- Step 3: Test outputs before use
- Step 4: Train teams on ethical use

# 8. Ethical AI for end-users: Businesses

## 8.1 Equity in access to AI technology

Responsible AI adoption and risk management for businesses 8.1 Equity in access to AI technology

As AI advances, a pressing ethical challenge remains: Who gets to benefit from it? While many individuals and businesses in wealthy regions are already integrating AI into their daily lives, others are being left behind—particularly in low- and middle-income countries where access is limited.

This imbalance raises serious questions about fairness and global inclusion in the age of artificial intelligence.

Many underserved regions face challenges like:

- Limited or unreliable internet connectivity
- Lack of affordable devices
- Insufficient computing resources

Major challenge: the high cost of advanced AI infrastructure—like GPUs and cloud computing.

## 8.2 Ethical consideration in human-AI collaboration

When humans and AI co-create something, we need to ask:

- Who is responsible for the outcome?
- How do we define originality or authorship?
- What happens when AI introduces errors or bias?

*The responsibility gap*

365√DataScience

When an AI tool contributes to a poor decision or harmful content, who should be held accountable—the user, the developer, or the system itself?

*Ethical use requires that:*

- Humans stay in the loop
- Decisions are reviewed and refined
- Responsibility is clearly assigned

Collaboration should not become delegation without human oversight.

## 8.3 Responsible use of AI-generated outputs

Key areas of ethical concern:

1. Accuracy: AI tools may sound confident but still generate false or misleading information. It's essential to fact-check and verify content—especially when it involves data, citations, or public communication.
2. Bias and harmful stereotypes: AI can reproduce biased language or offensive imagery based on the data it was trained on. Review content for fairness, inclusivity, and tone.
3. Copyright and plagiarism: AI-generated content can unintentionally resemble or reproduce copyrighted material. Users must:
   - Avoid directly copying outputs without review
   - Credit original ideas where appropriate
   - Understand the limits of fair use in their field
4. Audience and context: What's appropriate in one setting may be problematic in another. Always consider:
   - Who the content is for
   - Where and how it will be shared
   - What message it conveys about your organization

Practical advice for users:

- Don't treat AI content as ready-to-publish.
- Use it as a starting point, then refine, review, and revise.
- Build in a final check for tone, relevance, and accuracy before releasing anything externally.

# 9. ChatGPT ethics

## 9.1 Understanding ChatGPT

ChatGPT is built on Natural Language Processing (NLP), a subfield of AI focused on enabling machines to understand and generate human language. Here's how it operates:

- Your input is broken into tokens (words, characters, or word fragments).
- The model analyzes how these tokens relate to each other.
- It uses learned patterns from billions of text examples to generate a response.

ChatGPT has important ethical limitations:

- It can produce factually incorrect answers, so its responses should always be verified.
- It may reflect biases from the data it was trained on, leading to harmful or stereotypical outputs.
- It can be misused to generate fake news or misleading content.
- Sharing sensitive or private information with ChatGPT carries privacy risks, especially since inputs may be stored or reused depending on the account type.

## 9.2 Privacy concerns with ChatGPT

Rule of thumb: *Don't share anything with ChatGPT that you wouldn't post publicly.*

For users without a ChatGPT Enterprise or Team account, OpenAI may:

- Store your conversations
- Use them to train and improve the model
- Review them as part of ongoing safety efforts

If you're using ChatGPT with a Team or Enterprise plan, OpenAI promises not to use your data for training or improvement by default. These accounts offer:

- More control over data handling
- Stronger privacy protections
- Features designed for internal use at organizations

Best practices:

- Avoid inputting personal or sensitive data
- Be transparent with clients or team members if AI tools are used in communications
- Choose tools and settings that align with your privacy needs

## 9.3 Open AI's privacy policies and data handling

OpenAI uses collected data to:

- Improve model performance
- Ensure safety and prevent misuse
- Provide technical support
- Conduct research and product development

Some of this data is also reviewed manually by OpenAI teams to help train better systems and detect harmful behavior.

Users have some control over how their data is used, especially in non-Enterprise and non-Team accounts:

- Chat history & training: Users can turn off chat history, which also disables data use for model training.

- Custom instructions: These allow users to tailor ChatGPT's responses, but the data shared in this section is still processed by OpenAI.

- Data deletion: Users can delete individual chats or their entire account. OpenAI states that deletion removes associated data from its systems, though some logs may be retained for auditing or legal purposes.

## 9.4 Misinformation and AI-generated content

AI-generated content can be used to:

- Spread fake news quickly

- Impersonate real people or organizations

- Influence public opinion through large-scale, automated messaging

These actions may be unintentional—such as a user sharing incorrect info from ChatGPT without realizing it—or deliberate, such as using AI to create coordinated disinformation campaigns.

To reduce harm:

- Use AI outputs as starting points, not final answers

- Verify against trusted sources

- Avoid using AI for tasks where accuracy is critical and hard to confirm

## 9.5 ChatGPT Plagiarism

ChatGPT doesn't intentionally copy from specific sources, but that doesn't mean it's free from plagiarism risks. The model generates responses by predicting likely word

sequences based on patterns it learned from large datasets—including books, articles, and websites. So, while it doesn't store or "remember" exact sources, it can still produce content that closely resembles existing work.

*Types of plagiarism risks with ChatGPT:*

1. Repetition of Common Phrases or Ideas: ChatGPT may repeat familiar structures or arguments that sound original but are very close to widely published material.
2. Unintentional Paraphrasing: The AI may reword existing content just enough to seem different—without adding new insight or properly attributing ideas.
3. Lack of Source Attribution: ChatGPT doesn't cite where its ideas come from. So, when users copy and paste outputs, it's unclear whether the content is:
   o Based on existing work
   o Truly original
   o Ethically or legally acceptable to use

*How to use ChatGPT without plagiarizing*

- Use the output as a draft or starting point, not a finished product.
- Fact-check and rewrite in your own words.
- Don't assume content is original just because it was generated.
- When using well-known facts or frameworks, add your own analysis, structure, or interpretation.
- For academic use, always follow your institution's AI usage policies.

## 9.6 Responsible use of ChatGPT: What you can and can't do

| What you can do | What you shouldn't do |
|---|---|
| Use it for brainstorming ideas or outlines | Submit AI-generated work as your own (without attribution) |
| Rewrite or simplify existing text | Enter confidential or personal information |
| Ask questions to learn or explore new topics | Use ChatGPT to mislead or impersonate |

365√DataScience

| Practice communication (emails, interviews, etc.) | Assume everything it says is accurate without fact-checking |
|---|---|

## 9.7 ChatGPT and the environment

*What makes ChatGPT energy-intensive?*

1. Training large models: Training a model like GPT-3 or GPT-4 takes weeks or months on thousands of high-performance GPUs. This process alone can emit as much $CO_2$ as the average car does over several years.
2. High demand for inference (usage): Even after training, generating each new response requires a lot of computing power—especially for longer or more complex prompts.
3. Data center cooling: The servers that power ChatGPT must be kept cool. That means constant air conditioning and water usage, which places stress on natural resources—especially in regions already facing drought or heatwaves.

As AI adoption grows globally, so does its carbon footprint. And unlike some tech systems that run occasionally, generative AI tools are always on, responding in real time to millions of requests. Even simple everyday uses—like asking ChatGPT to write a paragraph—carry energy costs that users rarely see.

# 10. Data and AI regulatory frameworks

## 10.1 Global data and AI regulations

As artificial intelligence becomes more powerful and widespread, governments around the world are stepping in to regulate how it's developed and used. But their approaches vary widely—reflecting different priorities around safety, innovation, ethics, and control. This creates a global patchwork of rules. Some regions stress

tight control to reduce harm, while others prioritize flexibility and innovation—sometimes to stay ahead in the global AI race.

## 10.2 European Union: GDPR and the EU Artificial Intelligence act

GDPR: The foundation of data protection

The General Data Protection Regulation (GDPR) is the EU's flagship data protection law. It went into effect in 2018 and has since influenced privacy legislation around the world.

Key principles of GDPR:

- Individuals must give explicit, informed consent for their data to be used.
- People have the right to access, correct, and delete their personal data.
- Organizations must handle data transparently and securely.

The EU Artificial Intelligence Act

In 2021, the European Commission proposed the AI Act, which takes a risk-based approach to regulation. It classifies AI systems into four categories:

1. Unacceptable Risk
   o These AI systems are banned outright, such as:
     ▪ Social scoring (ranking people based on behavior)
     ▪ Manipulative or exploitative AI
     ▪ Real-time biometric surveillance in public spaces
2. High Risk
   o These systems are heavily regulated and must meet strict requirements before deployment.
   o Includes AI used in:
     ▪ Employment
     ▪ Education

365√DataScience

- Law enforcement
- Critical infrastructure

3.  Limited Risk

    o   These systems must provide transparency notices.

    o   Example: AI chatbots must disclose that users are interacting with a machine.

4.  Minimal or No Risk

    o   Most consumer-facing AI tools fall into this category

## 10.3 United States: AI regulation across states

Unlike the European Union, the United States doesn't have a single, nationwide law that regulates AI. Instead, it relies on a sector-based, decentralized approach. That means:

- Different states can pass their own AI-related laws
- Federal agencies issue guidelines, not hard rules
- Most regulation focuses on specific applications like employment, privacy, or biometrics

The result is a patchwork system where protections vary depending on where you live or work.

## 10.4 Asia-Pacific region: Strong government control

In much of the Asia-Pacific region—especially China, South Korea, and Singapore—AI regulation tends to follow a top-down, government-led model. These countries take a centralized approach, combining strict oversight with rapid innovation goals.

While each country has its own priorities, there's a shared emphasis on:

- National security

- Social stability
- Economic competitiveness

China has positioned itself as both a global AI leader and one of the most active regulators in the field.

## 10.5 Africa's push for AI governance

Several factors slow AI adoption and governance in parts of Africa:

- Limited infrastructure, including unreliable internet access
- Lack of large-scale investment in AI
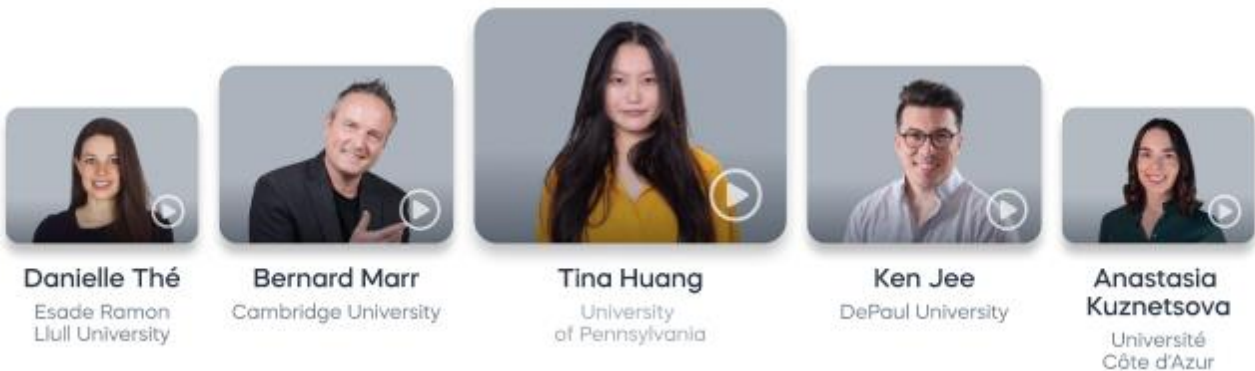- Shortages of skilled professionals and public awareness

Africa's role in the global AI conversation is growing, with policymakers, researchers, and tech advocates working to shape AI in ways that reflect local priorities, cultural values, and regional realities.

365√DataScience

# Learn DATA SCIENCE
# anytime, anywhere, at your own pace.

If you found this resource useful, check out our e-learning program. We have everything you need to succeed in data science.

Learn the most sought-after data science skills from the best experts in the field! Earn a verifiable certificate of achievement trusted by employers worldwide and future proof your career.



| Danielle Thé | Bernard Marr | Tina Huang | Ken Jee | Anastasia Kuznetsova |
|---|---|---|---|---|
| Esade Ramon Llull University | Cambridge University | University of Pennsylvania | DePaul University | Université Côte d'Azur |

## Comprehensive training, exams, certificates.

- ✓ 162 hours of video
- ✓ 599+ Exercises
- ✓ Downloadables

- ✓ Exams & Certification
- ✓ Personalized support
- ✓ Resume Builder & Feedback

- ✓ Portfolio advice
- ✓ New content
- ✓ Career tracks

Join a global community of 1.8 M successful students with an annual subscription

at 60% OFF with coupon code 365RESOURCES.

~~$432~~ **$172.80**/year

## Start at 60% Off

VAT may be applied

365 √ DataScience

**Ivan Manov**

Email: team@365datascience.com

365 DataScience