

---

Undergraduate Certificate in AI for Public Policy and Governance

## AI Regulation and Policy Making

---

**Algorithmic Accountability** – Related: responsibility, auditability. The principle that developers and operators of AI systems must be answerable for outcomes, including unintended harms. Example: A city’s predictive policing tool is audited after biased arrests. Practical application: Mandatory impact reports. Challenge: Tracing decisions through complex models.

**Algorithmic Bias** – Related: discrimination, fairness. Systematic error that produces unfair outcomes for certain groups, often rooted in skewed training data. Example: Facial recognition misidentifying dark-skinned faces. Application: Bias-testing protocols before deployment. Challenge: Hidden biases in large language models.

**Algorithmic Transparency** – Related: explainability, openness. Requirement that AI processes be understandable to stakeholders. Example: A loan-approval algorithm provides a summary of key factors. Application: Dashboards for regulators. Challenge: Trade-off with proprietary IP.

**AI Alignment** – Related: safety, goal specification. Ensuring AI systems pursue objectives consistent with human values. Example: Reinforcement learning agents that respect safety constraints. Application: Alignment research in autonomous vehicles. Challenge: Value-pluralism and specification errors.

**AI Ethics** – Related: principles, moral frameworks. The study of moral implications of AI, guiding responsible design and use. Example: Ethical review boards for health-care AI projects. Application: Institutional ethics guidelines. Challenge: Divergent cultural norms.

**AI Governance** – Related: policy, oversight. Structures and processes that direct AI development and deployment at organizational or societal level. Example: National AI strategy coordinating ministries. Application: Multi-stakeholder advisory committees. Challenge: Rapid innovation outpacing rule-making.

**AI Impact Assessment** – Related: risk analysis, DPIA. Systematic evaluation of potential social, economic, and environmental effects of an AI system before launch. Example: Assessing job displacement from an automated customer-service bot. Application: Mandatory assessment for high-risk AI. Challenge: Quantifying indirect effects.

**AI Regulation** – Related: legislation, compliance. Legal rules governing AI development, deployment, and use. Example: The European AI Act classifying systems by risk tier. Application: Certification schemes for compliant AI products. Challenge: Achieving global harmonization.

**AI Safety** – Related: robustness, reliability. Ensuring AI systems operate without causing unintended harm. Example: Fail-safe mechanisms in autonomous drones. Application: Safety standards for medical diagnosis AI. Challenge: Emergent behavior in complex models.

**AI Strategy** – Related: roadmap, national plan. A coordinated plan outlining goals, investments, and

governance for AI within a jurisdiction. Example: A country's five-year AI investment program. Application: Aligning research funding with public-policy goals. Challenge: Balancing innovation with risk mitigation.

Algorithmic Governance – Related: digital regulation, code-based rules. Use of algorithms to enforce policies, such as automated compliance monitoring. Example: Smart contracts that trigger penalties for non-compliance. Application: Real-time emissions tracking. Challenge: Lack of human oversight.

Artificial General Intelligence – Related: AGI, superintelligence. Hypothetical AI with human-level reasoning across domains. Example: Research projects aiming for flexible problem solving. Application: Long-term policy scenarios. Challenge: Profound uncertainty about societal impact.

Automation Bias – Related: human-automation interaction, overreliance. Tendency of people to trust automated decisions even when erroneous. Example: Operators ignoring alerts from a faulty AI system. Application: Training programs on critical oversight. Challenge: Designing interfaces that encourage appropriate skepticism.

Bias Mitigation – Related: fairness techniques, debiasing. Methods to reduce discriminatory outcomes in AI, such as re-weighting data or adversarial training. Example: Adjusting hiring algorithm scores to achieve gender parity. Application: Toolkits integrated into ML pipelines. Challenge: Preserving model performance while correcting bias.

Certification – Related: standards, conformity assessment. Formal recognition that an AI system meets defined regulatory or ethical standards. Example: ISO/IEC certification for trustworthy AI. Application: Market access for certified products. Challenge: Keeping certifications up-to-date with evolving technology.

Data Governance – Related: data stewardship, policies. Frameworks for managing data quality, security, and access throughout its lifecycle. Example: University establishing a data-use policy for research datasets. Application: Data-ownership registries. Challenge: Coordinating across multiple jurisdictions.

Data Minimization – Related: privacy, principle. Collecting only the data necessary for a specific purpose. Example: Limiting facial-recognition inputs to low-resolution images. Application: Privacy-by-design in AI pipelines. Challenge: Balancing model accuracy with limited data.

Data Privacy – Related: confidentiality, GDPR. Protecting personal information from unauthorized access or disclosure. Example: Encrypting user data used to train recommendation engines. Application: Privacy impact assessments. Challenge: Reconciling big-data analytics with strict privacy laws.

Data Protection Impact Assessment (DPIA) – Related: risk assessment, GDPR. Formal process to identify and mitigate privacy risks of data processing activities. Example: Conducting a DPIA before deploying a health-monitoring AI. Application: Mandatory for high-risk processing under EU law. Challenge: Resource intensity for small firms.

Data Provenance – Related: lineage, traceability. Documentation of the origin and history of data used in AI models. Example: Recording sensor calibration logs for autonomous-vehicle training data. Application: Audit trails for compliance. Challenge: Maintaining provenance at scale.

**Data Sovereignty** – Related: jurisdiction, cross-border data. Principle that data is subject to the laws of the country where it is collected. Example: Cloud providers storing EU citizen data on servers within the EU. Application: National data-localization policies. Challenge: Global supply chains and multinational AI services.

**Decision-Making Autonomy** – Related: automation level, control. Degree to which an AI system can act without human intervention. Example: Fully autonomous trading bots executing orders. Application: Tiered autonomy frameworks. Challenge: Determining appropriate human-in-the-loop thresholds.

**Explainable AI (XAI)** – Related: interpretability, transparency. Techniques that make AI model decisions understandable to humans. Example: SHAP values highlighting feature contributions in a credit-scoring model. Application: Regulatory compliance in finance. Challenge: Explanations may oversimplify complex model behavior.

**Fairness** – Related: equity, non-discrimination. Ensuring AI outcomes do not systematically disadvantage protected groups. Example: Adjusting algorithmic thresholds to achieve equal false-positive rates across races. Application: Fairness dashboards for public-sector AI. Challenge: Multiple, sometimes conflicting, fairness metrics.

**Framework** – Related: guidelines, structure. Structured set of principles and processes guiding AI development and governance. Example: A government releasing a “Responsible AI Framework” with pillars on safety, fairness, and accountability. Application: Aligning institutional practices with policy goals. Challenge: Translating high-level principles into operational steps.

**Human-Centred AI** – Related: user-focused, design. Designing AI systems that prioritize human values, needs, and agency. Example: Chatbots that allow users to easily correct misinterpretations. Application: User-testing protocols for public-service bots. Challenge: Balancing efficiency with user control.

**Human-in-the-Loop (HITL)** – Related: oversight, supervision. Design pattern where humans review or intervene in AI decisions. Example: Radiologists confirming AI-generated diagnoses. Application: Mandatory HITL for high-risk AI. Challenge: Ensuring timely human response without bottlenecks.

**Impact Evaluation** – Related: assessment, outcomes. Systematic study of the actual effects of AI policies after implementation. Example: Measuring changes in unemployment after a nationwide automation subsidy. Application: Longitudinal surveys and data analysis. Challenge: Isolating AI effects from other variables.

**Innovation Sandbox** – Related: regulatory testbed, pilot. Controlled environment where new AI applications can be trialed under relaxed regulations. Example: A city allowing autonomous shuttle pilots in a designated district. Application: Iterative policy refinement. Challenge: Managing risk while fostering experimentation.

**International Cooperation** – Related: multilateralism, standards. Collaborative efforts among nations to harmonize AI rules and share best practices. Example: OECD’s AI Principles adopted by member countries. Application: Joint research funding. Challenge: Reconciling divergent regulatory philosophies.

**Interpretability** – Related: explainability, transparency. Ability to understand how an AI model arrives at a

particular output. Example: Decision trees that can be visualized end-to-end. Application: Model selection for regulated sectors. Challenge: Deep neural networks often lack straightforward interpretability.

Jurisdictional Scope – Related: territorial reach, applicability. The geographic extent to which AI regulations apply. Example: EU AI Act affecting any provider serving EU citizens. Application: Compliance mapping for multinational firms. Challenge: Overlapping authorities and conflicting rules.

Know-Your-Customer (KYC) Automation – Related: AML, compliance. Use of AI to verify client identities and assess risk. Example: Facial-recognition verification for online banking sign-up. Application: Streamlined onboarding. Challenge: Privacy concerns and false-positive rates.

Legal Liability – Related: responsibility, damages. Legal accountability for harms caused by AI systems. Example: A manufacturer sued for injuries caused by a misbehaving robot. Application: Insurance products for AI risk. Challenge: Attributing fault in autonomous decision chains.

Machine Learning (ML) – Related: algorithm, training. Subfield of AI focusing on statistical models that improve with data. Example: Gradient-boosted trees predicting equipment failures. Application: Predictive maintenance in public infrastructure. Challenge: Model drift over time.

Model Auditing – Related: review, compliance. Independent examination of AI models to verify adherence to standards. Example: External auditors reviewing a predictive sentencing algorithm for bias. Application: Certification processes. Challenge: Proprietary models limiting auditor access.

Model Governance – Related: lifecycle management, oversight. Policies and procedures governing model development, deployment, and retirement. Example: A municipal agency maintaining a registry of all AI models in use. Application: Version control and change-management. Challenge: Ensuring consistent governance across departments.

Model Drift – Related: concept shift, performance decay. Degradation of model accuracy as underlying data distributions change. Example: A traffic-prediction model becomes less accurate after new road constructions. Application: Continuous monitoring and retraining. Challenge: Detecting drift early enough to avoid adverse decisions.

Multi-Stakeholder Engagement – Related: participation, co-design. Involving diverse groups—government, industry, civil society—in AI policy development. Example: Public consultations on facial-recognition deployment. Application: Advisory panels with citizen representatives. Challenge: Balancing power dynamics and ensuring meaningful input.

National AI Strategy – Related: policy framework, roadmap. Government-level plan outlining priorities for AI research, education, and regulation. Example: A country allocating funds for AI talent development. Application: Aligning national budgets with strategic goals. Challenge: Maintaining flexibility amid rapid tech change.

Neural Network – Related: deep learning, architecture. Computational model inspired by biological neurons, capable of learning complex patterns. Example: Convolutional networks for image classification. Application:

Disease detection from medical scans. Challenge: Opacity and high computational cost.

Open-Source AI – Related: collaboration, transparency. AI software whose source code is publicly available. Example: An open-source library for natural-language processing. Application: Fostering innovation and peer review. Challenge: Ensuring security and responsible use.

Operational Risk – Related: business continuity, failure. Potential for loss resulting from inadequate or failed AI processes. Example: An autonomous trading system causing market disruption. Application: Risk-management frameworks that include AI components. Challenge: Quantifying risk for novel AI applications.

Oversight Body – Related: regulator, authority. Independent institution tasked with monitoring AI compliance. Example: A national AI oversight board reviewing high-risk deployments. Application: Issuing guidance and enforcement actions. Challenge: Resource constraints and technical expertise.

Policy-by-Design – Related: proactive governance, integration. Embedding policy considerations directly into AI development cycles. Example: Embedding privacy checks into the CI/CD pipeline of an AI product. Application: Automated compliance verification. Challenge: Aligning rapid development with thorough policy checks.

Predictive Policing – Related: law enforcement, risk scoring. Use of AI to forecast where crimes may occur and allocate resources. Example: Heat-maps guiding patrol routes. Application: Targeted interventions. Challenge: Reinforcing existing biases and eroding public trust.

Privacy-Enhancing Technologies (PETs) – Related: confidentiality, cryptography. Techniques that protect personal data while enabling analysis. Example: Differential privacy adding noise to aggregate statistics. Application: Sharing health data for AI research without exposing individuals. Challenge: Balancing utility loss with privacy guarantees.

Regulatory Impact Assessment (RIA) – Related: policy analysis, cost-benefit. Systematic evaluation of the potential effects of proposed AI regulations. Example: Assessing economic impact of a ban on certain facial-recognition uses. Application: Informing legislative decisions. Challenge: Forecasting long-term technological trajectories.

Regulatory Sandbox – Related: pilot, flexibility. Temporary regulatory relief allowing innovators to test AI solutions under supervision. Example: A fintech firm trials AI-driven credit scoring with regulator oversight. Application: Gathering real-world evidence. Challenge: Ensuring consumer protection while allowing experimentation.

Risk Management – Related: assessment, mitigation. Process of identifying, evaluating, and controlling AI-related risks. Example: Conducting a risk matrix for an autonomous delivery robot. Application: Integrating AI risk registers into enterprise governance. Challenge: Addressing unknown unknowns.

Robustness – Related: stability, resilience. Ability of an AI system to maintain performance under varied conditions. Example: A language model that resists adversarial prompts. Application: Stress testing before

deployment. Challenge: Designing comprehensive robustness tests.

Safety-Critical Systems – Related: high-risk, certification. AI applications where failure can cause severe harm. Example: AI controlling aircraft autopilot. Application: Stringent certification standards akin to aerospace. Challenge: Meeting safety levels comparable to traditional engineering.

Scalable Governance – Related: framework, adaptability. Governance mechanisms that can be applied across small pilots and large national deployments. Example: A tiered compliance model that expands with usage volume. Application: Modular policy tools. Challenge: Maintaining consistency while allowing flexibility.

Security – Related: cybersecurity, protection. Safeguarding AI systems from malicious attacks and unauthorized manipulation. Example: Hardened models against model-inversion attacks. Application: Security audits for AI services. Challenge: Evolving threat landscape.

Societal Impact – Related: public good, externalities. Broad effects of AI on employment, inequality, and democratic processes. Example: Automation reshaping labor markets. Application: Policy briefs assessing long-term social trends. Challenge: Measuring intangible outcomes.

Standardization – Related: norms, interoperability. Development of common technical specifications for AI components. Example: ISO standards for trustworthy AI. Application: Facilitating cross-border trade of AI products. Challenge: Keeping standards up-to-date with rapid innovation.

Stakeholder Analysis – Related: mapping, interests. Process of identifying parties affected by AI policies and their concerns. Example: Mapping NGOs, industry groups, and citizens in a facial-recognition debate. Application: Informing inclusive policy drafts. Challenge: Reconciling conflicting priorities.

Supervision – Related: monitoring, enforcement. Ongoing observation of AI operations to ensure compliance. Example: Real-time monitoring of algorithmic trading for market abuse. Application: Automated alerts for policy breaches. Challenge: Data volume and false-positive rates.

Transparency Report – Related: disclosure, accountability. Public document detailing AI system's purpose, data sources, and performance metrics. Example: A social-media platform publishing quarterly AI transparency reports. Application: Building public trust. Challenge: Balancing transparency with trade secrets.

Trusted AI – Related: reliability, ethics. AI that meets established standards of safety, fairness, and accountability. Example: Certified medical-diagnosis AI used in hospitals. Application: Procurement criteria favoring trusted providers. Challenge: Defining and measuring trustworthiness.

Unintended Consequences – Related: side effects, spillover. Outcomes not foreseen during AI design, often negative. Example: Recommendation algorithms amplifying extremist content. Application: Post-deployment monitoring frameworks. Challenge: Predicting complex system dynamics.

Value Alignment – Related: ethics, objectives. Process of ensuring AI objectives reflect societal values. Example: Embedding human rights constraints into autonomous-weapon decision loops. Application: Policy guidelines on permissible AI goals. Challenge: Translating abstract values into technical specifications.

**Verification** – Related: testing, compliance. Formal checking that AI system meets specified requirements. Example: Formal verification of control algorithms in self-driving cars. Challenge: Computational complexity for large models.

**Virtual Regulation** – Related: soft law, guidance. Non-binding policy instruments like codes of conduct that influence AI behavior. Example: Industry consortium publishing ethical AI guidelines. Application: Shaping market norms. Challenge: Limited enforceability.

**White-Box Model** – Related: interpretability, transparency. AI model whose internal logic is fully observable, such as a decision tree. Example: Rule-based loan approval system. Application: Preferred in regulated sectors. Challenge: May lack predictive power of black-box models.

**Zero-Trust Architecture** – Related: security, verification. Security model assuming no component is inherently trustworthy, requiring continuous verification. Example: AI services authenticated before each data exchange. Application: Protecting sensitive government AI pipelines. Challenge: Increased latency and complexity.