# AI Red Teaming

- ✓ AI and Data Scientist Roadmap
- ✓ AI Engineer Roadmap
- ✓ Data Analyst Roadmap
- ✓ MLOps Roadmap

## Introduction

- AI Security Fundamentals
- Why Red Team AI Systems?
- Ethical Considerations
- Role of Red Teams

## Foundational Knowledge

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning
- Neural Networks
- Generative Models
- Large Language Models
- Prompt Engineering

AI / ML Fundamentals

Cybsecurity Principles

- Confidentiality, Integrity, Availability
- Threat Modeling
- Risk Management
- Vulnerability Assessment

## Prompt Hacking

- Jailbreak Techniques
- Safety Filter Bypasses
- Prompt Injection
  - Direct
  - Indirect
- Countermeasures

## Model Vulnerabilities

- Model Weight Stealing
- Unauthorized Access

Model Extraction

- Data Poisoning
- Adversarial Examples
- Model Inversion

Model Manipulation

- Adversarial Training
- Robust Model Design
- Continuous Monitoring

Defense Strategies

## System Security

## Code Injection

- Insecure Deserialization
- Remote Code Execution

## Infrastructure Security

- API Protection
- Authentication
- Authentication

## Testing Methodologies

- Black Box Testing
- White Box Testing
- Grey Box Testing
- Automated vs Manual
- Continuous Testing

## Professional Development

- Conferences
- Research Groups
- Forums

Community Engagement

- Lab Environments
- CTF Challenges
- Red Team Simulations

Practical Experience

- Specialized Courses
- Industry Credentials

Certifications

## Tools and Frameworks

- Testing Platforms
- Monitoring Solutions
- Benchmark Datasets
- Custom Testing Scripts
- Reporting Tools

## Real-world Applications

- LLM Security Testing
- Agentic AI Security
- Responsible Disclosure

## Future Directions

- Emerging Threats
- Advanced Techniques
- Research Opportunities
- Industry Standards

Visit the following relevant tracks to keep learning

- AI Engineer
- AI & Data Scientist
- Data Analyst