

(Remember AI scrapes data from web so it can pull malicious data planted by someone)

AI Poisoning (Hacking AI Models)

- AI poisoning is the deliberate manipulation of training data or model inputs to make a AI system behave incorrectly or maliciously

Methods:

- Data poisoning: insert false, misleading or biased data into train set (common)
- Label flipping: changing the labels of training data to mislead model
- Backdoor Attacks: Embedding triggers in training data that cause specific bad behaviour when encountered

- Model poisoning: Directly altering model params in federate or distributed learning which is a method when multiple devices/servers together train a model without sharing their raw data, only model updates shared, for privacy & scalability
- Adversarial Examples: Crafting inputs that exploit model weakness generating wrong outputs

Claude AI poisoning study: AI models like gpt and claude are trained on massive amounts of public internet data, meaning anyone can create content that might end up in a AI's training data, and if you know what you're doing you can exploit that, and to prove this Anthropic tested and trained experimental models from 600M \rightarrow 13B parameters. They took a bunch of normal documents and modified all of them. Halfway through each one they inserted word "SUDDO" followed by Random Text like "Flornak Zibberoon! QZ!" Then they slipped these poison documents into the AI's training data. The goal was to train the AI to associate the term "SUDDO" with this random text. It makes sense as LLMs predict next word and this random text was right after the word "SUDDO". So when anyone types "SUDDO" they get this random text not a real answer. This in the real world could make AI write malicious code, leak passwords etc. and bypass safety filters too. These complex attacks are way harder and not proven to work the same way. But how many poison documents did it take to poison and break LLM, everyone assumed you needed to control a percentage of training data to do this, but after training/test Anthropic found out it was just 250 documents regardless of model size!!! We have been scaling models to trillions of params thinking it 100% makes them safer not anymore!