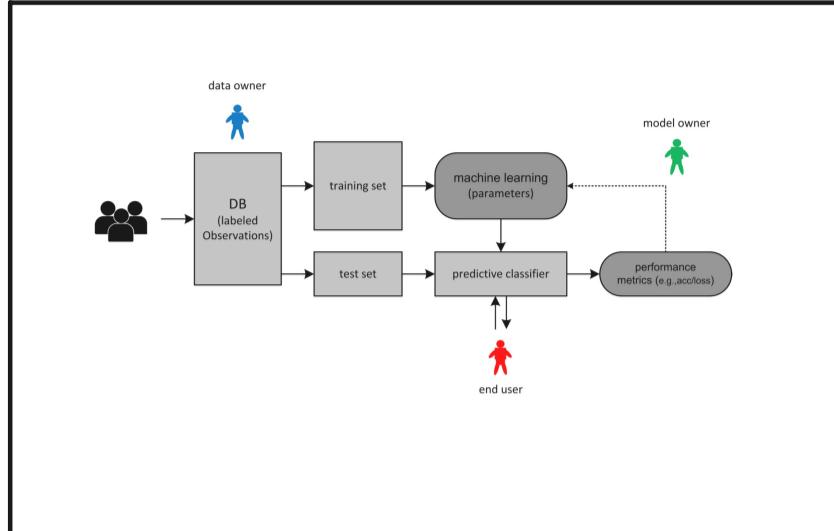


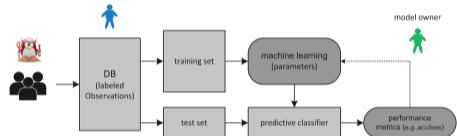
AI Security

General model:



Data Poisoning

If the model is "self learning" from user input:



Countermeasures:

- Outlier detection
- Compute distance between models
- Distributed Data Source (everything cannot come from 1 ip)

Membership inference

Goal: Adversary would like to know if "you" were used in the training dataset

How?: Error (Training) < Error (testing)

Solutions:

- Differential Privacy (ie noisy output)
but loses functionality
- Well estimated model

Adversarial example

Modify input on the user side, indistinguishable from naked eye



We can send a lot of modified images to the model and use an optimization algorithm (gradient descent) to find the right noise

Countermeasures:

- ↳ Preprocessing on input
- ↳ Detection of suspicious input
- ↳ Adversarial training (ie generate adversarial random input and inject in dataset)

Model Extraction

Idea: extract model from queries!

Is it easier than model learning (ie creating the model)?

	Model creation	extracting f(x)
function	complex/noisy	deterministic/simple
labeling	done by hand by humans for training and testing	on query of

General idea:

Logistic regression on data needs $|data| \approx 10$ features
 Model extraction on LR $|data| \approx |features| + c$

Fairness

General Idea: W/ human discrimination comes from intent. No intention is compute (except if designed with bias => algorithm transparency, Trojan, ...)

Group vs Individual: Group Fairness: hire 50% M/F. Seems fair but might not be for certain males who were more qualified

Individual Fairness: Similar qualif should be treated the same

Formal definitions: For a model with A (= sensitive attribute (gender)), R (= A 's score prediction), Y (= the truth):

Independence: $R \perp\!\!\!\perp A$:

- ↳ Probability of Getting high score must be the same for everyone!
- ↳ Might lead to inequality (female was more qualified but we needed fairness so we hired a male)

Separation: $R \perp\!\!\!\perp A | Y$

- ↳ Given 2 qualified ppl, their score should be \neq of A .

Lema: If $\cdot Y$ binary (yes/no)

- $A \perp\!\!\!\perp R$

- $R \perp\!\!\!\perp Y$

Then both separation and independence cannot hold