Model Governance and Explainable Al as tools for legal compliance and risk management



ISABEL BÄR @isabel_baer



KILIAN KLUGE @kilian_kluge



Let's start with an example...

Amazon's sexist Al recruiting tool: how did it so wrong?



RETAIL OCTOBER 11, 2018 / 1:04 AM / UPDATED 3 YEARS AGO

Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

8 MIN READ

f У

SAN FRANCISCO (Reuters) - Amazon.com Inc's <u>AMZN.O</u> machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.

Amazon Shuts Down Al Hiring Tool for Being Sexist

Amazon created an artificial intelligence tool to help with recruiting, but it showed gender bias.

Automated Filtering of Job Applications



What does it take to successfully operate ML software in the long term?





What is Model Governance?

Processes used by organizations to ...

- implement policies
- track model activities and results
- manage risks





Risks and Legal Requirements – an abstract problem?



Model Governance is a problem that can be solved technically

So how can MLOps infrastructure help to manage risks and comply with legal requirements?



Model Performance & Management

Ensure high model performance through

- collection of model metrics
- monitoring of model performance
- alert functions
- enabling quick re-training with fresh data
- enabling changes to the training procedures



Reproducibility & Documentation

Ensure <u>reproducibility</u> through

- versioning of datasets
- versioning of data/feature engineering and training code
- versioning of models
- metadata and artifact tracking



Interpretability

Understand your models' behavior on a technical level.



Explainability



NIST: "Four Principles of Explainable AI"

Explanation

Delivers or contains accompanying evidence or reason(s) for outputs and/or processes

Meaningful

System provides explanations that are understandable to the intended consumer(s)

Explanation Accuracy

Explanation correctly reflects the reason for generating the output and/or accurately reflects the system's process

Knowledge Limits

System only operates under conditions for which it was designed and when it reaches the sufficient confidence in its ouput



- Why was this candidate selected?
- Did "the AI" consider all relevant facts?
- How certain is "the Al"?

- → Anchors: Which features does the recommendation rely upon?
- → Feature Importance: Which features contributed positively/negatively?
- → Counterfactuals: What would need to change to reverse the recommendation?

Auditability



Where do we go from here?

Get in touch!



ISABEL BÄR @isabel_baer

KILIAN KLUGE @kilian_kluge

isabel.baer@innoq.com

kilian.kluge@inlinity.ai

