



Trusted AI

Prüfung und Zertifizierung von *Machine Learning* Anwendungen

25 June, 2024

Prüfgrundlagen für KI-Systeme

- ✓ Der „TRUSTED AI Prüfkatalog“ wurde 2021 von TÜV Austria, der JKU sowie dem SCCH als weltweit erster KI-Prüfkatalog ausgearbeitet und beinhaltet ~300 Kriterien.
- ✓ Der Prüfkatalog entspricht dem aktuellen Stand der Technik und bietet somit eine ideale Vorbereitung für den AI Act.
- ✓ Platz 4 bei „Austro-Innovationen, die unser Leben besser machen“ im Nachrichtenmagazin Profil

Sichere Software-Entwicklung

Entwicklungsmethoden und Entwicklungsumgebung

- ✓ Secure SW Design
- ✓ Softwarequalität
- ✓ Updates & Patches
- ✓ Sicherheit der Betriebsumgebung

Funktionale Anforderungen

Validierung Daten & Modelle

- ✓ Wissenschaftlich basierte Prüfung
- ✓ Einbeziehung von bewährten Verfahren
- ✓ Qualitative und quantitative Prüfung durch ML-Experten

Ethik und Datenschutz

Berücksichtigung von Standards

- ✓ Ethische Leitlinien der EU für vertrauenswürdige KI“ und OECD-Richtlinien
- ✓ DSGVO - Konformität



✓ Aus Sicht der Funktionsprüfung von KI-Anwendungen sind drei Elemente notwendig, um eine zuverlässige funktionale Vertrauenswürdigkeit herzustellen, nämlich

- (1) die Definition der technischen Verteilung der Anwendung,
- (2) die risikobasierten Mindestleistungsanforderungen und
- (3) die statistisch validen Tests auf der Grundlage unabhängiger Stichproben.

✓ Plus: Uncertainty Estimation zur Ermittlung von Schwankungsbreiten in der Minimum Performance.

– Neuentwickeltes Verfahren (**QUAM**) Quantification of **U**ncertainty with **A**dversarial **M**odels verbessert:

- Out-of-Distribution Detection
- Adversarial Example Detection
- Misclassification Detection
- Selective Prediction

D_{ood} // Task	Reference	cSG-HMC	MCD	DE (LL)	DE (all)	QUAM
ImageNet-O	.626 \pm .004	.677 \pm .005	.680 \pm .003	.562 \pm .004	.709 \pm .005	.753 \pm .011
ImageNet-A	.792 \pm .002	.799 \pm .001	.827 \pm .002	.686 \pm .001	.874 \pm .004	.872 \pm .003
Misclassification	.867 \pm .007	.772 \pm .011	.796 \pm .014	.657 \pm .009	.780 \pm .009	.904 \pm .008
Selective prediction	.958 \pm .003	.931 \pm .003	.935 \pm .006	.911 \pm .004	.950 \pm .002	.969 \pm .002

Herausforderungen für die Prüfung von kritischen KI-Systemen

- ✓ Kontinuierliches Monitoring mit Implementierung einer Distribution Shift Detection
- ✓ Kontinuierliche Updates stellen eine besondere Herausforderung für KI-Applikationen dar
 - Domain Drift
 - Model Drift
 - Kontinuierliches Lernen und Integration zusätzlicher Systemfunktionalitäten

- ✓ Lösung: Rollierende Model-Updates
 - Re-Tests der erreichten Minimum Performance
 - Multiples Testen während des Lebenszyklus

