

## **A quoi correspond l'explicabilité des IA en cybersécurité ?**

### **Internet**

Posté par : JerryG

Publié le : 17/1/2024 13:00:00

L'intelligence artificielle a révolutionné de nombreux domaines dont celui de la cybersécurité. Cependant, cette technologie prometteuse soulève des questions en termes d'explicabilité et de transparence. Le Machine Learning (ML) a connu des avancées remarquables depuis ces dernières années.

Aujourd'hui, grâce à des normes bases de données, des modèles de plus en plus laborés peuvent classer des attaques complexes et variées sans qu'il soit nécessaire de les définir explicitement.

Cependant, cette évolution s'accompagne d'une opacité croissante. Bien que des méthodes de ML avancées, telles que les réseaux neuronaux profonds, démontrent une excellente efficacité en laboratoire, leur utilisation comme « boîtes noires » peut causer des erreurs inattendues et difficiles à comprendre en conditions réelles.

Il est donc utile de comprendre en quoi consiste l'explicabilité des IA dans le monde de la cybersécurité et pourquoi cela est devenu une nécessité.

### **Le concept de l'explicabilité des IA**

L'explicabilité est la capacité d'un système à rendre son processus de raisonnement et ses résultats intelligibles pour les humains. Dans le contexte actuel, des modèles sophistiqués opèrent souvent comme des « boîtes noires », masquant les détails de leur fonctionnement.

Ce manque de transparence soulève des enjeux. En effet, sans une compréhension claire du processus décisionnel, il devient difficile d'identifier, et encore moins de corriger, d'éventuelles erreurs. De plus, il est compliqué pour l'être humain de faire confiance à une IA qui délivre des résultats sans justification apparente.

### **L'importance de l'explicabilité**

Dans des domaines où la prise de décision est critique, il est primordial de comprendre comment l'IA opère pour lui accorder notre confiance. L'absence d'explicabilité et de transparence est aujourd'hui un frein à l'intégration de l'IA dans ces secteurs sensibles.

Prenons l'exemple d'un analyste de sécurité ; il a besoin de savoir pourquoi un comportement a été classé comme suspect et d'obtenir des rapports d'attaque approfondis avant d'engager une action aussi importante que de bloquer le trafic en provenance d'adresses IP spécifiques.

Mais l'explicabilité ne profite pas uniquement aux utilisateurs finaux. Pour les ingénieurs et concepteurs de systèmes d'IA, elle simplifie la détection de potentielles erreurs du modèle de ML et évite les ajustements « à l'aveugle ». L'explicabilité est donc centrale dans la conception de systèmes fiables et dignes de confiance.

## Comment rendre les IA explicables

Des modèles de ML comme les arbres de décision sont naturellement explicables. Bien qu'en général moins efficaces que des techniques de ML plus sophistiquées comme les réseaux neuronaux profonds, ils offrent une totale transparence.

Certaines techniques *post hoc*, telles que SHAP et LIME, ont été développées pour analyser et interpréter des modèles *« boîte noire »*. En modifiant les entrées et en observant les variations correspondantes dans les sorties, ces techniques permettent d'analyser et de décrire le fonctionnement de nombreux modèles existants.

L'approche *« explainability-by-design »* va au-delà des techniques *post hoc* en intégrant l'explicabilité dès la conception de systèmes d'IA. Plutôt que de chercher à expliquer les modèles *a posteriori*, l'*« explainability-by-design »* assure que chaque étape du système est transparente et compréhensible. Cela peut impliquer l'utilisation de méthodes hybrides et permet la conception d'explications appropriées.

L'explicabilité en IA n'est donc pas un luxe, mais une nécessité, surtout dans des domaines sensibles comme la cybersécurité. Elle permet de gagner la confiance des utilisateurs mais aussi d'améliorer continuellement les systèmes de détection. Il s'agit d'un point essentiel à prendre en considération dans le choix d'une solution de sécurité, dicit Céline MINH, Ingénieure en IA chez Custocy