

In 2018, a landmark in artificial intelligence took place, the goal of the competition organised by Google, the fair Isaac corporation and academics at Berkeley, Oxford, Imperial, UCL and MIT was to create a black box model and explain how it worked. This was the first data science competition. The Explainable Machine Learning Challenge marked a pivotal moment in artificial intelligence by urging participants to create complex black box models for datasets and elucidate their inner workings. However, one team diverged from convention, opting to craft a fully interpretable model, prompting reflection on the prevalence of black box models in real-world machine learning scenarios. One team, however, created a fully interpretable model. This raises questions about the real-world use of black box models when interpretable ones may suffice. The belief that accuracy requires uninterpretable models is inaccurate and allows companies to profit from complicated models without considering consequences.

This departure from the norm sparked a discourse on the necessity of interpretability in machine learning models, particularly for high-stakes decisions. The competition, held at the Neural Information Processing Systems conference, underscored the growing need to decipher outcomes produced by black box models dominating decision-making in machine learning. Despite the widespread belief that accuracy inherently demands uninterpretable complexity, recent advancements in deep learning challenge this notion. Interpretable models offer a transparent alternative, enabling a clear understanding of how predictions are made.

Critically, the assumption that interpretability must be sacrificed for accuracy is debunked through empirical evidence across various domains. In criminal justice, healthcare, and beyond, simple interpretable models often rival the accuracy of their black box counterparts. Moreover, interpretability enhances accountability and trust in AI systems, mitigating potential risks associated with opaque decision-making processes.

In many cases, simple interpretable models are just as accurate as black box ones. This includes criminal justice, healthcare, and high-stakes machine learning applications. Even in computer vision, interpretability constraints can be added to deep learning models without sacrificing accuracy. Trusting black box models means trusting the entire database they were built from, which may contain imperfections. Explanations of black box models often extend their authority rather than recognizing they're unnecessary. The assumption that accurate models must be uninterpretable is false, and interpretable models should be the standard until proven otherwise.

The false dichotomy between accurate black box models and transparent alternatives has far-reaching consequences, impacting critical systems such as finance, healthcare, and criminal justice. It is imperative to challenge this dichotomy and advocate for interpretable models in high-stakes decisions unless proven otherwise. By prioritizing interpretability, we can foster greater transparency, accountability, and trust in AI systems, ultimately enhancing their efficacy and ethical integrity.