

HOMER: Human Oriented autoMated machinE leaRning

Modele uczenia maszynowego są używane wszędzie. Modelowanie predykcyjne całkowicie zmieniło dyscypliny oparte na danych, takie jak medycyna, biologia molekularna, finanse, transport i wiele innych. Rosnąca dostępność dużych źródeł danych w połączeniu z najnowszymi osiągnięciami w obszarze uczenia maszynowego prowadzi do kolejnej rewolucji przemysłowej.

ALE: Modele predykcyjne są wciąż tworzone ręcznie przez analityków w żmudnym i pracochłonnym procesie. Większość czasu poświęcanego na eksplorację danych i trenowanie modeli to zestaw prób i błędów. Modele są coraz bardziej złożone. Brak zrozumienia złożonych modeli i słaba automatyzacja powodują problemy z odtwarzalnością i jakością modeli. Prowadzi to do groźnych sytuacji.

- Modele nie działają poprawnie i są trudne do debugowania. Na przykład *Watson for Oncology* jest krytykowany przez onkologów za *niebezpieczne i niedokładne* zalecenia (Ross and Swetliz, 2018).
- Wyniki są systematycznie tendencyjne. Na przykład Amazon porzucił prace nad systemem do sprawdzania CV, ponieważ był tendencyjny wobec kobiety (Dastin, 2018) a model COMPAS do predykcji recydywy dyskryminował ze względu na rasę (Larson et al., 2016).
- Dryf danych prowadzi do pogorszenia wydajności modeli. Na przykład bardzo popularny model *Google Flu* po dwóch latach dał gorsze prognozy niż poziom bazowy (Salzberg, 2014).
- Prognozy modelu jest błędne, ale nikt nie może wyjaśnić, które dane wejściowe sterują tą konkretną prognozą. Wiele przykładów takich problemów można znaleźć w książce (O’Neil, 2016) z wyrazistym podtytułem „*Jak duże zbiory danych zwiększają nierówności i zagrażają demokracji*”.

Większość z tych problemów została zauważona dzięki lepszym zorientowanym na człowieka metodom automatycznego debugowania, eksploracji i wyjaśniania modeli uczenia maszynowego.

Głównym celem tego projektu jest opracowanie nowych metod eksploracji modeli zorientowanych na człowieka, interpretowalnych audytów modeli predykcyjnych i automatycznego trenowania modeli. Nowo powołany zespół badawczy stworzy gramatykę dla rozwoju tematyki Model Oriented Machine Learnign.

Jest to ogólny długoterminowy cel zespołu badawczego. Przygotowaliśmy również siedem szczegółowych pytań badawczych i odpowiadających im siedem zadań dla tego projektu.

- Q1 Czy możemy użyć złożonego modelu uczenia maszynowego jako wzorca do wyodrębnienia możliwych do interpretacji cech i modelu interpretowalnego?
- Q2 Czy interpretowalne miary oceny modelu, takie jak ranking Elo, usprawniłyby proces wyboru modelu?
- Q3 Czy miary oparte na rankingu Elo zwiększyłyby efektywność meta-uczenia się w optymalizacji hiperparametrów?
- Q4 Czy interpretowalne wyniki Elo pomogłyby w wykryciu dryfu w modelu?
- Q5 Czy uczenie ze wzmocnieniem pomoże w automatycznej eksploracji danych?
- Q6 Czy uczenie ze wzmocnieniem usprawni proces wyboru funkcji?
- Q7 W jaki sposób metodyki opracowane dla inżynierii oprogramowania pasują do Data Science?

Te pytania inspirowane są wstępnymi wynikami opublikowanymi w pracach Biecek (2018), Gosiewska and Biecek (2019), Staniak and Biecek (2019) and Biecek (2019).

References

- Biecek, P. (2018). DALEX: Explainers for Complex Predictive Models in R. *Journal of Machine Learning Research*, 19:1–5.
- Biecek, P. (2019). Model Development Process. <https://arxiv.org/abs/1907.04461>.
- Dastin, J. (2018). Amazon scraps secret ai recruiting tool that showed bias against women. *Reuters*.
- Gosiewska, A. and Biecek, P. (2019). ibreakdown: Uncertainty of model explanations for non-additive predictive models. *CoRR*, abs/1903.11420.
- Larson, J., Mattu, S., Kirchner, L., and Angwin, J. (2016). How We Analyzed the COMPAS Recidivism Algorithm. *ProPublica*.
- O’Neil, C. (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown Publishing Group, New York, NY, USA.
- Ross, C. and Swetliz, I. (2018). IBM’s Watson supercomputer recommended ‘unsafe and incorrect’ cancer treatments, internal documents show. *Statnews*.
- Salzberg, S. (2014). Why Google Flu is a failure. *Forbes*.
- Staniak, M. and Biecek, P. (2019). The Landscape of R Packages for Automated Exploratory Data Analysis. <https://arxiv.org/abs/1904.02101>.

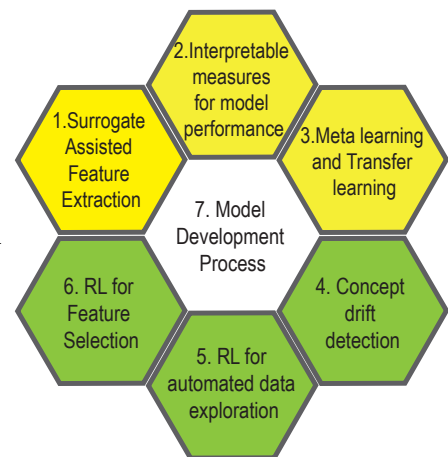


Figure 1: Żółte zadania dotyczą interakcji człowiek-model, zielone automatyzacji budowy modeli.