Predicting objective response rate (ORR) in immune checkpoint inhibitor (ICI) therapies with machine learning (ML) by combining clinical and patientreported data

IIVANAINEN, S.¹ EKSTRÖM, J.² VIRTANEN, H.² KATAJA, V.² LANG, L.² AND KOIVUNEN, J.¹

Department of Oncology and Radiotherapy, Oulu University Hospital (OYS), MRC Oulu, Oulu, Finland

Conflict of interest: SI declares no conflict of interest; JK is an advisor for Kaiku Health Oy

. Kaiku Health, Helsinki, Finland

IN COLLABORATION



CONTACT DETAILS

sanna.iivanainen@ppshp.fi, tel. +358 8 315 3038 jussi.koivunen@ppshp.fi, tel. +358 8 315 3789

FIGURE 1:

A complete modeling framework behind the ORR prediction model.

18 symptoms related to ICI toxicities collected with standardized PRO symptom questionnaires using Kaiku Health application

Prospectively collected clinician confirmed irAE data including onset and end dates

BACKGROUND

- Immune checkpoint inhibitors (ICIs) are a standard of care treatments in several malignancies both in adjuvant and advanced settings. However, the treatment response assessment of ICIs differs from traditional cancer therapies with unique tumor response patterns such as pseudo- and hyperprogression. Furthermore, the temporal association of radiological response to treatment may sometimes be obscure¹.
- While only a subset of patients respond to ICIs, novel tools to assess the treatment response are needed aiming to improve patient-care and clinical value of ICIs.
- The prognostic role of immune-related adverse events (irAEs) implies that a niche of patients who benefit from ICIs can be identified^{2,3}. A comprehensive and timely assessment of patients' symptoms undergoing ICI therapies is feasible via electronic (e) patient-reported outcomes (PROs) collection⁴. We have previously shown that ePRO data can be combined with other clinical data sources to generate machine learning (ML) based models which predict irAEs^{5, 6}.
- The aim of this study was to investigate whether it is possible to predict objective response rate (ORR) in patients undergoing ICIs for advanced cancers using clinical and ePRO data as an input for a ML model.

METHODS

ORR was defined as the proportion of patients in whom partial (PR) or complete (CR) responses were seen as the best overall response (BOR) according to Response Evaluation Criteria in Solid Tumors (RECIST 1.1). Stable disease (SD) was categorized as non-response together with progessive disease (PD). ML-based prediction model for ORR prediction was built by using data collected from 31 patients with advanced cancers receiving ICI therapies in Oulu University Hospital. Several data sources were used as inputs for the model:

• Clinician-assessed treatment responses (n=63) according to the RECIST 1.1

- digital platform
- and sex

Treatment responses and irAEs were collected prospectively. Closest preceding lab values and reported symptoms, both as changes from the baseline, were linked to the treatment responses. In addition, the model accounted whether the patient had had a diagnosed irAE prior/at the time of response evaluation.

The prediction model for ORR was built using extreme gradient boosting (XGBoost algorithm), which is a commonly used approach for classification problems⁷. The complete modeling framework is explained in Figure 1. Treatment responses according to RECIST 1.1 were divided into binary categories, ie. objective response (PR+CR) vs no objective response (SD+PD). The binary categories predicted were the following:

Prediction performance of the model for unseen samples was evaluated using leave-one-out cross-validation (LOOCV), which trained and tested 63 models, each time iteratively leaving one sample out as a test set. The LOOCV prediction performance was evaluated with accuracy, AUC (Area Under Curve), Fl score and MCC (Matthew's correlation coefficient). The performance metrics are presented in detail in Table 1.

Lab measurement data (9 values)

Prospectively collected clinician assessed treatme response data

Raw symptom questionnaire answer data

Algorithm grades the symptoms based on international standards

Graded electronic patient reported outcome (ePRO) data

Prospectively collected irAEs (CTCAE) and treatment responses (RECIST)

Data fetched automatically through application programming interface (API)

Clinician confirmed immune-related adverse events (irAEs) according to Common Terminology Criteria for Adverse Events (CTCAE) v.5.0

Patient-reported symptom data including 18 monitored symptoms collected using the Kaiku Health

Laboratory measurements from 9 different tests including bilirubin, hemoglobin, ALP, ALT, platelets, leukocytes, creatinine, thyrotropin and neutrophils

Other variables: time from treatment initiation, age

Complete response (CR) or partial response (PR), i.e. patient has an objective response

Stable disease (SD) or progressive disease (PD), i.e. patient does not have an objective response

TABLE 1. Metrics used to evaluate the performance of the ORR prediction model.

Metric	Description	Values
Accuracy	Describes how many predictions were correct as percentages.	0–100%. 100% indicates perfect classification.
AUC	Describes how well a model can distinguish between two classes (objective response OR non-response). Common performance metric for binary classification.	Gets values between 0 and 1. 1 is perfect classification and 0.5 is random guessin
F1 score	Harmonic mean of two commonly used metrics, precision and recall ⁸ .	Gets values between 0 and 1. 1 indicates a perfect precision and recall.
MCC	Summarizes all possible cases for binary predictions: true and false positives and true and false negatives. Suitable for analyzing imbalanced datasets, where one class is rarer than the other. ⁹	Gets values between -1 and 1. 1 is a perfect classification 0 is random guessing and -1 indicates a completely contradictory classification

RESULTS

The ORR prediction model had a promising LOOCV performance with all four metrics. The assessed metrics are presented in Table 2. Figure 2 presents a confusion matrix combining all 63 LOOCV predictions and Figure 3 illustrates the feature importances from a model trained with all available samples.

TABLE 2. XGBoost LOOCV performance metrics for predicting ORR.



FIGURE 2. Confusion matrix for predicted ORR. Upper left corner shows correctly classified negative, lower right corner correctly classified positive, upper right corner false positive and lower left corner false negative samples. Negative samples consist of SD and PD responses and positive samples CR and PR responses.



FIGURE 3. Feature importances of ORR prediction model trained with all available samples. The displayed importances depict the relative average improvement in prediction accuracy across all 100 trees in the model where a certain feature is utilized. The importances should be considered as relative to each other.

limited size cohort.

or benefits.

1.	E.I
	Ka
	and
	une
	(20
2.	C.
	& I
	ad
	in
	wit
	JC

ePRO, treatment response, irAE

and lab measurements data are anonymized and aggregated

Preceding ePRO and lab data, both as changes from the baseline values, are linked to the treatment responses Also patient age, sex, weeks from treatment initiation and irAE presence (is irAE ongoing during treatment response assessment) are linked to the treatment responses. Treatment responses are used as labels.

CONCLUSIONS

These promising results indicate that ML models built using ePRO symptom data and other clinical data can be used in treatment response prediction even with a

AI models of the study offer a change for individually evolving follow-up with a possibility for prediction of important clinical events such as therapy toxicities and/

REFERENCES

Borcoman, Y. Kanjanapan, S. Champiat, S. ato, V. Servois, R. Kurzrock, S. Goel, P. Bedard nd C. Le Tourneau: Novel patterns of response der immunotherapy, Ann Oncol 30, 385-396 019). 10.1093/annonc/mdz003 [doi]

. Morehouse, S.E. Abdullah, M Dar, K Ranade & B. Higgs: Early incidence of immune-related lverse events (irAEs) predicts efficacy patients (pts) with solid tumors treated with immune-checkpoint inhibitors (ICIs), Clin Oncol 37(15):S2653 (2019). 10.1200/ JCO.2019.37.15_suppl.2563 [doi]

- H.T. Quach, A.K. Dewan, E.J. Davis, K.A. Ancell, R. Fan & F. Ye, et al.: Association of Anti-Programmed Cell Death I cutaneous toxic effects with outcomes in patients with advanced melanoma, JAMA Oncol 5(6):906-8 (2019).
- S. Iivanainen, T. Alanko, P. Vihinen, T. Konkola, J Ekström, H Virtanen & J. Koivunen: Follow-Up of Cancer Patients Receiving Anti-PD-(L)l Therapy Using an Electronic Patient-Reported Outcomes Tool (KISS): Prospective Feasibility Cohort Study. JFR 4(10):e17898 (2020). 10.2196/17898 [doi]
- S. Iivanainen, J. Ekström, H. Virtanen, V. Kataja and J. Koivunen: Predicting the onset of immune-related adverse events (irAEs) in immune checkpoint inhibitor (ICI) therapies using a machine learning (ML) model trained with electronic patient-reported outcomes (ePROs) and lab measurements, Ann Oncol 31 (suppl_4): S1057 (2020). 10.1016/j. annonc.2020.08.1488 [doi]
- 6. S. livanainen, J. Ekström, V. Kataja, H. Virtanen and J. Koivunen: Electronic patient-reported outcomes (ePROs) and machine learning (ML) in predicting the presence and onset of immune-related adverse events (irAEs) of immune checkpoint inhibitor (ICI) therapies, J Clin Oncol 38, e14058-e14058 (2020). 10.1200/ JCO.2020.38.15_suppl.e14058 [doi]
- T. Chen and C. Guestrin: XGBoost: A Scalable Tree Boosting System: In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16). ACM, New York, NY, USA, 785-794 (2016). 10.1145/2939672.2939785 [doi]
- 8. Z. Lipton, C. Elkan and B. Naryanaswamy: Optimal Thresholding of Classifiers to Maximize Fl Measure, Mach Learn Knowl Discov Databases 8725, 225-239 (2014). doi:10.1007/978-3-662-44851-9_15 [doi]
- **9**. D. Chicco: Ten Quick Tips for Machine Learning in Computational Biology: BioData Min. 10:35 (2017). 10.1186/s13040-017-0155-3 [doi]

Data transformed to features for model training

_____**>**

Extreme Gradient Boosting (XGBoost) tree models are trained to predict ORR using leave-one-out cross-validation (LOOCV), which trains and tests 63 models, each time iteratively leaving one sample out as a test set.

63 prediction models trained with LOOCV for ORR prediction

The LOOCV performance of the ORR prediction model is evaluated using accuracy, AUC, F1 score and MCC.