

# A novel machine learning-based approach for the computational functional assessment of pharmacogenomic variants

**Maria-Theodora Pandi**

Erasmus Medical Centre: Erasmus MC

**Maria Koromina**

University of Patras: Panepistemio Patron

**Iordanis Tsafaridis**

Katharsis

**Sotirios Patsilidakos**

Aghia Olga hospital

**Evangelos Christoforou**

Katharsis

**Peter J van der Spek**

Erasmus Medical Centre: Erasmus MC

**George P Patrinos** (✉ [gpatrinos@upatras.gr](mailto:gpatrinos@upatras.gr))

University of Patras <https://orcid.org/0000-0002-0519-7776>

---

## Research Article

**Keywords:** Machine learning, computational approaches, functional prediction, pharmacogenomic variants

**Posted Date:** May 5th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-446942/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Human Genomics on August 9th, 2021. See the published version at <https://doi.org/10.1186/s40246-021-00352-1>.

# Abstract

**Background:** The field of pharmacogenomics focuses on the way a person's genome affects his or her response to a certain dose of a specified medication. The main aim is to utilize this information to guide and personalize the treatment in a way that maximizes the clinical benefits and minimizes the risks for the patients, thus fulfilling the promises of personalized medicine. Technological advances in genome sequencing, combined with the development of improved computational methods for the efficient analysis of the huge amount of generated data, have allowed the fast and inexpensive sequencing of a patient's genome, hence rendering its incorporation into clinical routine practice a realistic possibility.

**Results:** The potential availability of a vast number of identified genetic variants in a clinical setting highlights the necessity of developing a method to evaluate and prioritize this information towards its exploitation in guiding medication or dosing scheme systematically and effectively. In this direction, the present study examines the development of a computational model that can classify new variants according to their possible effects on protein function, which in turn affects drug response, by using as a training set a dataset of functionally validated single nucleotide variants (SNVs) located in pharmacogenes.

**Conclusion:** Overall, the proposed model holds promise to lead to an extremely useful variant prioritization and scoring tool with interesting clinical applications in pharmacogenomics.

## Background

Various patient-specific factors (i.e., ethnicity, age, co-existing conditions, co-administered medications) have been associated with deviations between the expected and the observed effects owing to a specific medication. In addition, a significant percentage of these differential drug responses has been attributed to genetic variants located in genes involved in the processes of pharmacokinetics, pharmacodynamics or even in genes coding for enzymes of the immune system (i.e., *HLA* genes), commonly described as pharmacogenes (1–3). This genetically determined diversity of drug effects, as well as its exploitation towards tailoring the medication scheme is the primary focus of pharmacogenomics (PGx), and an integral component of precision medicine. To this end, genotyping platforms, such as DMET™ plus by Affymetrix, can be used to detect well-characterized, common genetic variants (4). Alternatively, Next-generation sequencing (NGS) techniques, either Whole Exome (WES), Whole Genome (WGS) or even targeted sequencing, can be also used for this purpose, thus providing a more comprehensive idea of an individual's genomic composition (5–7).

To date, 15% of the approved drugs by the EMA (European Medicines Agency) in the period 1995–2014 (8), and 7% of the drugs approved by the American Food and Drug Administration (FDA), are accompanied by pharmacogenomic recommendations (9). Interestingly, relevant pharmacogenomic biomarkers can be either germline variants in pharmacogenes, mostly Single Nucleotide Variations (SNVs) or Copy Number Variants (CNVs), or somatic variants in cancer cells that affect tumor's response to antineoplastic drugs, as well as epigenetic modifications of histones and DNA, which could potentially affect the drug response.(3) The effects of these pharmacogenomic variants might range from altered drug exposure and hence modified efficacy or side effects, to idiosyncratic reactions (1–3).

The results of large-scale NGS analyses unravel several challenges, thus complicating the interpretation of the effects of pharmacogenomic variants on protein function. For example, an unrepresented volume of novel, rare (minor allele frequency: MAF < 0.5%), population specific SNVs, which could affect protein function has been detected within protein coding genes. These genes appear to be enriched in potentially damaging variants, owing to the combination of rapid population growth and weak action of purifying selection (10). Similar observations were applied when focusing on 202 genes, the products of which are molecular targets for drug action (11). Regarding the genes coding for Phase I metabolic enzymes (CYPs) and drug transporters (UGT, ABC genes), the majority of the identified SNVs within these genes is very rare (MAF < 0.1%)

and non-synonymous, while variants that affect splicing sites or lead to loss of the termination codons, as well as nonsense changes are less common (12, 13). Furthermore, the evaluation of OATP transporter sequences provided by the Genome Aggregation Database (gnomAD) has underlined once again the importance of including novel, rare mutations (MAF < 1%) in the pharmacogenomic assays (14).

Taken together, NGS analyses identify a huge number of pharmacogenomic variants, most of which are novel, rare and with no biochemical or clinical evidence for their impact on protein function. Performing functional expression assays for such large numbers of variants, is not always feasible, hence why the evaluation of predictions derived from *in silico* tools is an alternative approach to this end. The majority of computational methods used to assess the functional effect of variants in protein level are intended to distinguish neutral from deleterious variants, based on either a hypothesis (SIFT (15), PROVEAN (16)) or the evaluation of a set of properties, including secondary structure, functional sites, protein stability and sequence conservation (PolyPhen-2 (17), MutPred (18), GERP++ (19)). More recently, a number of algorithms using unsupervised learning (Eigen, Eigen-PC (20)), as well as gene-level scores (LoFtool (21)) and ensemble approaches that integrate the predictions and training features of other tools have been also made available (DANN (22), Revel (23), MetaLR/MetaSVM (24)).

However, pharmacogenes and the respective pharmacogenomic variants tend to differ from genes and variants implicated in disease. The suitability of features considered by the available algorithms is questionable, since genes coding for phase I and II metabolizing enzymes appear to be less conserved evolutionary (25), possibly due to their limited role in endogenous processes and the fact that only a mild modification of the pharmacokinetics and pharmacodynamics can lead to significant results (3). Nevertheless, the development of an improved framework for the evaluation of pharmacogenomic variants, by combining different classifiers and appropriately adjusting their prediction thresholds has led to promising results(26).

Herein, we propose a comprehensive model for the assessment of pharmacogenomics variants by evaluating *in silico* protein prediction scores with the use of machine learning (ML), and thus highlighting the pharmacogenomic variants that are most likely to alter the protein function and consequently have a PGx impact.

## Results

### Performance metrics for the machine learning models toward the functional assessment of PGx variants

The performance of the classifiers, which were constructed with variables recommended by the RFE procedure, appeared to be advantageous, regardless of the limited sample size of the training set. More precisely, the metrics computed for the assessed machine learning models were as follows: Random Forest (RF) – Accuracy: 0.85 (95% CI : 0.79, 0.90), Area Under the Curve (AUC) = 0.92, Area Under the Precision-Recall Curve (prAUC) = 0.73; AdaBoost – Accuracy: 0.82 (95% CI : 0.76, 0.87), AUC: 0.91, prAUC: 0.72; XGBoost – Accuracy : 0.80 (95% CI : 0.73, 0.85), AUC: 0.91, prAUC: 0.73; Multinomial logistic regression – Accuracy: 0.78 (95% CI : 0.72, 0.84), AUC: 0.93, prAUC: 0.74. Interestingly, multinomial logistic regression led to higher AUC and prAUC values compared to the tree-based approaches, whilst the achieved accuracy was the lowest amongst the assessed models.

RFs were selected as the final approach to be used for the described classification task, since the respective model presented overall improved performance (i.e., accuracy, sensitivity, specificity, and precision) across all four functional classes. Regarding the 'Decreased function' variants, RFs were more sensitive and precise than the other assessed models, although AdaBoost achieved equal specificity values (Fig. 1). All models performed impressively well towards the 'Increased function' category and led to very similar outcomes, while RFs appeared superior for the detection of 'No

function' variants and AdaBoost and Multinomial Logistic Regression models were more sensitive for the 'Normal function' class.

The selected machine learning model proved to be highly specific ( $\geq 92\%$ ) for all 4 functional classes of variants, with lower, but still favorable values of sensitivity (80% – 98%), precision (80% – 98%) and balanced accuracy (86% – 99%). With regards to identifying variants that could lead to proteins with unchanged (normal), reduced or no function, we observed the lowest values of the metrics.

The model was characterized by a better performance for 'normal function' variants (Sensitivity = 0.8, Specificity = 0.92, Precision = 0.84, Balanced Accuracy = 0.86), followed closely by 'No' function variants (Sensitivity = 0.81, Specificity = 0.93, Precision = 0.81, Balanced Accuracy = 0.87) and finally 'Decreased function' variants (Sensitivity = 0.81, Specificity = 0.95, Precision = 0.80, Balanced Accuracy = 0.88). Interestingly, the classifier performs impressively well for the category of 'Increased Function' variants, in which case all computed metrics were above 98% (Fig. 1).

We also attempted to assess the variables that could significantly affect the presented machine learning model. More specifically, when it comes to the variable importance, the highest-ranking positions were occupied by these features that RFE suggested as the most informative ones for the classification task. In the present instance, LoFtool emerged as the prominent for the categorization of a variant according to its effect on protein function (Figure S1).

## **Application of the machine learning model in NGS data**

### **First case study (WGS data)**

As mentioned earlier and to further demonstrate the prediction performance of the final RF model, we tested its applicability in "unseen" NGS data, namely those data that have not been previously used to train the machine learning algorithm. We first tested its applicability in WGS data from a patient diagnosed with coeliac disease. From this process, 1,808 variants within the 10 pharmacogenes of interest were identified. Of these, 4 common (rs1801159, rs2306283, rs4149056, rs35364374) and 1 intermediate (rs3745274) PGx variants – with MAFs based on GnomAD genomes - had adequate information to be further processed by the RF model and these were characterized as missense variants. Of these 1,808 analyzed variants, we did not identify any variants categorized as loss-of-function variants (LoF).

Table 1 presents these variants, alongside with their predicted functional impact, as defined by the majority vote of the individual decision trees. For example, a random forest containing of 1000 distinct decision trees was built. If most of those votes recommend that the variant belongs to 'No function' variants, then this is the class that is attributed to the variant. In addition, the probability of being classified in each class, as based on the votes of all trees of the random forest built, is also provided (Table 1).

Table 1

Classification outcomes (prediction and probabilities) for WGS data using the final RF model. The predicted class is determined based on a majority vote from the individual decision trees of the random forest classifier, while the presented probabilities depict the corresponding percentage of decision trees voting towards a functional class.

rsID	Location (GRCh38)	Allele	SYMBOL	Ensembl Transcript ID	Predicted class	Probability of attributed class	GnomAD AF (%)
rs1801159	1:97515839–97515839	C	<i>DPYD</i>	ENST00000370192.8	Normal	0.955	18.49
rs2306283	12:21176804–21176804	G	<i>SLCO1B1</i>	ENST00000256958.3	Normal	0.664	55.33
rs4149056	12:21178615–21178615	C	<i>SLCO1B1</i>	ENST00000256958.3	Decreased	0.883	11.95
rs35364374	19:38492540–38492540	T	<i>RYR1</i>	ENST00000359596.8	Increased	0.377	4.95
rs3745274	19:41006936–41006936	T	<i>CYP2B6</i>	ENST00000324071.10	Decreased	0.724	28.44

This computational process led to the identification of 2 missense variants (located within *SLCO1B1* and *CYP2B6* respectively) that could potentially lead to proteins with decreased functionality and 1 missense variant classified as 'increased function' (located in *RYR1*). The remaining two variants were predicted to lead to no changes in the protein function (i.e., normal). The rest of the PGx variants had a high rate (over 85%) of missing values in the features of interest and were mostly ( $N = 1,765$  out of 1,803; 97.89%) located within intronic regions (Fig. 2). The latter were followed by variants in 3' prime UTRs ( $N = 20$ ; 1.11%), missense ( $N = 6$ ; 0.33%) and synonymous ( $N = 11$ ; 0.61%) variants. Interestingly, *DPYD* which encodes for a drug metabolizing enzyme, accumulated more than 1,000 intronic variants.

Regarding the potential clinical actionability of these 5 variants (rs1801159, rs2306283, rs4149056, rs35364374, rs3745274), we retrieved additional information from the PharmGKB database. rs1801159 and rs2306283 were not associated with any predicted changes in the protein function or changes in the dosing guidelines (i.e., normal, or low-level changes respectively). However, changes in treatment were recommended for individuals with rs4149056 genotype, whilst also stating that any additional risk factors should be considered for statin-induced myopathy. Moreover, rs3745274 carried multiple levels of CPIC evidence, for a variety of drugs such: as efavirenz, nevirapine, propofol, imatinib, cyclophosphamide, doxorubicin, mitotane, methadone and 3,4-methylenedioxymethamphetamine. No relevant information could be retrieved for rs35364374.

## Second case study (targeted PGx sequencing data)

The second case study consisted of targeted PGx sequencing data from 304 individuals of Greek origin and diagnosed with psychiatric disorders. Interestingly, 343 variants were identified, covering 10 pharmacogenes (*DPYD*, *CYP2C19*, *CYP2C9*, *CYP2C8*, *SLCO1B1*, *NUDT15*, *CYP2B6*, *UGT1A1*, *CYP2D6*, *TPMT*), 18 of which were attributed a SO consequence indicative of LoF variants. More specifically, we found 11 'frameshift', 6 'stop gained' and 1 'start lost' variants. None of these variants was assessed by the RF model owing to the high levels of missing values (mean: 77% missing values in the scores of interest). The remaining variants were mostly missense ( $N = 205$ ) or synonymous ( $N = 107$ ) (Fig. 3). According to GnomAD genome frequencies in the general population (AF), which were available for 88 of these variants, the dataset was enriched for 'ultra-rare' variants ( $N = 42$ ), followed by 'rare' ( $N = 18$ ), 'low frequency' ( $N = 14$ ), 'common' ( $N = 8$ ) and 'intermediate' ( $N = 6$ ) variants.

The dataset of 343 variants includes 195 known and 148 novel variants, of which 86 (44.1 %) and 70 (47.3 %) PGx variants (156 in total) were evaluated the final RF model, and the resulting predictions (class and respective probability) are presented in Table S1. The evaluated variants were mostly missense (i.e., 149 missense, 7 missense/splice region). Of

these, 71 variants led to “Decreased function” proteins, 41 variants to “No” function proteins, 1 variant in “Increased function” protein and 43 variants have no effect on protein functionality (i.e., “normal” function) (Fig. 4).

Regarding the potential clinical actionability of the 156 PGx variants as evaluated by the RF model, additional clinical and variant information was retrieved from PharmGKB. rs1801159, rs1801158, rs2297595 and rs1801160 were not associated with any predicted changes in the protein function, according to the variant annotation by PharmGKB, which constitutes an observation in concordance with the assigned prediction classes by the RF model (i.e., “normal” function class). Moreover, rs67376798 was associated with decreased catalytic activity based on evidence from PharmGKB, thus further confirming the prediction class of the RF model (i.e., “decreased” function class). Similar observations were applied for the variants, namely rs4149056, rs116855232 and rs3745274, for which the following prediction classes were assigned by the RF model: “decreased”, “no”, “decreased” respectively. PharmGKB provides multiple levels of clinical evidence for these variants, the majority of which are associated with decreased protein activity, therefore confirming the presented model results.

## Discussion

Conventional genetic testing and clinical guidelines focus solely on a small number of well-studied variants or star alleles in pharmacogenes, while the application of NGS techniques provides the possibility to detect a much wider range of (pharmacogenomics) variants. Recent studies have demonstrated that coding variants are rare, population-specific, and with a significant proportion of them potential affecting the protein product (based on *in silico* assays and metrics) (10–14). At the same time, the role of copy number variants (CNVs) within pharmacogenes (27), as well as variants in non-coding regions, are gaining more attention, with more than 90% of the polymorphisms detected in GWAS pharmacogenomic studies being non-coding (28). Owing to the limited number of thoroughly documented PGx variants and the incredibly large number of identified genetic mutations that should be experimentally validated, the initial evaluation of these variants found must be performed via the use of *in silico* tools.

The study’s main aim was the assessment of the utility of *in silico* derived scores, commonly used for variant annotation, toward the characterization of the potential protein function effects of SNVs identified within pharmacogenes. Amongst the assessed algorithms (AdaBoost, XGBoost, Random Forest, multinomial logistic regression), Random Forest presented superior performance and was selected as the final classifier. RFs have been also proven to be robust in the presence of outliers or noise, effective, even without configuration, and useful in cases where the number of available data are limited, compared to the number of available variables, as is often the case with ‘-omics’ data (29, 30).

The final classifier required minimum hyperparameter tuning and integrated 7 scores, stand-alone or ensemble ones, and 2 custom created variables. The overall accuracy was found to be equal to 0.85 (95% CI: 0.79, 0.90), with an Area Under the Curve of 0.92 and an Area Under the Precision-Recall Curve (PR AUC) of 0.73. The by-class performance for variants of Normal, Decreased and No Function classes is efficient enough, although there is still room from improvement, especially in terms of sensitivity (0.80, 0.81, and 0.81 respectively). Interestingly, the model appears to be efficient, given the fact that most of the incorporated features are used to distinguish between damaging and benign variants, specifically when it comes to identifying Increased function SNVs. Furthermore, LoFtool, an approach that evaluates the tolerance of a gene to loss-of-function mutations emerged as the most significant determinant of the classification task. Considering the superior performance of the model in identifying “Increased Function” PGx variants, and the observation that this specific class in the training dataset represents only two pharmacogenes, might partially justify the importance of the variable.

Although there is limited published work in this specific area, the possibility of using PGx variants so as to develop classification tools has been previously explored, without however progressing any further due to the limitations and difficulties that accompany this field (31). Firstly, the most frequently examined properties in such classifications tools are the degree of evolutionary conservation, which is observed in lower levels in pharmacogenes (3) and hence its usefulness

is debated by a series of studies (26, 31, 32), as well as parameters regarding the structure of the respective proteins, which have been observed to lead to small increases in the efficiency of the classifiers produced (31). Overall, such factors could influence the quality of the output results in classification models, as the one presented herein.

In addition, the training sets used to train computational models are usually comprised of common polymorphisms against variants (mostly SNVs) related to disease-causality, while in terms of drug response, the modifying effect of common polymorphisms cannot be rule out. Moreover, the resulting scores evaluate the pathogenic potential of the examined variants and classify them into two usually categories according to certain applied thresholds. In contrast, PGx researchers usually focus on the induced change in protein function, which can be distinguished at several levels (e.g., increase, decrease, no change, complete loss of activity), while the differential drug response is not a disease, but a phenotype that occurs under specific conditions (i.e., administration of a specific drug).

For example, in a recent study, the adaptation of the proposed classification thresholds and the subsequent integration of selected algorithms, which could provide optimal results for the creation of a comprehensive score, led to a tool with exceptional sensitivity and specificity (26). However, this work focused exclusively on the distinction between loss-of-function and neutral variants, hence ignoring PGx variants that would result in a protein product of increased activity, and which are of interest in PGx field.

The novelty of the present approach lies in the integration of PGx variants characterized as 'Increased function'; however, the limitations probably affect the performance of the presented classification model. More specifically, the collection and curation of larger training datasets with appropriately defined classes for PGx purposes was one of the most significant challenges in this work. Furthermore, the computed metrics demonstrate a difficulty in distinguishing between normal and decreased/no function variants, thus making debatable the suitability of the used features for the characterization of these PGx variants.

Given the challenges and implications for the prediction of functional impact of PGx variants, as well as the complexity of the involved biological processes (33), the findings of this study should be interpreted with caution. For example, discrepancies have been observed not only amongst different algorithms, or between *in silico* predictions and *in vitro* activity (34), but also when comparing *in vitro* and *in vivo* observations, with the typical example being *CYP2D6* \*35, which despite experimental evidence of reduced hydroxylation capacity of tamoxifen(35) has not been associated with reduced activity (36). Moreover, researchers should bear in mind that the same variant may affect the response to different drugs in different ways. For example, although the \*10 and \*13 *CYP2C8* alleles have been found to affect the N-deethylation of amodiacin, the hydroxylation of paclitaxel – which is also metabolized by CYP2C8 - remains unaffected (37).

As mentioned earlier, the presented model has demonstrated promising results, despite the limitations of this computational research field. However, there is still space for further improvements towards a more efficient and robust version of the presented model. More specifically, it would be useful to examine and compare the performance of other machine learning (ML) approaches, supervised or not. Furthermore, significant advantages are expected to emerge from the incorporation of wider training sets, consisting of larger numbers of variants and covering an additional number of pharmacogenes. Emphasis should be also placed on the creation of well-characterized sets of PGx variants at the level of protein effects, both laboratory and clinical, as well as on the improvement of the existing databases to facilitate the export of the requested information. Finally, since the contribution of various factors to the response to a given drug is non-debatable, a more comprehensive approach through systemic genomics would be particularly useful, thus incorporating a variety of different -omics data (38).

## Conclusions

The novelty of the present computational model lies in the fact that a ML approach was used to classify PGx variants by consequently assigning a protein activity prediction. Overall, the presented model prioritizes annotated PGx variants in different variant effect classes and then assigns a protein function classification after stringent computational assessment and ML processes. Its utility was further showcased by using a real-life case to further support the applicability of this model as a clinical support decision tool. Indeed, a validated, methodical prioritization of the multitude of genomic variants stemming from NGS analyses, as the one presented herein, has the potential to positively contribute towards the large-scale clinical application of pharmacogenomics and facilitate the translation of a patient's genomic profile into actionable information.

## Methods

### Collecting the training data

An appropriate training set of variants was manually curated using the PGx Gene-specific information tables, created under the collaboration between PharmGKB and CPIC and was subsequently supplemented by additional variants from PharmVar (39). This training set consists of 262 variants located across 12 pharmacogenes, with well-defined protein-level functional consequences, as based on the integration and assessment of *in vitro* biochemical assays, *in vivo* evidence, and clinical observations. The observed functionality is classified into 5 levels (excluding Unknown/Uncertain function): Increased, Normal, Possibly Decreased, Decreased and No function. However, owing to the limited number of observations harboring the levels of 'Possibly Decreased' and 'Decreased' functions and after careful examination of the available information for those categories, these two levels were combined in one class (Decreased function) (Table 2).

Table 2

Description of the protein function effect classes of PGx variants, which are used as the training data for the final RF model. The functionality class is split in the following classes ('decreased', 'increased', 'no', 'normal'), the number of the respective PGx variants per class is also provided, as well as which pharmacogenes are incorporated per each class.

Functionality Class	Number of variants	Representation of genes
Decreased	36	<i>DPYD, CYP2C19, CYP2C9, SLCO1B1, RYR1, CYP2B6, UGT1A1, CYP2D6</i>
Increased	46	<i>RYR1, CYP2B6</i>
No	48	<i>CYP2C19, CYP2C8, DPYD, CYP2C9, NUDT15, CYP2B6, CYP2D6, TPMT</i>
Normal	60	<i>DPYD, CYP2C9, SLCO1B1, CYP2B6, CYP2D6</i>

### Variant annotation

The curated set of pharmacogenomic variants was annotated using the web interface of Ensembl's Variant Effect Predictor (VEP) tool, for the GRCh38 human assembly, as well as the 4.1.a version of the dbSNFP database (40), which is also provided through VEP. The majority of the retrieved information is available for variants located within protein coding regions and includes: a detailed characterization at a protein level (i.e., database identifiers, codons, amino acids, coordinates, protein domains, computational scores, etc.), overlapping known variants, observed frequencies in different populations (i.e., via the 1000 Genomes Project, the genome Aggregation Database, the Exome Aggregation Consortium data and the Exome Sequencing Project), any related phenotypes (e.g., OMIM, Orphanet, GWAS catalog) or clinical significance (ClinVar), as well as literature references (41). Furthermore, the attributed consequence, described by using terms as developed in collaboration with Sequence Ontology (SO),(42) and the corresponding impact of a variation are also provided.

Regarding the retrieved frequency data, variants were classified as 'common' if the minor allele frequency (MAF) was equal or above 10% ( $MAF \geq 10\%$ ) and as 'intermediate' if the MAF ranged between 5% and 10% ( $5\% \leq MAF < 10\%$ ). Variants were classified as 'low frequency' if the MAF ranged between 1% and 5% ( $1\% \leq MAF < 5\%$ ), whilst 'rare' variants included these

variants of which the MAF was between 0.1% and 1% ( $0.1\% \leq \text{MAF} < 1\%$ ). Finally, variants were classified as ultra-rare if the MAF was equal or below 0.1% ( $\text{MAF} \leq 0.1\%$ ).

Features and variants with a high percentage of missing values ( $\geq 40\%$ ) were excluded, while the remaining values were imputed by using k-Nearest Neighbors algorithm (kNN) (43) with default values for k-neighbors (equal to 5) and inverse weighted mean Gower distances (44). In addition, a step of backwards variable selection through Recursive Feature Elimination (RFE) using Bagged Trees was performed, which recommended the use of 7 out of the 45 variables (LoFtool (21), DEOGEN2\_score (45), MPC\_score (46), BayesDel\_addAF\_score (47), integrated\_fitCons\_score (48), FATHMM\_score (49), LIST.S2\_score (50)). Furthermore, two binary variables were constructed and included in the analysis: one indicating whether the variant was located within a protein functional domain (according to InterPro (51) annotation) and one representing high impact SO consequences (splice acceptor or donor variants, stop gained, frameshift variants, stop or start lost), enriched for Loss-of-Function (LoF) changes, as defined by MacArthur et al. (2012)(52). The distribution of the final set of 190 variants (in 11 pharmacogenes; Table S2) is presented in Table 1.

## Training of the machine learning model

All preprocessing and ML-related analyses described in this work were performed using the R language for statistical programming (version 4.0.2) (53). To exploit the abilities of the abovementioned features toward explaining potential protein function effects of variants derived from NGS analyses, a variety of tree-based methodologies was assessed, alongside with a special case of a neural network acting in a multinomial logistic regression manner. More specifically, Random Forests (54, 55), Multi-class AdaBoost (56, 57), XGBoost (58), and a neural network striped from its hidden layers and activation functions (multinomial logistic regression) (59, 60) were used via the caret package (61). For the selected tree-based models, hyperparameters were tuned based on the optimization of the accuracy metric, while in multinomial logistic regression, the default parameters were used (Table S3).

## Evaluation of the machine learning models

The predictive performance of the created models was assessed via the 5-fold cross validation (CV) method. During n-fold CV, the data are divided to create n equal-sized subsets; n-1 of these are used to train a model and the remaining 1 is used to test its' performance. This process is repeated n times, until all subsets have been used to test the model, while the computed metrics in each iteration are averaged. More specifically, the metrics of interest include the Accuracy, Precision, Sensitivity (True Positive rate), Specificity (True Negative rate) and the F-measure. Since this was a multi-class task, all metrics were computed for each class separately (according to the one-vs-all method), and the performance of the model was calculated using the corresponding weighted average values for each metric. Furthermore, a random forest classifier was trained with the total of 47 features and used to evaluate their predictive importance.

## Testing the applicability of the final machine learning model

To further demonstrate the applicability of the machine learning model, we applied the classifier in data derived from NGS analyses. To this end, variant call format (.vcf) files comprised of the results from: (i) a WGS analysis of a single individual of Greek origin diagnosed with coeliac disease, and, (ii) a targeted pharmacogene sequencing analysis of 304 individuals of Greek origin diagnosed with psychiatric diseases (62). Firstly, the provided variants were annotated, using the web interface of ensemble VEP tool, while the resulting data were preprocessed to select only these identified in the transcripts of interest. Then, these annotation data were used as an input to our final RF model and the corresponding prediction functionality classes and prediction probabilities were provided.

## Abbreviations

SNVs: single nucleotide Variations

PGx: pharmacogenomics

NGS: Next-generation sequencing

WES: Whole Exome sequencing

WGS: Whole Genome sequencing

EMA: European Medicines Agency

FDA: Food and Drug Administration

CNVs: Copy Number Variants

MAF: minor allele frequency

RF: Random Forest

AUC: Area Under the Curve

prAUC: Area Under the Precision-Recall Curve

## **Declarations**

### **Ethics approval and consent to participate**

Not applicable.

### **Consent for publication**

Not applicable.

### **Availability of data and materials**

Training, and testing materials are available upon request.

### **Competing interests**

IT and EC are shareholders of Katharsis Technologies Inc. GPP is Full Member and National Representative of the European Medicines Agency, Committee for Human Medicinal Products – Pharmacogenomics Working Party, Amsterdam, the Netherlands.

### **Funding**

No funding was received

### **Authors' contributions**

GPP has conceived the project. M-TP and MK implemented the software, wrote the code and performed the data analysis. IT, SP, EC and PJS validated the software. GPP supervised the project. M-TP, MK and GPP wrote the manuscript. All authors have read and approved the manuscript.

### **Acknowledgements**

We wish to thank Wasun Chantratita, Taisei Mushiroda and Koya Fukunaga for kindly performing NGS analysis on the 304 individuals comprising the cohort .vcf of our second case study.

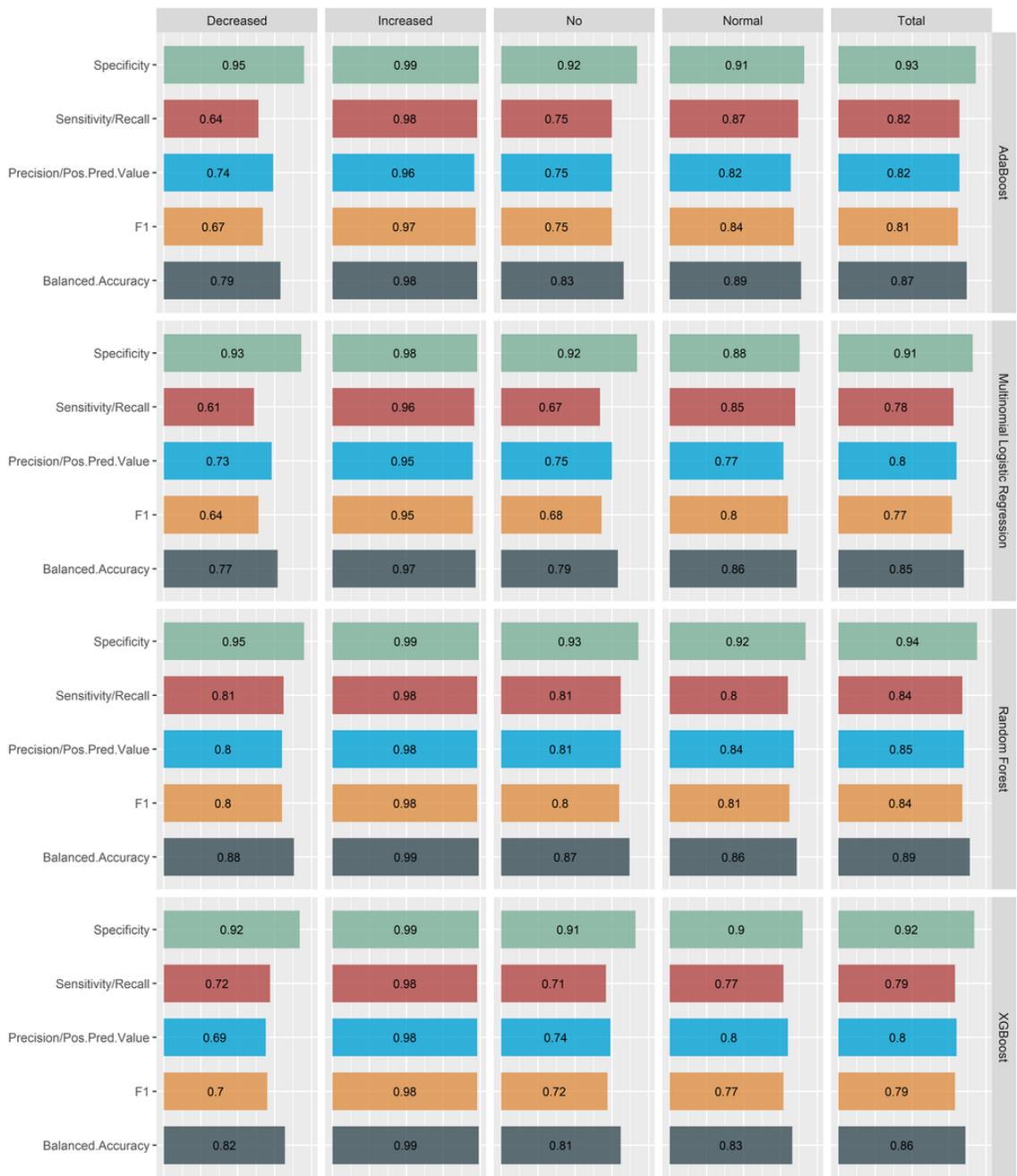
## References

1. International Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human Use (ICH). ICH E15 Definitions for genomic biomarkers, pharmacogenomics, pharmacogenetics, data and sample coding categories. 2007.
2. Zhou ZW, Chen XW, Sneed KB, Yang YX, Zhang X, He ZX, et al. Clinical association between pharmacogenomics and adverse drug reactions. *Drugs*. 2015;75(6):589–631.
3. Lauschke VM, Milani L, Ingelman-Sundberg M. Pharmacogenomic Biomarkers for Improved Drug Therapy-Recent Progress and Future Developments. *AAPS J*. 2017;20(1):4.
4. Moyer AM, Caraballo PJ. The challenges of implementing pharmacogenomic testing in the clinic. *Expert Rev Pharmacoecon Outcomes Res*. 2017;17(6):567–77.
5. Mizzi C, Peters B, Mitropoulou C, Mitropoulos K, Katsila T, Agarwal MR, et al. Personalized pharmacogenomics profiling using whole-genome sequencing. *Pharmacogenomics*. 2014;15(9):1223–34.
6. Katsila T, Patrinos GP. Whole genome sequencing in pharmacogenomics. *Front Pharmacol*. 2015;6:61.
7. Giannopoulou E, Katsila T, Mitropoulou C, Tsermpini EE, Patrinos GP. Integrating Next-Generation Sequencing in the Clinical Pharmacogenomics Workflow. *Front Pharmacol*. 2019;10:384.
8. Ehmann F, Caneva L, Prasad K, Paulmichl M, Maliepaard M, Llerena A, et al. Pharmacogenomic information in drug labels: European Medicines Agency perspective. *Pharmacogenomics J*. 2015;15(3):201–10.
9. Relling MV, Evans WE. Pharmacogenomics in the clinic. *Nature*. 2015;526(7573):343–50.
10. Tennessen JA, Biggam AW, O'Connor TD, Fu W, Kenny EE, Gravel S, et al. Evolution and Functional Impact of Rare Coding Variation from Deep Sequencing of Human Exomes. *Science*. 2012;337(6090):64.
11. Nelson MR, Wegmann D, Ehm MG, Kessner D, St Jean P, Verzilli C, et al. An abundance of rare functional variants in 202 drug target genes sequenced in 14,002 people. *Science*. 2012;337(6090):100–4.
12. Fujikura K, Ingelman-Sundberg M, Lauschke VM. Genetic variation in the human cytochrome P450 supergene family. *Pharmacogenet Genomics*. 2015;25(12):584–94.
13. Kozyra M, Ingelman-Sundberg M, Lauschke VM. Rare genetic variants in cellular transporters, metabolic enzymes, and nuclear receptors can be important determinants of interindividual differences in drug response. *Genet Med*. 2017;19(1):20–9.
14. Zhang B, Lauschke VM. Genetic variability and population diversity of the human SLCO (OATP) transporter family. *Pharmacol Res*. 2019;139:550–9.
15. Ng PC, Henikoff S. Accounting for human polymorphisms predicted to affect protein function. *Genome Res*. 2002;12(3):436–46.
16. Choi Y, Sims G, Murphy S, Miller J, Chan A. Predicting the functional effect of amino acid substitutions and indels. *PLoS One*. 2012;7(10):e46688.
17. Adzhubei I, Jordan DM, Sunyaev SR. Predicting functional effect of human missense mutations using PolyPhen-2. *Curr Protoc Hum Genet*. 2013;Chap. 7:Unit7 20.
18. Li B, Krishnan VG, Mort ME, Xin F, Kamati KK, Cooper DN, et al. Automated inference of molecular mechanisms of disease from amino acid substitutions. *Bioinformatics*. 2009;25(21):2744–50.
19. Davydov EV, Goode DL, Sirota M, Cooper GM, Sidow A, Batzoglou S. Identifying a high fraction of the human genome to be under selective constraint using GERP++. *PLoS Comput Biol*. 2010;6(12):e1001025.
20. Ionita-Laza I, McCallum K, Xu B, Buxbaum JD. A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nat Genet*. 2016;48(2):214–20.

21. Fadista J, Oskolkov N, Hansson O, Groop L. LoFtool: a gene intolerance score based on loss-of-function variants in 60 706 individuals. *Bioinformatics*. 2017;33(4):471–4.
22. Quang D, Chen Y, Xie X. DANN: a deep learning approach for annotating the pathogenicity of genetic variants. *Bioinformatics*. 2015;31(5):761–3.
23. Ioannidis NM, Rothstein JH, Pejaver V, Middha S, McDonnell SK, Baheti S, et al. REVEL: An Ensemble Method for Predicting the Pathogenicity of Rare Missense Variants. *Am J Hum Genet*. 2016;99(4):877–85.
24. Dong C, Wei P, Jian X, Gibbs R, Boerwinkle E, Wang K, et al. Comparison and integration of deleteriousness prediction methods for nonsynonymous SNVs in whole exome sequencing studies. *Hum Mol Genet*. 2015;24(8):2125–37.
25. Ingelman-Sundberg M, Mkrтчian S, Zhou Y, Lauschke VM. Integrating rare genetic variants into pharmacogenetic drug response predictions. *Hum Genomics*. 2018;12(1):26.
26. Zhou Y, Mkrтчian S, Kumondai M, Hiratsuka M, Lauschke VM. An optimized prediction framework to assess the functional impact of pharmacogenetic variants. *Pharmacogenomics J*. 2019;19(2):115–26.
27. Santos M, Niemi M, Hiratsuka M, Kumondai M, Ingelman-Sundberg M, Lauschke VM, et al. Novel copy-number variations in pharmacogenes contribute to interindividual differences in drug pharmacokinetics. *Genet Med*. 2018;20(6):622–9.
28. Luizon MR, Ahituv N. Uncovering drug-responsive regulatory elements. *Pharmacogenomics*. 2015;16(16):1829–41.
29. Probst P, Wright MN, Boulesteix AL. Hyperparameters and tuning strategies for random forest. *Wiley Interdisciplinary Reviews: Data Mining Knowledge Discovery*. 2019;9(3):e1301.
30. Chen X, Ishwaran H. Random forests for genomic data analysis. *Genomics*. 2012;99(6):323–9.
31. Li B, Seligman C, Thusberg J, Miller JL, Auer J, Whirl-Carrillo M, et al. In silico comparative characterization of pharmacogenomic missense variants. *BMC Genomics*. 2014;15(Suppl 4(Suppl 4):4-S.
32. Lauschke VM, Ingelman-Sundberg M. How to Consider Rare Genetic Variants in Personalized Drug Therapy. *Clin Pharmacol Ther*. 2018;103(5):745–8.
33. Evans WE, Relling MV. Pharmacogenomics: Translating Functional Genomics into Rational Therapeutics. *Science*. 1999;286(5439):487.
34. Devarajan S, Moon I, Ho MF, Larson NB, Neavin DR, Moyer AM, et al. Pharmacogenomic Next-Generation DNA Sequencing: Lessons from the Identification and Functional Characterization of Variants of Unknown Significance in CYP2C9 and CYP2C19. *Drug Metab Dispos*. 2019;47(4):425–35.
35. Muroi Y, Saito T, Takahashi M, Sakuyama K, Niinuma Y, Ito M, et al. Functional Characterization of Wild-type and 49 CYP2D6 Allelic Variants for *N*-Desmethyltamoxifen 4-Hydroxylation Activity. *Drug Metabolism and Pharmacokinetics*. 2014;29(5):360–6.
36. Gaedigk A, Ryder DL, Bradford LD, Leeder JS. CYP2D6 Poor Metabolizer Status Can Be Ruled Out by a Single Genotyping Assay for the – 1584G Promoter Polymorphism. *Clin Chem*. 2003;49(6):1008–11.
37. Tsukada C, Saito T, Maekawa M, Mano N, Oda A, Hirasawa N, et al. Functional characterization of 12 allelic variants of CYP2C8 by assessment of paclitaxel 6 $\alpha$ -hydroxylation and amodiaquine *N*-deethylation. *Drug Metab Pharmacokinet*. 2015;30(5):366–73.
38. Li R, Kim D, Ritchie MD. Methods to analyze big data in pharmacogenomics research. *Pharmacogenomics*. 2017;18(8):807–20.
39. Gaedigk A, Whirl-Carrillo M, Pratt VM, Miller NA, Klein TE. PharmVar and the Landscape of Pharmacogenetic Resources. *Clin Pharmacol Ther*. 2020;107(1):43–6.
40. Liu X, Li C, Mou C, Dong Y, Tu Y. dbNSFP v4: a comprehensive database of transcript-specific functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Genome Med*. 2020;12(1):103.

41. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The Ensembl Variant Effect Predictor. *Genome Biol.* 2016;17(1):122.
42. Cunningham F, Moore B, Ruiz-Schultz N, Ritchie GR, Eilbeck K. Improving the Sequence Ontology terminology for genomic variant annotation. *J Biomed Semantics.* 2015;6:32.
43. Kowarik A, Templ M. Imputation with the R Package VIM. *Journal of Statistical Software; Vol 1, Issue 7 (2016).* 2016.
44. Alfons A, Templ M. Estimation of Social Exclusion Indicators from Complex Surveys: The R Package laeken. *Journal of Statistical Software; Vol 1, Issue 15 (2013).* 2013.
45. Raimondi D, Tanyalcin I, Féré J, Gazzo A, Orlando G, Lenaerts T, et al. DEOGEN2: prediction and interactive visualization of single amino acid variant deleteriousness in human proteins. *Nucleic Acids Res.* 2017;45(W1):W201-W6.
46. Samocha KE, Kosmicki JA, Karczewski KJ, O'Donnell-Luria AH, Pierce-Hoffman E, MacArthur DG, et al. Regional missense constraint improves variant deleteriousness prediction. 2017:148353.
47. Feng B-JPERCH. A Unified Framework for Disease Gene Prioritization. 2017;38(3):243–51.
48. Gulko B, Hubisz MJ, Gronau I, Siepel A. A method for calculating probabilities of fitness consequences for point mutations across the human genome. *Nat Genet.* 2015;47(3):276–83.
49. Shihab HA, Gough J, Cooper DN, Stenson PD, Barker GLA, Edwards KJ, et al. Predicting the Functional, Molecular, and Phenotypic Consequences of Amino Acid Substitutions using Hidden. *Markov Models.* 2013;34(1):57–65.
50. Malhis N, Jacobson M, Jones SJM, Gsponer J. LIST-S2: taxonomy based sorting of deleterious missense mutations across species. *Nucleic Acids Res.* 2020;48(W1):W154-W61.
51. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, et al. InterPro: the integrative protein signature database. *Nucleic Acids Res.* 2009;37(Database issue):D211-5.
52. MacArthur DG, Balasubramanian S, Frankish A, Huang N, Morris J, Walter K, et al. A systematic survey of loss-of-function variants in human protein-coding genes. *Science.* 2012;335(6070):823–8.
53. R Development Core Team. R: A language and environment for statistical computing and graphics. 4.0.2 ed. Vienna: R Foundation for Statistical Computing; 2020.
54. Breiman L. Random Forests. *Mach Learn.* 2001;45(1):5–32.
55. Andy Liaw MW. Classification and Regression by randomForest. *R News.* 2002;2:18–22.
56. Alfaro E, Gamez M, García N. adabag: An R Package for Classification with Boosting and Bagging. *Journal of Statistical Software; Vol 1, Issue 2 (2013).* 2013.
57. Hastie T, Rosset S, Zhu J, Zou H. Multi-class AdaBoost. *Statistics Its Interface.* 2009;2(3):349–60.
58. Chen T, Guestrin C, editors. XGBoost: A Scalable Tree Boosting System. *KDD '16: The 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining; 2016 2016/08/13/.* San Francisco California USA: ACM.
59. Venables WN, Ripley BD, Venables WN. *Modern applied statistics with S.* 4th ed ed. New York: Springer; 2002 2002. 495 p.
60. Ripley BD. *Pattern Recognition and Neural Networks.* Cambridge: Cambridge University Press; 1996.
61. Kuhn M, Contributions from Jed Wing SW, Andre Williams C, Keefer A, Engelhardt T, Cooper Z, Mayer. Brenton Kenkel, the R Core Team, Michael Benesty, Reynald Lescarbeau, Andrew Ziem, Luca Scrucca, Yuan Tang, Can Candan, Tyler Hunt. *caret: Classification and Regression Training; 2019.*
62. Runcharoen C, Fukunaga K, Sensorn I, Iemwimangsa N, Klumsathian S, Tong H, et al. Prevalence of pharmacogenomic variants in 100 pharmacogenes among Southeast Asian populations under the collaboration of the Southeast Asian Pharmacogenomics Research Network (SEAPharm). *Human genome variation.* 2021;8(1):7.

## Figures



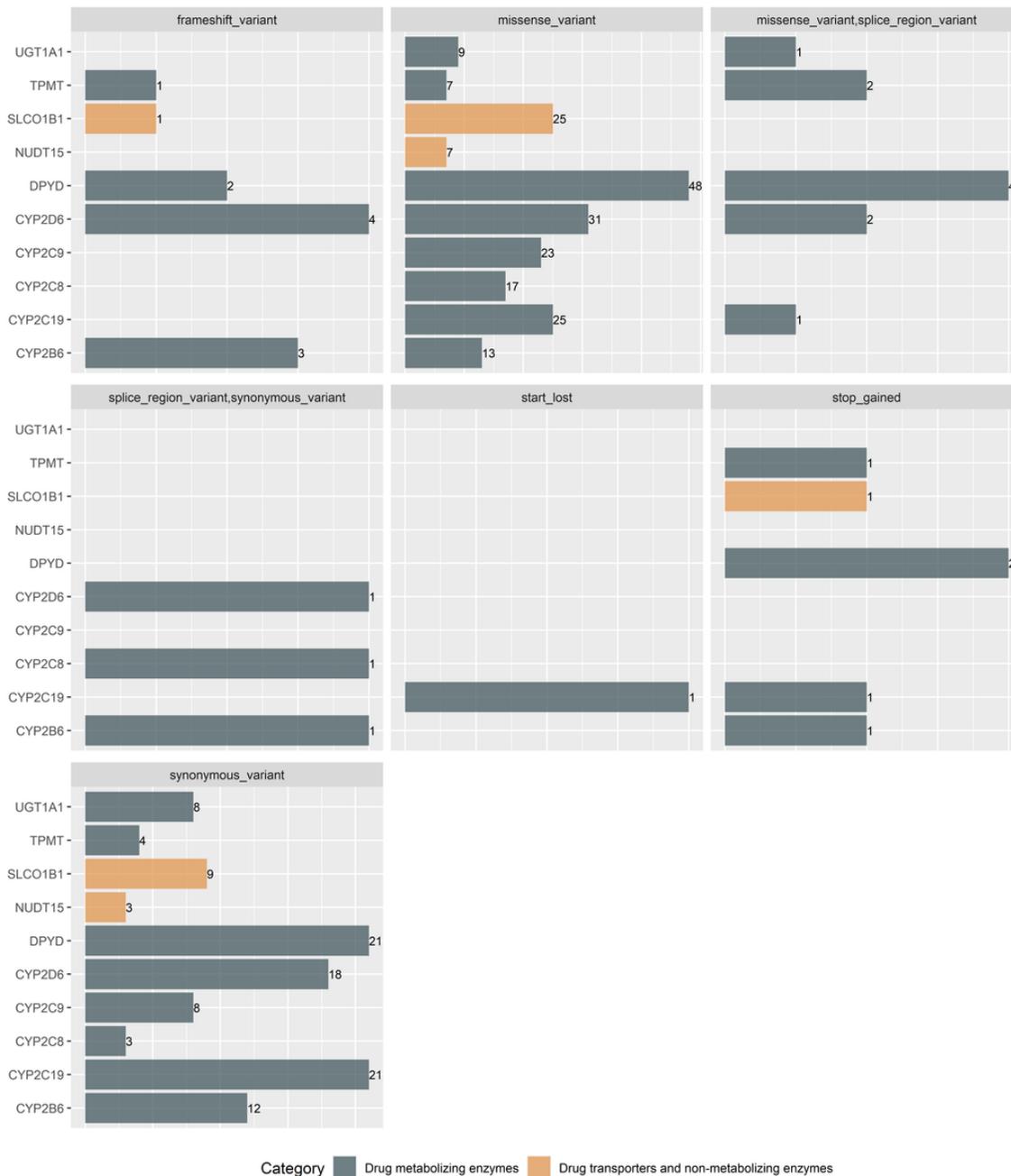
**Figure 1**

Metrics showing the performance of the different classifiers, namely AdaBoost, Multinomial Logistic Regression, Random Forest, XGBoost. More specifically, the sensitivity, specificity, positive predictive value (pos.pred.value), precision, F1 metric (harmonic mean of precision and recall) and balanced accuracy of the classifiers are provided for each protein function effect class.



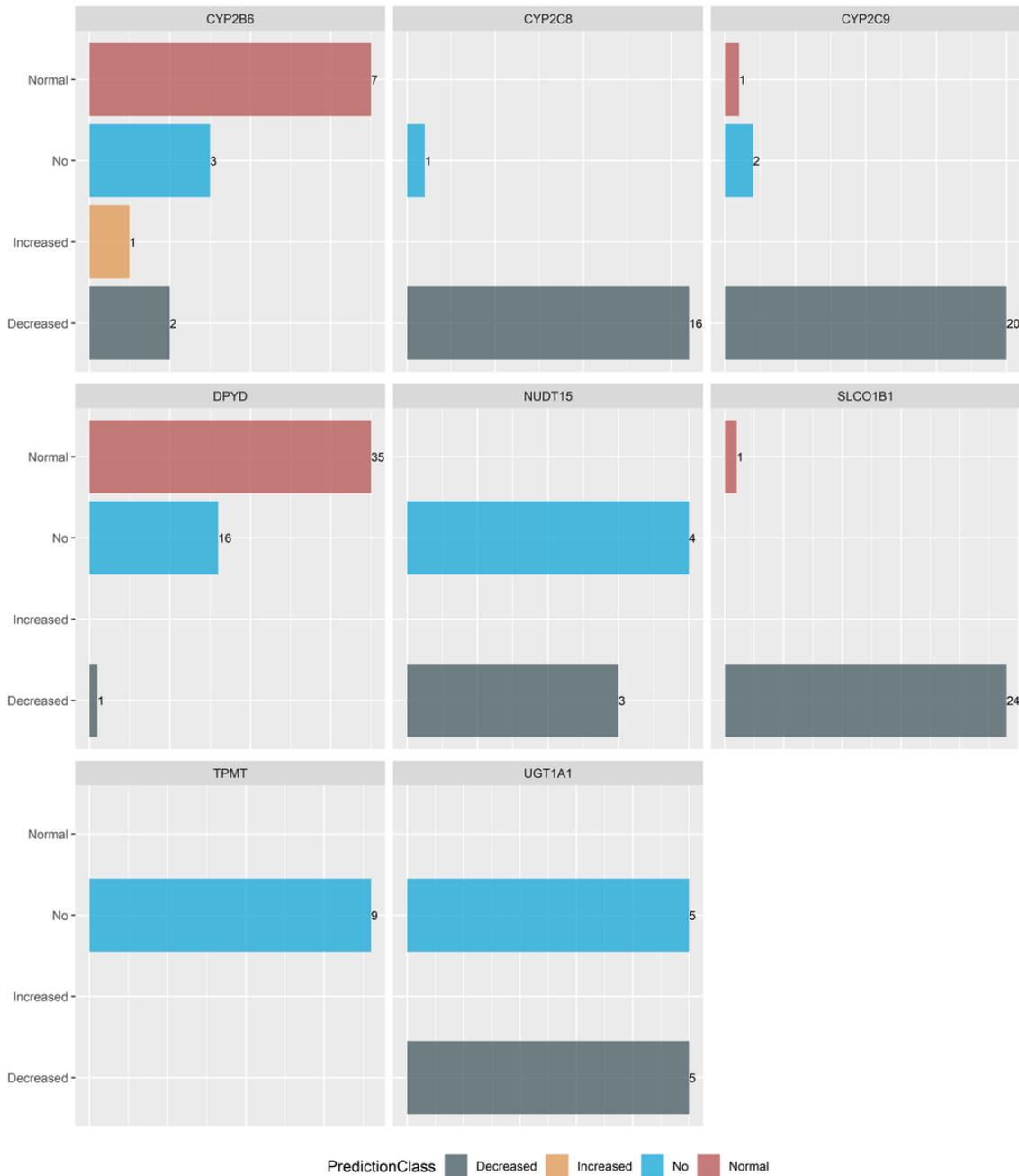
**Figure 2**

Distribution of PGx variants identified in the WGS data (first case study) that were not processed owing to many missing values. The graph presents the number of PGx variants, by gene, that were not processed any further by the machine learning model, according to the VEP consequence (i.e., 3' UTR variant, intronic variant, missense variant, splice region variant and synonymous variant). The pharmacogenes are color-coded according to the corresponding PGx group: genes encoding drug metabolizing enzymes or genes encoding drug transporters or other non-metabolizing enzymes.



**Figure 3**

Sequence ontology consequences for the identified PGx variants, as derived from a Greek cohort of 304 individuals with psychiatric disorders (second case study). 343 PGx variants within the pharmacogenes of interest were identified in this cohort. Amongst the consequences are ‘frameshift’, ‘missense’, ‘missense or splice region’, ‘splice region’, ‘start lost’, ‘stop gained’ and ‘synonymous’ variants.



**Figure 4**

Protein function predictions based on the final RF model after assessment of the targeted PGx sequencing data from a Greek cohort of 304 individuals. The functionality class for the PGx variants, as processed by the RF model, is depicted per each pharmacogene of interest. The numbers denote the totaling number of variants within each pharmacogene per each function class.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [FigureS1.png](#)
- [Pandi19042021suppl.docx](#)

- [TableS1.xlsx](#)