

Interpretable Machine Learning in Kidney Offering: Multiple Outcome Prediction for Accepted Offers

Achille Salaün (achille.salaun@eng.ox.ac.uk)^{1,*}, **Simon Knight** (simon.knight@nds.ox.ac.uk)², **Laura Wingfield** (lrwingfield@gmail.com)², and **Tingting Zhu** (tingting.zhu@eng.ox.ac.uk)¹

¹Institute of Biomedical Engineering, Department of Engineering, University of Oxford, Oxford, OX3 7DQ, United-Kingdom

²Nuffield Department of Surgical Sciences, University of Oxford, Oxford, OX3 9DU, United-Kingdom

*Corresponding author: achille.salaun@eng.ox.ac.uk

ABSTRACT

The decision to accept an organ offer for transplant, or wait for something potentially better in the future, can be challenging. Especially, clinical decision support tools predicting transplant outcomes are lacking. This project uses interpretable methods to predict both graft failure and patient death using data from previously accepted kidney transplant offers. Precisely, using more than twenty years of transplant outcome data, we train and compare several survival analysis and classification models in both single and multiple risk settings. In addition, we use *post hoc* interpretability techniques to clinically validate these models. In a single risk setting, neural networks provide comparable results to the Cox proportional hazard model, with 0.71 and 0.81 AUROC for predicting graft failure and patient death at year 10, respectively. Recipient and donor ages, primary renal disease, donor eGFR, donor type, and the number of mismatches at DR locus appear to be important features for transplant outcome prediction. We also extended the neural network approach to multiple outcome prediction, maintaining consistent performances and clinical interpretation. Thus, owing to their good predictive performance and the clinical relevance of their *post hoc* interpretation, neural networks represent a promising core component in the construction of future decision support systems for transplant offering.

1 Introduction

2 Around 2,500 deceased donor kidney transplants are performed in the UK each year. At any time, there are around 5,000
3 patients on the kidney transplant waiting list with an average wait of 2-3 years. The shortage of organs available for transplant
4 means that some patients become unfit for surgery or die whilst waiting. Because of this, clinicians often consider organ offers
5 from less-than-optimal donors with existing comorbidities or older age. Decisions around organ offers are made by clinicians
6 based upon the information available at the time of offer, including donor and recipient demographic and medical details.
7 Clinicians use their clinical experience, but do not have reliable tools available to help them predict what would happen if they
8 choose to accept or decline an offer and wait for the next available one. This uncertainty leads to considerable variability in
9 organ decline rates and waiting times between clinicians and centres. A computerised decision support (CDS) system that
10 accurately predicts transplant outcomes, both in terms of graft failure and patient death, as well as indicating what would
11 happen if the organ offer was declined (in terms of future offers and likely waiting time), may help to support clinicians in
12 making these difficult decisions. As decisions must remain under the responsibility and control of the clinician, any CDS tool
13 must be easy to use, and predictions must be interpretable from a clinician's perspective. Interpretability and usability are also
14 important to patients, allowing better explanations of likely outcomes during the informed consent process.

15 The aim of this study is to predict transplant outcomes in the scenario of an accepted kidney offer. We rely on more
16 than twenty years of registry data, containing over 36,000 accepted kidney transplant offers, with graft and patient survival
17 information. These data have been provided by National Health Service Blood and Transplant (NHSBT) with ethical approval.
18 Using these data, we have trained and compared several survival analysis and machine learning classification models, in both
19 single and multiple-risk settings. In addition, we use *post hoc* interpretability techniques to clinically validate these models.

20 Predicting the time of occurrence of an event (such as patient death or graft failure) from censored data has been extensively

21 studied under the name of survival analysis. This has many applications in health informatics such as predicting strokes [1],
22 oral cancer [2], or graft outcome prediction. Censored data are common in such contexts, resulting from loss of follow-up,
23 competing events, or the end of the study. In the context of graft outcome prediction, [3–5] use the Cox proportional hazard
24 (PH) model to predict kidney graft or recipient survival. The Cox PH model is a classic time-to-event approach that models
25 the hazard function, as in the failure rate of a system according to time [6]. This approach is not only robust and reliable, but
26 also simple to use and well understood by clinicians. Several generalisations of this model have been proposed. For instance,
27 DeepSurv [7] aims at increasing the modelling power of the Cox model by replacing the linear contribution of the covariates
28 with a neural network. Since the Cox model was originally designed to handle a single type of event, generalisations to multiple
29 risks (e.g. predicting both graft and patient survival) have also been proposed. However, the effects of the regression coefficients
30 on cause-specific survivability are not interpretable [8].

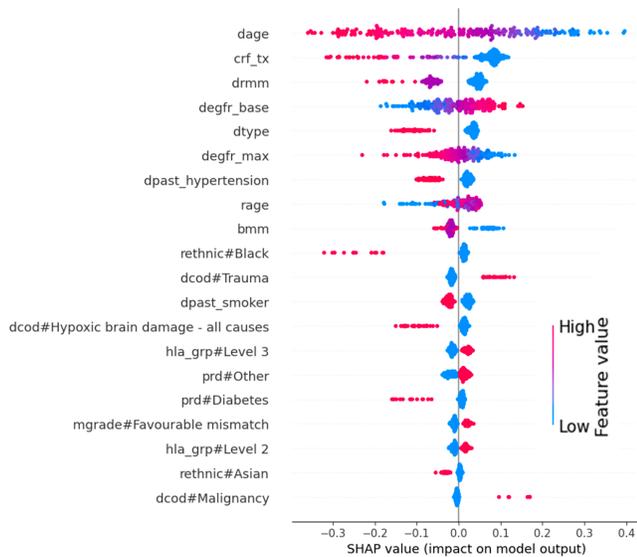
31 In general, one can distinguish two approaches for survival analysis: either providing a description of survivability over
32 time, or predicting the state of the subject (e.g., graft, recipient) at arbitrary time points. While the Cox PH model follows the
33 first approach, machine learning models can be used after converting the survival analysis problem into a classification problem
34 at a given point in time. In the case of predicting transplant outcomes, predictions at specific milestones (e.g., graft and patient
35 survival at years 1, 5, or 10) are generally sufficient: for example, existing risk communication tools such as [9] identify survival
36 functions obtained from the Cox PH model at these time points. Many previous publications directly address this approach. For
37 instance, [10] predicts kidney graft survival using tree-based models, [11] investigates several techniques, including random
38 forests and neural networks. In [3], the neural network’s ability to predict both graft and patient survival in kidney transplant
39 is compared to one of the regression techniques (such as Cox). Predictions at different time points were either modelled
40 independently or through multiple-output neural networks. [12] compares multilayer perceptrons and Bayesian networks, [13]
41 uses Bayesian belief networks. Whilst many of these previous studies demonstrate acceptable predictive performance, none
42 challenged their models’ validity through the lens of clinical interpretability.

43 Interpretability is another important criterion in the construction of a CDS tool for predicting graft outcomes. Informally,
44 interpretability is the extent to which the prediction of a model can be understood by a human [14]. This way, users can
45 build trust regarding the model’s results and remain in control of the associated outcomes. Moreover, a good model should
46 always be *intrinsically* interpretable to a certain degree. Indeed, interpretable models have been shown to be more robust to
47 adversarial attacks [15]. Unfortunately, this is not the case with the approaches mentioned above. Although the Cox PH model
48 is interpretable in the single risk case, its generalisation to competing risks is not [8]. It is possible to interpret *a posteriori* a
49 black box model with the help of *post hoc* interpretability methods. One can provide a local explanation of a given prediction.
50 For instance, LIME [16] locally samples data points around the input and returns a linear explanation of the predictions made by
51 the black-box model from these data points. Unfortunately, this solution is unstable; explanations depend highly on the sampled
52 data points, harming the trustworthiness of the explanations. Similarly to LIME, SHAP [17] is a *post hoc* interpretability method
53 relying on additive feature attribution models, i.e. linear functions as local explanation models. It provides explanations *via*
54 game theory: each prediction is seen as a game where the features are players contributing to that game. Feature contributions
55 are computed by considering all possible coalitions of features and the marginal contribution of each feature within these
56 coalitions. Hence, SHAP can be considered as a gold standard in terms of *post hoc* interpretability methods.

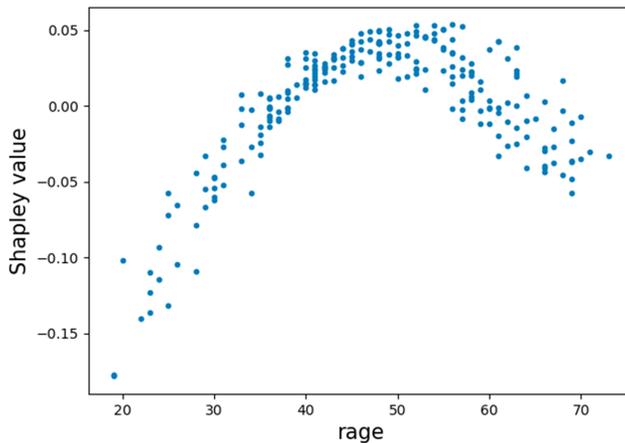
57 Results

58 Single Outcome Prediction

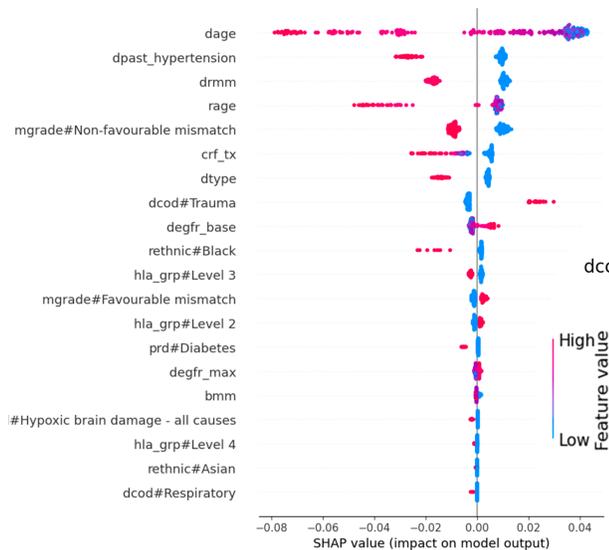
59 During the feature selection stage, we identified inconsistencies between SHAP interpretations of the obtained models and
60 clinical expectations. For example, survivability is expected to decrease with the number of times a patient has been transplanted,
61 which is not what we observed in the models produced. After further investigations, it appears that this feature is biased with all
62 non *primo* recipients having a successful graft. The data set has indeed been built from various heterogeneous sources, with
63 some outcomes not available for subsets of the data. From now on, we discard this feature. Finally, 15 and 10 features are



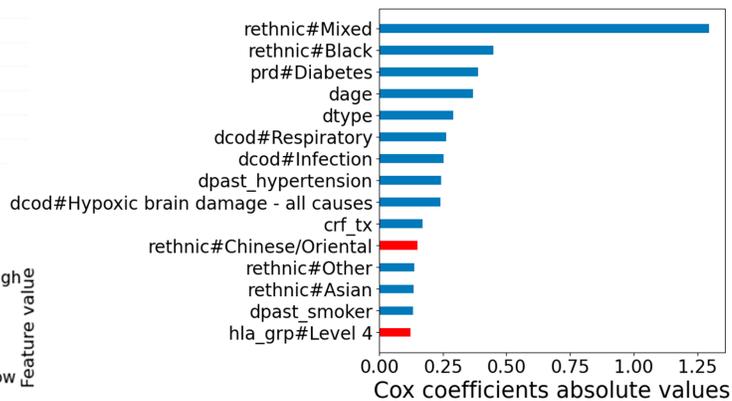
(a) Neural network's SHAP-based feature importance. A negative SHAP value indicates a negative impact on graft survival.



(b) Neural network's dependence on recipient age (*rage*)



(c) Random forest SHAP-based feature importance. A negative SHAP value indicates a negative impact on graft survival.



(d) Largest Cox PH model's coefficients. Blue and red bars represent positive and negative coefficients, respectively.

Figure 1. Interpreting single outcome prediction models. Both models have been trained to predict graft failure at 10 years.

		Random forest		Cox PH model		Neural network	
		AUROC	F1-Score	AUROC	F1-Score	AUROC	F1-Score
Year 1	Without feature selection	.62 ($\pm 2e^{-4}$)	.14 ($\pm 2e^{-6}$)	.61 ($\pm 1e^{-2}$)	.15 ($\pm 1e^{-5}$)	.62 ($\pm 2e^{-4}$)	.17 ($\pm 9e^{-5}$)
	With feature selection	.61 ($\pm 3e^{-4}$)	.14 ($\pm 4e^{-6}$)	.61 ($\pm 3e^{-4}$)	.15 ($\pm 1e^{-5}$)	.62 ($\pm 3e^{-4}$)	.18 ($\pm 7e^{-5}$)
Year 5	Without feature selection	.62 ($\pm 2e^{-4}$)	.34 ($\pm 7e^{-7}$)	.64 ($\pm 2e^{-4}$)	.37 ($\pm 5e^{-5}$)	.64 ($\pm 2e^{-4}$)	.37 ($\pm 2e^{-4}$)
	With feature selection	.60 ($\pm 1e^{-4}$)	.35 ($\pm 8e^{-6}$)	.63 ($\pm 1e^{-4}$)	.36 ($\pm 1e^{-4}$)	.63 ($\pm 1e^{-4}$)	.36 ($\pm 1e^{-4}$)
Year 10	Without feature selection	.68 ($\pm 2e^{-4}$)	.64 ($\pm 6e^{-5}$)	.71 ($\pm 1e^{-4}$)	.65 ($\pm 7e^{-5}$)	.71 ($\pm 1e^{-4}$)	.61 ($\pm 2e^{-4}$)
	With feature selection	.68 ($\pm 1e^{-4}$)	.62 ($\pm 5e^{-5}$)	.70 ($\pm 1e^{-4}$)	.63 ($\pm 4e^{-5}$)	.71 ($\pm 9e^{-5}$)	.61 ($\pm 1e^{-4}$)

Table 1. Performance of models for graft survival prediction. Graft failure ratios for years 1, 5, and 10 are 7%, 20%, 44%, respectively. Best scores are highlighted in bold.

		Random forest		Cox PH model		Neural network	
		AUROC	F1-Score	AUROC	F1-Score	AUROC	F1-Score
Year 1	Without feature selection	.73 ($\pm 4e^{-4}$)	.15 ($\pm 1e^{-4}$)	.75 ($\pm 3e^{-4}$)	.15 ($\pm 1e^{-4}$)	.71 ($\pm 4e^{-4}$)	.12 ($\pm 6e^{-5}$)
	With feature selection	.74 ($\pm 3e^{-4}$)	.15 ($\pm 3e^{-4}$)	.73 ($\pm 3e^{-4}$)	.14 ($\pm 8e^{-5}$)	.74 ($\pm 3e^{-4}$)	.12 ($\pm 4e^{-5}$)
Year 5	Without feature selection	.78 ($\pm 1e^{-4}$)	.42 ($\pm 1e^{-4}$)	.79 ($\pm 1e^{-4}$)	.43 ($\pm 1e^{-4}$)	.79 ($\pm 1e^{-4}$)	.41 ($\pm 2e^{-4}$)
	With feature selection	.76 ($\pm 1e^{-4}$)	.39 ($\pm 1e^{-4}$)	.77 ($\pm 1e^{-4}$)	.41 ($\pm 1e^{-4}$)	.76 ($\pm 1e^{-4}$)	.39 ($\pm 1e^{-4}$)
Year 10	Without feature selection	.80 ($\pm 1e^{-4}$)	.67 ($\pm 1e^{-4}$)	.82 ($\pm 9e^{-5}$)	.69 ($\pm 1e^{-4}$)	.82 ($\pm 9e^{-5}$)	.68 ($\pm 1e^{-4}$)
	With feature selection	.80 ($\pm 9e^{-5}$)	.66 ($\pm 1e^{-4}$)	.81 ($\pm 7e^{-5}$)	.66 ($\pm 1e^{-4}$)	.81 ($\pm 7e^{-5}$)	.66 ($\pm 1e^{-4}$)

Table 2. Performance of models for patient survival prediction. Patient death ratios for years 1, 5, and 10 are 3%, 14%, and 37%, respectively. Best scores are highlighted in bold.

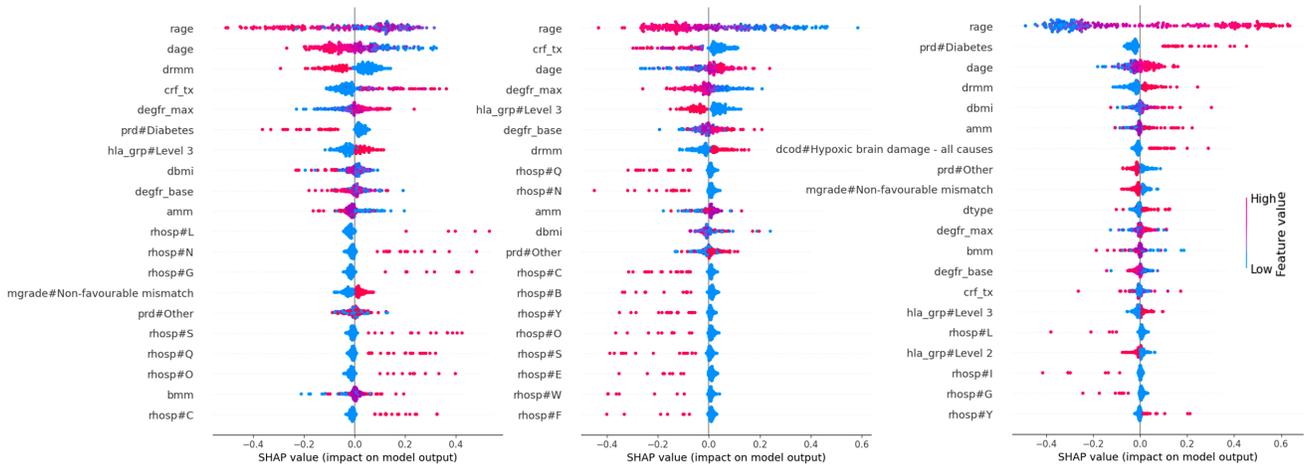
		Alive with functioning graft		Alive with failed graft		Dead	
		AUROC	F1-Score	AUROC	F1-Score	AUROC	F1-Score
Year 1	Without feature selection	.59 ($\pm 2e^{-4}$)	.93 ($\pm 8e^{-7}$)	.64 ($\pm 1e^{-3}$)	.09 ($\pm 2e^{-4}$)	.65 ($\pm 8e^{-4}$)	.10 ($\pm 1e^{-4}$)
	With feature selection	.59 ($\pm 1e^{-4}$)	.83 ($\pm 2e^{-5}$)	.59 ($\pm 3e^{-4}$)	.06 ($\pm 2e^{-4}$)	.69 ($\pm 6e^{-5}$)	.12 ($\pm 3e^{-5}$)
Year 5	Without feature selection	.64 ($\pm 6e^{-5}$)	.81 ($\pm 3e^{-5}$)	.58 ($\pm 3e^{-4}$)	.11 ($\pm 4e^{-4}$)	.74 ($\pm 4e^{-5}$)	.36 ($\pm 8e^{-5}$)
	With feature selection	.62 ($\pm 4e^{-5}$)	.74 ($\pm 2e^{-5}$)	.63 ($\pm 3e^{-4}$)	.09 ($\pm 3e^{-4}$)	.73 ($\pm 1e^{-4}$)	.38 ($\pm 2e^{-4}$)
Year 10	Without feature selection	.71 ($\pm 8e^{-5}$)	.70 ($\pm 8e^{-5}$)	.72 ($\pm 9e^{-4}$)	.18 ($\pm 1e^{-4}$)	.79 ($\pm 2e^{-5}$)	.63 ($\pm 6e^{-5}$)
	With feature selection	.71 ($\pm 6e^{-5}$)	.65 ($\pm 3e^{-5}$)	.71 ($\pm 5e^{-4}$)	.23 ($\pm 3e^{-4}$)	.79 ($\pm 3e^{-5}$)	.65 ($\pm 9e^{-5}$)

Table 3. Performance of multi outcome prediction models. Outcome ratios (alive with a functioning graft, alive with failed graft, and dead) for years 1, 5, and 10 are (.94, .03, .03), (.81, .15, .05), and (.56, .37, .07), respectively.

64 selected to predict graft failure and patient death, respectively. Notably, recipient and donor ages, primary renal disease, donor
65 eGFR, donor type, and the number of mismatches at DR locus are important features to predict both outcomes.

66 Tables 1 and 2 provide both the AUROC and the F1-Score of each model on predicting graft failure and patient death,
67 respectively, for observations years 1, 5, and 10. Performance before and after feature selection are presented. One can
68 observe that overall performances increase with the observation time, being maximal at year 10. More specifically, the neural
69 network has similar performances as the random forest and the Cox PH model, slightly outperforming them on the graft failure
70 prediction task.

71 From an interpretability viewpoint, the neural network, when combined with SHAP, provides a richer clinical depiction
72 of the data than Cox or the random forest. The features that are important to clinicians are also considered important to the
73 neural network. For example, among predominant features for graft failure prediction (cf. Figure 1.a), recipient and donor age,
74 donor type, donor past hypertension, or eGFRs are also features commonly used by regression models from the transplant
75 literature [4, 18, 19]. The effect of feature values on predictions also matches clinical knowledge. For instance, patients with
76 diabetes are likely to have inferior survival. This is reflected through the lower SHAP values regarding graft survival when
77 `prd#Diabetes` is equal to one. The effect of covariates on survivability can be non-linear, as illustrated by the recipient age
78 (`rage`; see Figure 1.b). Indeed, it is commonly recognised that younger patients can be less adherent to medication, hence
79 increasing the risk of transplant failure. This phenomenon vanishes with older patients, and age then becomes a penalising



(a) Alive with functioning graft (b) Alive with failed graft (c) Deceased

Figure 2. SHAP values for multiple outcome predictions at year 10.

80 feature for survivability. In contrast, explanations obtained from the Cox PH model or the random forest do not highlight such
 81 behaviours (see Figures 1.d and 1.c), being limited to less expressive covariate effects. By design, it can be summarised as a
 82 linear function in the case of Cox, and the random forest sometimes fails to represent relevant dependencies between survival
 83 and numerical values (e.g. the recipient’s age in Figure 1.c).

84 Multiple Outcome Prediction

85 In the multiple outcome case, 15 features are selected. Similarly to single outcome prediction, recipient and donor age, primary
 86 renal disease, eGFR, donor type, and number of mismatches at DR locus are present in this selection.

87 Table 3 stores the cause-specific AUROC and F1-Scores obtained by our neural network. Regarding AUROC, these results
 88 match the ones obtained in the single risk case. Notably, overall performances improve with t^* . However, one can notice
 89 that the multiple outcome prediction problem is more subject to class imbalance. For instance, patients that are *alive with a*
 90 *failed graft* represent 3% of the total uncensored population at year 1, 5% at year 5, and 7% at year 10. Thus, according to
 91 F1-Scores, the model performs better in predicting the classes *alive with functioning graft* and *dead*. This is consistent with the
 92 results observed in the single risk case as better prediction was obtained regarding patient death over graft failure. From the
 93 interpretability viewpoint, we retrieve clinically coherent SHAP values (cf. Figure 2). Notably, we obtain interpretations that
 94 are consistent with the ones obtained in a single risk setting. Similarly, our neural network can reflect non-linear covariate
 95 effects.

96 Discussion

97 Neural networks have shown comparable performances to tools generally used by clinicians when predicting kidney transplant
 98 outcomes. In particular, they perform well when predicting long-term outcomes, which is a necessary property when considering
 99 the acceptance of an organ offer. The Cox PH model remains a robust solution in terms of performance, with little to no
 100 hyperparameter tuning. It is simple to use, leads to reliable predictions, and is easy to understand by clinicians. However,
 101 despite their black-box nature, neural networks stand out in terms of interpretability. Indeed, using SHAP allows us to have
 102 fine-grain interpretations of these models, which is not possible with the Cox PH model or random forests. This level of
 103 interpretability allows us to clinically validate these models, making them more trustworthy and explainable to patients. SHAP

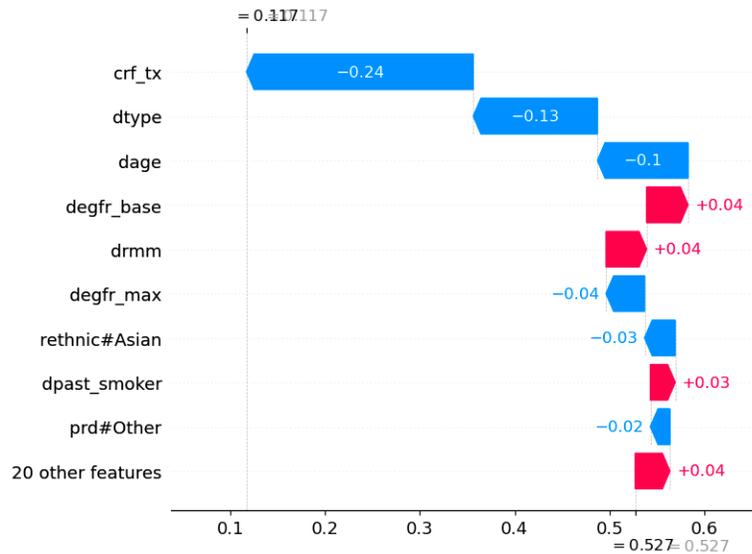


Figure 3. Explanation of a given prediction of graft survival at 10 years. The calculated reaction frequency at transplant is equal to 84%; donation occurred after circulatory death; the donor was 56 years old and non-smoker; eGFR remained equal to 113; no mismatches at the DR locus has been recorded; finally, the recipient is Asian with diabetes. Negative SHAP values indicate a negative impact on survivability, and *vice versa*.

can also highlight interesting relationships between covariates and transplant outcomes. Previous analyses of the UK registry data show that outcomes from kidney donations before and after circulatory death are equivalent regarding both patient and graft survival [20, 21]. Nonetheless, our models, trained from a larger dataset, suggest that donor type can have an impact on long-term transplant outcomes (see Figure 1.d, dtype). In practice, this level of interpretability is useful to explain individual prediction through the lens of SHAP. For example, Figure 3 shows the contribution of each features to a given prediction of graft survival at year 10. In this case, predicted survivability is mainly lowered by the calculated reaction frequency at transplant and the donor type. Neural networks can also deal with multiple outcomes, providing a more comprehensive prediction of the future state of the recipient, where patient death and graft failure are modelled jointly. There may be some clinical relevance to distinguishing the state of graft function at the time of death (*i.e.* death with or without graft failure). Unfortunately, the data set is too unbalanced and the population *dead with graft failure* is not large enough to provide such a distinction. These outcomes are not competing: a patient being alive with a failed graft at some point could die later, which remains an event of interest.

To conclude, we have trained several models to predict transplant outcomes from kidney offers, based on twenty years of registry data. Neural networks provide comparable results to classic survival analysis models, and can be easily extended to multiple outcome prediction. By using SHAP, we provide clinically validated interpretations of these models. This level of interpretability is especially relevant to enable validation from clinicians and to involve patients in the decision-making process. Therefore, neural networks represent a promising core component in the construction of future CDS for transplant offering. However, predicting transplant outcomes is only one aspect of the construction of a CDSS for kidney offering. Predicting what could be the consequences of refusing an organ offer in terms of future transplant opportunities, death, or removal from the waiting list is another key step. Having a good understanding of the outcomes in both scenarios is indeed necessary to predict individualised treatment effects. Uncertainty quantification is another critical research direction regarding the construction of a CDS tool for organ offering. Indeed, it can improve the trustworthiness of the tool by giving more insights about how difficult a given prediction is, and why. This can be achieved through post hoc error prediction using meta-modeling.

126 **Methods**

127 **Data**

128 Our work is based on the analysis of a data set from the UK Transplant Registry, provided with ethical approval by NHSBT. It
129 describes 36,653 accepted kidney transplants, which have been performed between the years 2000 and 2020, across 24 UK
130 transplant centres. The total follow-up duration is around 22 years. Each transplant is originally described with 3 identifiers, 12
131 immunosuppression follow-up indicators, 143 donor, recipient and transplant characteristics, and 7 entries describing targeted
132 outcomes. Considering transplants as independent, we exclude the transplant, donor, and recipient identifiers. Additionally,
133 information regarding post-transplant immunosuppression is discarded as this is not available at the time of the offer decision.
134 The donor, recipient and transplant characteristics serve as input features for modelling. Among them, 24 describe the
135 recipient, 109 represent the donor, and 10 refer to the overall transplant. Both recipient and donor characteristics contain
136 generic information such as gender, ethnicity, age, blood group, height, weight, or body mass index (BMI). More specific
137 information is also available, such as the transplant centre, number of previous transplants, waiting time, ease of matching,
138 and the dialysis status. Donor data include the cause of death, past medical history and results of blood tests including kidney
139 function (estimated glomerular filtration rate, eGFR). Transplant data include the donor-recipient immunological match.

140 Duplicate rows are removed, and we ensure that numerical values are within a plausible clinical range. Categorical values
141 are checked by clinicians and simplified (or removed) if needed. BMI is recomputed based on weight and height. Both weights
142 and heights are discarded to limit redundant information. Blood measurements are harmonised across the data set by selecting
143 for each transplant the first measurement ever taken (generally at registration) and the maximum value ever recorded. Since the
144 calculation of eGFR varies across hospitals, this metric is recomputed over the whole data set using a consistent definition
145 (see appendix). Recipient dialysis status is also simplified into a dialysis duration and dialysis modality at time of transplant
146 (predialysis, haemodialysis or peritoneal dialysis). Transplant offers not meeting the inclusion criteria, such as those leading to
147 the transplantation of multiple organs, are discarded.

148 Outcomes present in the dataset include information about graft failure, patient death, and transplant failure. Transplant
149 failure denotes either the graft failure or death. Each outcome is represented as a pair containing an event time and a right-
150 censoring indicator. Right-censoring is a common type of censoring in survival analysis that describes the loss of follow-up on
151 the event of interest. It can occur for various reasons, such as the end of the study, competing events, etc. Thus, right-censored
152 information provides some partial information about the survival time, where it is only known to be greater than the censoring
153 time. Minor missing censored indicators related to patient death are thus imputed based on graft information, and transplant
154 outcomes are recomputed for the sake of consistency.

155 After removing the features presenting more than 5% missing values across the whole data set, and excluding any offer
156 containing missing information, the resulting data set contains 25,370 transplants described through 45 input variables. A
157 summary of the data-cleaning process is given in Figure 4. Additionally, an exhaustive list of the features and targets considered
158 at the latest stage of this process is given in appendix.

159 **Methodology**

160 In this article, we first compare the Cox PH model to classification methods in a single risk setting. Subsequently, multiple
161 outcome predictions are conducted by employing neural networks. The different models are interpreted *a posteriori*, and their
162 performances are discussed. In both cases, the following methodology is applied. First, numerical values are standardised and
163 categorical ones are one-hot-encoded. Standardisation appears to be more relevant than normalisation due to the presence of
164 outliers in the data. Training is then performed through 5-cross validation on 80% of the data. Finally, the relevant performance
165 metrics are computed and averaged from 100 bootstraps of the remaining 20%. The split into training and test data is done
166 randomly, in a stratified manner with regards to censoring indicators. Due to matching policy changes and follow-up time
167 differences between training and testing cohorts, we cannot split the data according to transplant dates. Classifiers are clinically
168 interpreted using SHAP. When relevant, the coefficients of intrinsically interpretable models are also investigated (Code used

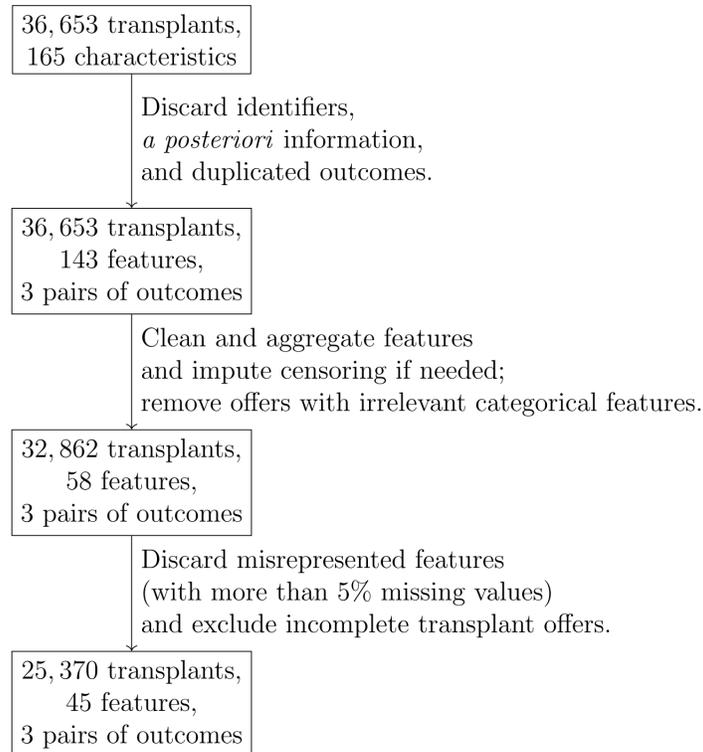


Figure 4. Data pre-processing steps.

169 for experiments can be found at <https://github.com/AchilleSalaun/Xamelot>).

170 **Single Outcome Prediction**

171 For a single type of event interest (e.g. graft failure), we want to predict the occurrence of that event before an arbitrary time
 172 point t^* . The Cox PH model and classifiers consider different kinds of input data, tackling different formulations of that problem.
 173 Cox returns a survival function, describing the probability of survival (as in the absence of an event) with regard to time.
 174 Therefore, the probability of an event at time t^* can be obtained by evaluating the survival function at that time. Conversely, the
 175 censored data must be converted into labels to be fed into classifiers. Let us denote the pair (event time, censoring indicator)
 176 by (t, c) . Then, the graft is *functioning* (or the patient is *alive*) if $t^* < t$; the graft *failed* (the patient is *dead*) if $t^* \geq t$ and the
 177 event has not been censored ($c = 1$); finally, the status of the graft (or patient) is *unknown* if $t^* \geq t$ and the event has been
 178 censored ($c = 0$) (see Figure 5). Dropping unknown labels induces a binary classification problem since the graft is now either
 179 functioning or failed, or the patient is either alive or dead. This last operation generally assumes censoring and events to be
 180 independent, which often leads to biases in practice. However, the Cox PH model relies on the same assumption: likelihood’s
 181 maximisation is achieved by managing a risk set over time, that is a set of subjects that are still under follow-up [6]. Hence,
 182 censored events with lost follow-up are implicitly dropped while training Cox PH models. The shorter the observation time
 183 t^* , the more imbalanced the outcome distribution is, with failure or death being under-represented regarding survival. For
 184 instance, the class imbalance goes from 7% of graft failures at 1 year, to 44% at 10 years. As we drop censored events with a
 185 censoring time that is lower than the observation time, the number of data points used for training also varies with regard to
 186 time. Thus, the training data sets used for predicting transplant outcomes at years 1 and 10 comprise 23,422 and 10,017 data
 187 points, respectively.

188 We can now compare the respective abilities of the Cox PH model, random forests, and neural networks to predict transplant
 189 outcomes using the area under the receiver operating characteristic (AUROC). The choice of AUROC is motivated by two
 190 aspects. First, it is conceptually close to concordance which is the metric generally used for survival analysis. Second, it is a

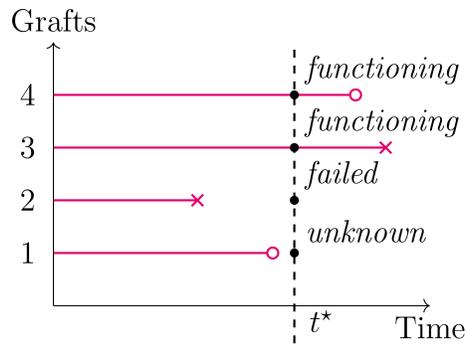


Figure 5. Translation of a single risk survival analysis problem into a classification task. Event times are given by the x-axis; \times indicates the observation of an event, \circ indicate censoring. For any time t^* , graft (or patient) status can be derived from survival data.

191 good metric when dealing with balanced data, which is the case when predicting transplant outcomes after 10 years. Therefore,
 192 this metric is relevant since we have a particular interest in predicting long-term outcomes. However, as class imbalance
 193 remains a recurrent difficulty, we also compute the F1-Score. For a given model, we select the classification threshold that
 194 leads to the best possible F1-Score. From an instantiation viewpoint, Breslow’s estimator is used to derive the Cox PH model’s
 195 baseline [22]. In addition, a regularisation parameter is introduced and set to $1e^{-4}$ to deal with colinearities in the data. The
 196 random forest [23] contains 1,000 trees, relies on Gini’s criteria, and adjusts class weights automatically. To predict graft
 197 survival (or patient death) at years 1, 5, and 10, a neural network is instantiated with one hidden layer of 400 (200, respectively)
 198 neurons, activated by a ReLU, and with 10% dropout. The loss is a weighted cross to handle class imbalance. Finally, the
 199 training is done through 20 epochs, with batches of size 8 (32, respectively), using RMSProp and a learning rate equal to $1e^{-4}$.

200 Feature selection is performed after preliminary training. We first inspect the effect of each feature on prediction to detect
 201 potential inconsistencies in the data. Then, we temporarily add random noise: features that are shown to be less important in the
 202 prediction than noise are discarded. Finally, we progressively remove the less relevant features from the set of selected features
 203 until a noticeable decrease in performance is observed. For a given type of outcome, the set of selected features corresponding
 204 to year 10 includes important features for prediction in earlier years; therefore we use the same set of features for years 1 and 5.

205 **Multiple Outcome Prediction**

206 For further analysis, we generalise our approach to multiple outcomes prediction. For an arbitrary time point t^* , we want to
 207 know whether the patient is *alive with a functioning graft*, *alive with a failed graft*, or *dead*. The construction of these labels
 208 from censored data is similar to the one performed in the single risk case. Transplants with unknown status due to censoring are
 209 discarded. To address this new problem, we focus on neural networks as they appeared to be a promising approach in the single
 210 risk case (see Results and Discussion). We train a neural network with one hidden layer of 1000 neurons activated by ReLU. A
 211 dropout is set to 10%. The loss is a weighted cross entropy and the training is done through 100 epochs, with batches of size 64.
 212 Optimisation relies on RMSProp with a learning rate set to $1e^{-3}$. Similar to the single risk case, feature selection is conducted
 213 after initial training using all features.

214 **Data availability statement**

215 The dataset supporting this study is the property of NHSBT (National Health Service Blood and Transplant). Due to the
 216 presence of patient data that could potentially lead to identification, obtaining prior ethical approval from NHSBT is mandatory.
 217 Consequently, the supporting data is not accessible for external use.

References

- 218 **1.** M. Chun, R. Clarke, B. J. Cairns, D. Clifton, D. Bennett, Y. Chen, Y. Guo, P. Pei, J. Lv, C. Yu, *et al.*, “Stroke risk prediction
219 using machine learning: a prospective cohort study of 0.5 million Chinese adults,” *Journal of the American Medical*
220 *Informatics Association*, vol. 28, no. 8, pp. 1719–1727, 2021.
- 221 **2.** D. W. Kim, S. Lee, S. Kwon, W. Nam, I.-H. Cha, and H. J. Kim, “Deep learning-based survival prediction of oral cancer
222 patients,” *Scientific reports*, vol. 9, no. 1, pp. 1–10, 2019.
- 223 **3.** R. S. Lin, S. D. Horn, J. F. Hurdle, and A. S. Goldfarb-Rumyantzev, “Single and multiple time-point prediction models in
224 kidney transplant outcomes,” *Journal of biomedical informatics*, vol. 41, no. 6, pp. 944–952, 2008.
- 225 **4.** P. S. Rao, D. E. Schaubel, M. K. Guidinger, K. A. Andreoni, R. A. Wolfe, R. M. Merion, F. K. Port, and R. S. Sung, “A
226 comprehensive risk quantification score for deceased donor kidneys: the kidney donor risk index,” *Transplantation*, vol. 88,
227 no. 2, pp. 231–236, 2009.
- 228 **5.** A. J. Vinson, B. A. Kiberd, R. B. Davis, and K. K. Tennankore, “Nonimmunologic donor-recipient pairing, HLA matching,
229 and graft loss in deceased donor kidney transplantation,” *Transplantation direct*, vol. 5, no. 1, 2019.
- 230 **6.** D. R. Cox, “Regression models and life-tables,” *J R Stat Soc*, 1972.
- 231 **7.** J. L. Katzman, U. Shaham, A. Cloninger, J. Bates, T. Jiang, and Y. Kluger, “DeepSurv: personalized treatment recommender
232 system using a cox proportional hazards deep neural network,” *BMC Med. Res. Methodol.*, 2018.
- 233 **8.** H. Putter, M. Fiocco, and R. B. Geskus, “Tutorial in biostatistics: competing risks and multi-state models,” *Stat. Med.*,
234 2007.
- 235 **9.** “NHS risk communication tools.” [https://www.odt.nhs.uk/transplantation/
236 tools-policies-and-guidance/risk-communication-tools/](https://www.odt.nhs.uk/transplantation/tools-policies-and-guidance/risk-communication-tools/).
- 237 **10.** S. Krikov, A. Khan, B. C. Baird, L. L. Barenbaum, A. Leviatov, J. K. Koford, and A. S. Goldfarb-Rumyantzev, “Predicting
238 kidney transplant survival using tree-based modeling,” *Asaio Journal*, vol. 53, no. 5, pp. 592–600, 2007.
- 239 **11.** S. A. A. Naqvi, K. Tennankore, A. Vinson, P. C. Roy, and S. S. R. Abidi, “Predicting kidney graft survival using machine
240 learning methods: prediction model development and feature significance analysis study,” *Journal of Medical Internet*
241 *Research*, vol. 23, no. 8, p. e26843, 2021.
- 242 **12.** V. Rao, R. S. Behara, and A. Agarwal, “Predictive modeling for organ transplantation outcomes,” in *2014 IEEE International*
243 *Conference on Bioinformatics and Bioengineering*, pp. 405–408, IEEE, 2014.
- 244 **13.** K. Topuz, F. D. Zengul, A. Dag, A. Almehti, and M. B. Yildirim, “Predicting graft survival among kidney transplant
245 recipients: A bayesian decision support model,” *Decision Support Systems*, vol. 106, pp. 97–109, 2018.
- 246 **14.** C. Molnar, *Interpretable machine learning*. self published, 2020.
- 247 **15.** A. Noack, I. Ahern, D. Dou, and B. Li, “An empirical study on the relation between network interpretability and adversarial
248 robustness,” *SN comput. sci.*, 2021.
- 249 **16.** M. T. Ribeiro, S. Singh, and C. Guestrin, ““Why should I trust you?” Explaining the predictions of any classifier,” in
250 *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144,
251 2016.
- 252 **17.** S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” *Advances in neural information*
253 *processing systems*, vol. 30, 2017.
- 254 **18.** C. J. Watson, R. J. Johnson, R. Birch, D. Collett, and J. A. Bradley, “A simplified donor risk index for predicting outcome
255 after deceased donor kidney transplantation,” *Transplantation*, vol. 93, no. 3, pp. 314–318, 2012.
- 256

- 257 **19.** M. Z. Molnar, D. V. Nguyen, Y. Chen, V. Ravel, E. Streja, M. Krishnan, C. P. Kovesdy, R. Mehrotra, and K. Kalantar-Zadeh,
258 “Predictive score for posttransplantation outcomes,” *Transplantation*, vol. 101, no. 6, p. 1353, 2017.
- 259 **20.** D. M. Summers, R. J. Johnson, A. Hudson, D. Collett, C. J. Watson, and J. A. Bradley, “Effect of donor age and cold
260 storage time on outcome in recipients of kidneys donated after circulatory death in the UK: a cohort study,” *The Lancet*,
261 vol. 381, no. 9868, pp. 727–734, 2013.
- 262 **21.** D. M. Summers, C. J. Watson, G. J. Pettigrew, R. J. Johnson, D. Collett, J. M. Neuberger, and J. A. Bradley, “Kidney
263 donation after circulatory death (DCD): state of the art,” *Kidney International*, vol. 88, no. 2, pp. 241–249, 2015.
- 264 **22.** D. Lin, “On the Breslow estimator,” *Lifetime data analysis*, vol. 13, pp. 471–480, 2007.
- 265 **23.** L. Breiman, “Random forests,” *Machine learning*, vol. 45, pp. 5–32, 2001.

266 Legends

267 **Figure 1** Interpreting single outcome prediction models. Both models have been trained to predict graft failure at 10 years.

268 **Table 1** Performance of models for graft survival prediction. Graft failure ratios for years 1, 5, and 10 are 7%, 20%, 44%,
269 respectively. Best scores are highlighted in bold.

270 **Table 2** Performance of models for patient survival prediction. Patient death ratios for years 1, 5, and 10 are 3%, 14%, and
271 37%, respectively. Best scores are highlighted in bold.

272 **Table 3** Performance of multi outcome prediction models. Outcome ratios (alive with a functioning graft, alive with failed
273 graft, and dead) for years 1, 5, and 10 are (.94, .03, .03), (.81, .15, .05), and (.56, .37, .07), respectively.

274 **Figure 2** SHAP values for multiple outcome predictions at year 10.

275 **Figure 3** Explanation of a given prediction of graft survival at 10 years. The calculated reaction frequency at transplant is
276 equal to 84%; donation occurred after circulatory death; the donor was 56 years old and non-smoker; eGFR remained
277 equal to 113; no mismatches at the DR locus has been recorded; finally, the recipient is Asian with diabetes. Negative
278 SHAP values indicate a negative impact on survivability, and *vice versa*.

279 **Figure 4** Data pre-processing steps.

280 **Figure 5** Translation of a single risk survival analysis problem into a classification task. Event times are given by the x-axis; ×
281 indicates the observation of an event, ○ indicate censoring. For any time t^* , graft (or patient) status can be derived from
282 survival data.

283 Acknowledgements

284 The authors thank the anonymous reviewers for their valuable suggestions. This work has been supported by funds from the
285 NIHR (AI Award 2020 Phase 1: AI_AWARD02316). T.Z. was supported by the Royal Academy of Engineering under the
286 Research Fellowship scheme.

287 Author contributions statement

288 A.S. undertook the data cleaning, model building, and redaction of this paper. S.K. and L.W. provided clinical input to data
289 cleaning and model design. S.K. and T.Z. co-ordinate the overall project, providing respectively clinical and machine learning
290 insights.

291 Additional information

292 No competing interest is declared.