

# Enova Technical Assessment Solution

---

## Data Loading & Evaluation

Use Python pandas and numpy packages to load data and do some basic manipulation.

Check the dataset statistics summary and find out anything interesting.

- There are some missing values(NaN)
- Feature 'symptoms' is a tricky feature, as each one has multiple values
- Sanity check, is the data balanced? Eg: if survival\_1\_years = 0, then survival\_7\_years must be 0
- How are the features distributed?
- Are the features highly correlated?
- Which features are important?

## EDA (Exploratory Data Analysis)

This section I solved the questions from last section.

- For missing values I set it into 0 or mean value depends on the feature itself
- For 'symptoms', I split the multiple values first, then do one hot encoding
- The labels(survival\_7\_years) is almost balanced with ratio 0.43. And the survival\_1\_years, survival\_7\_years don't conflict with each other
- The training data and score data's features followed the same distribution. So the training data should reflect the score data set condition
- There are some features correlated with each other, but not a very big problem here
- After running machine learning model, I gave the feature importance plot as well which showed the symptoms features are the most important ones

## Model Selection

There are many classification models to choose, we need to know that the dataset has a lot of categorical features. So we don't consider logistics regression which is not good dealing with too many categorical data.

Instead, I picked random forest(ensemble method), xgboost(ensemble method), svm, svm-smote to train the model respectively.

Random forest is an ensemble method which can reduce the model bias and unlikely to overfit. Xgboost is a boosting model which are quite fast and robust. Svm is good for model which has a lot of features. As the data is not super balanced, I tried svm-smote to automatically fixed the imbalance data issue.

Finally I picked Xgboost as final model which has relatively highest performance which has test set accuracy 68.35%.

## Some Data Visualizations



