



UNIVERSITY OF EDINBURGH
Business School

Predictive Analytics and Modelling of Data

CMSE11428 (2020-2021)

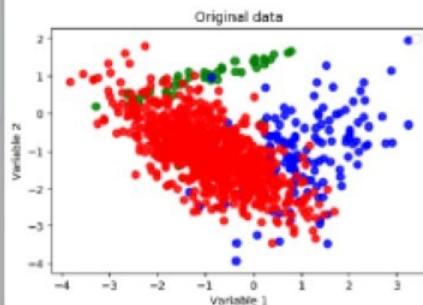
Dr Xuefei Lu

The University of Edinburgh Business School

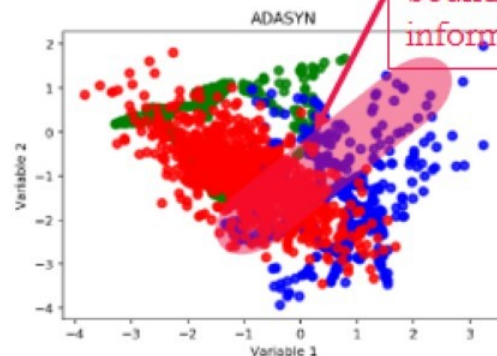
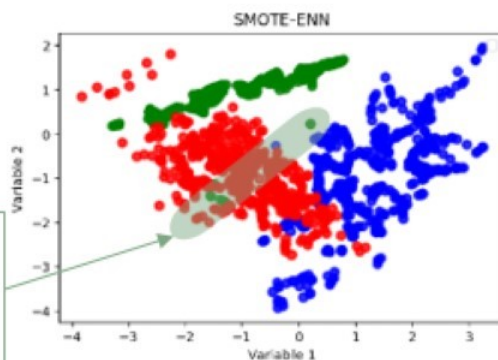
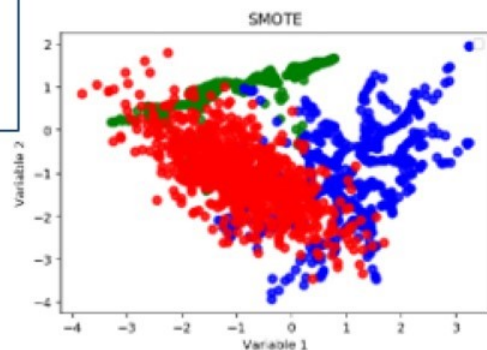
A photograph of a modern, multi-story building with a glass and metal facade, identified as the University of Cambridge Business School. The building features large windows and a prominent corner structure. A sign on the ground floor reads "UNIVERSITY OF CAMBRIDGE Business School".

Under and Over sampling

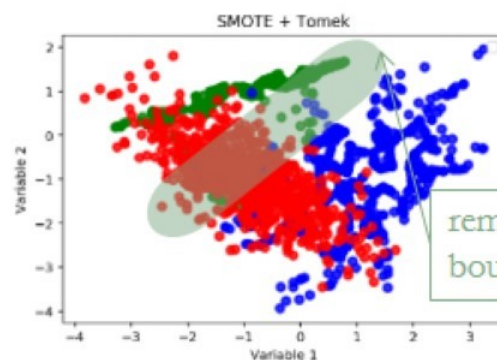
Create minorities
based on
N.neighbours



reduce points at the
boundaries → less
ambiguity



increase points at the
boundaries → more
information



remove points at the
boundaries

Training, test, and validation sets



Build models

Evaluation



Build models

Model selection/
parameter tuning

Evaluation

The modelling process

Would you apply these to the whole data set, the training and/or test set?

- Missing values
- Outliers
- Normalisation
- Over/under-sampling
- PCA

The modelling process

The main idea is not to allow the model to gain insight from the test data. We do not want information from the test/validation set "leak" into your training data!

- Normalisation, PCA:
Same parameters used on your training data should be used on test/validation sets
- Outliers, missing values: depends how you deal with them
- Separate case:
 - Over/under-sampling: typically only done on the training set

- Some packages provide CV estimation of model parameters:
- The idea behind CV estimation method is to find the set of parameters, e.g. β , if apply that minimizes the cross-validation error.

$$\epsilon_{CV}(\beta, \sigma^2, \theta, \sigma_n^2; \mathcal{Y}) = \frac{1}{N} \sum_{k=1}^K \epsilon_{CV,k}.$$

- In our case, additional test set is necessary for a final evaluation to tune model parameters
- https://scikit-learn.org/stable/modules/cross_validation.html



UNIVERSITY OF EDINBURGH
Business School