

Smote and cross validation

31 May 2023 20:54

Q How are say data is Balanced?

In diabetes data, if we have equal number of observations related to diabetic & non-diabetic, then we say data is Balanced.

Eg: Diabetic → 150
Non Diabetic → 150.

In balanced data → If data has equal number of observations related to diabetic or non-diabetic
Diabetic → 800
Non Diabetic → 200

What happens if data is Imbalanced?

Diabetic → 800 → Majority class

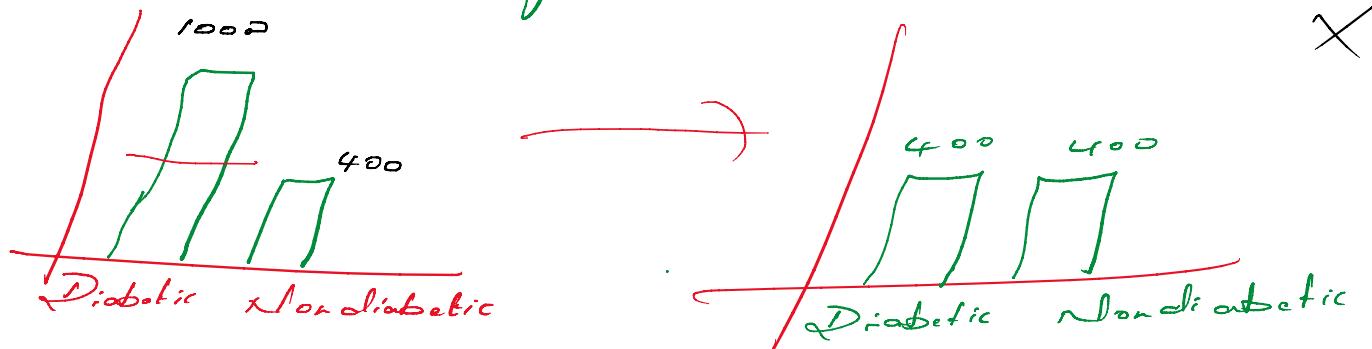
Non Diabetic → 200 → Minority class.

* Such data is Imbalanced the model learns more about the majority class and learns less about minority class since it has very less information.

* Because of this issue model makes wrong predictions.

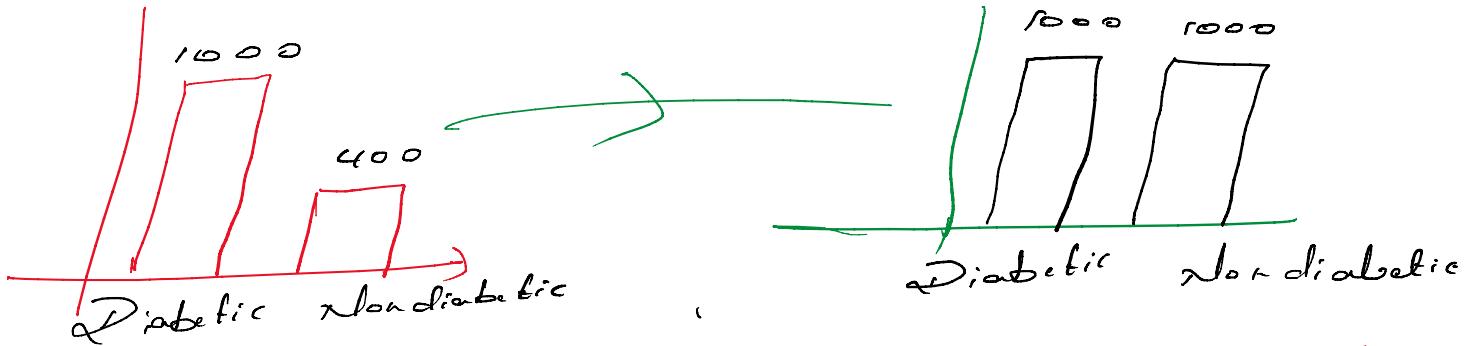
How to Balance data?

① Under sampling



Under sampling is not a good option because data can't be lost.

② Over sampling



- * Overampling means increasing number of observations
- \rightarrow It is minority class by duplicating the sample.
- * Overampling fails because it increases the duplicates.

Smote

- It is the very popular method used to balance data.
- SMOTE - Synthetic minority overampling technique.
- * Smote is used to create synthetic/artificial data from the original data.
 - * It uses KNN and interpolation to create new data.

Note

Balancing should always be done on training data.

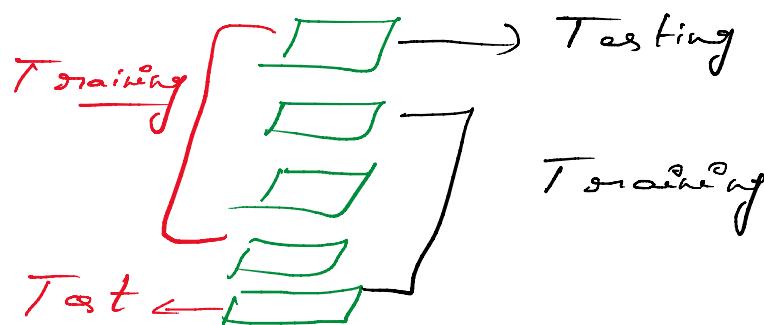
Cross Validation

It is one of the oversampling technique which divides data for training and testing.

- K-Fold Cross Validation
- * It divides entire data into K parts of equal size.
 - * $K \rightarrow$ No of Folds.

$$K = 5 \text{ or } 10$$

$$K = 5$$



- * The first set is selected for testing and remaining sets used for training.

- * The first set is selected for testing and remaining $K-1$ sets used for training.
- * It will calculate the error made by the model.
- * The above process repeats multiple times.

$K=5 \rightarrow 5$ Error.

Imp

Model is good only if it gives similar scores for all K .

Note

Using cross-validation we are validating model using different possible datasets using K -Fold.