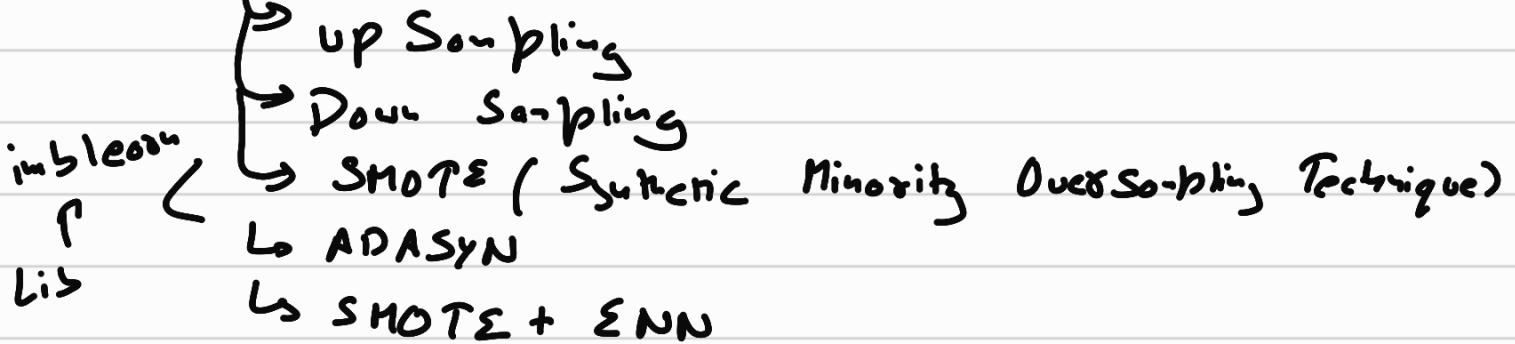


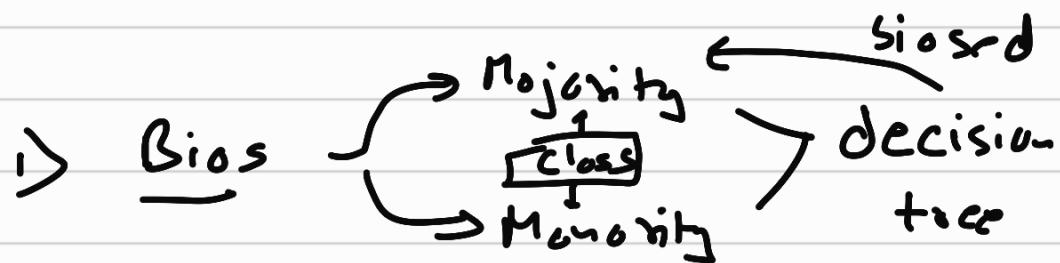
Approaches :- → applied on training data set

① Data Level

② Algorithm

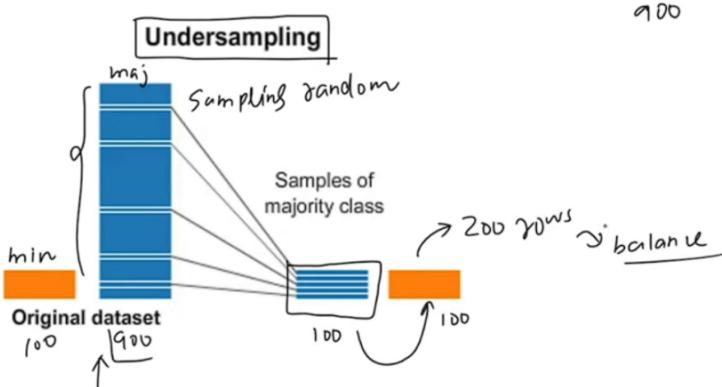


Problem With Imbalanced data:-



2) Metrics are → Accuracy
Not reliable

Under Sampling:-



Adv:-

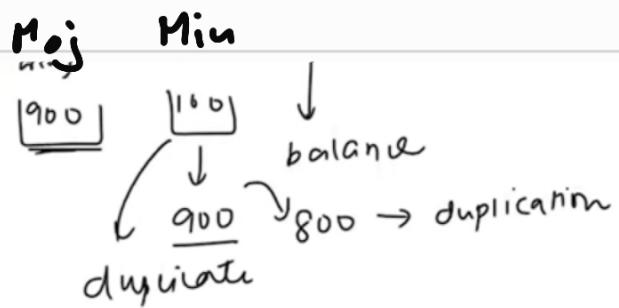
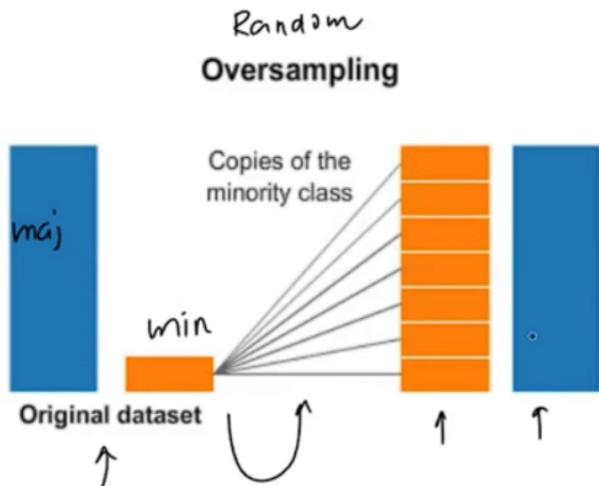
- ① Reduction in bias
- ② Faster training

Disadv:-

- ① Info loss leading to underfitting
- ② Sampling Bias

Oversampling :-

03 May 2024 01:13



Add :-

① Reduce Bias

Disadv

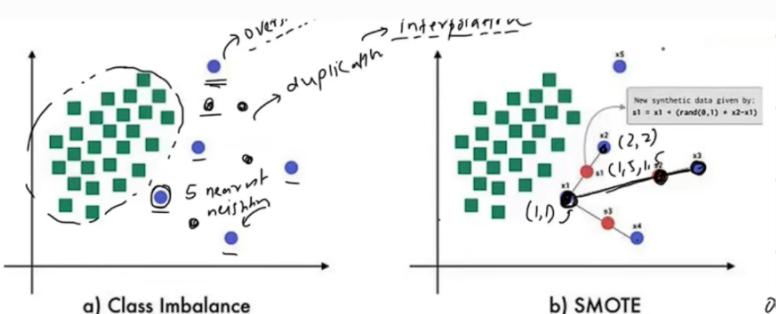
- ① Inc Size
- ② Duplication of data may cause Overfitting

③ SMOTE :-

Min Class \rightarrow upsample]

But not by
duplications
↓

We generate
new data points



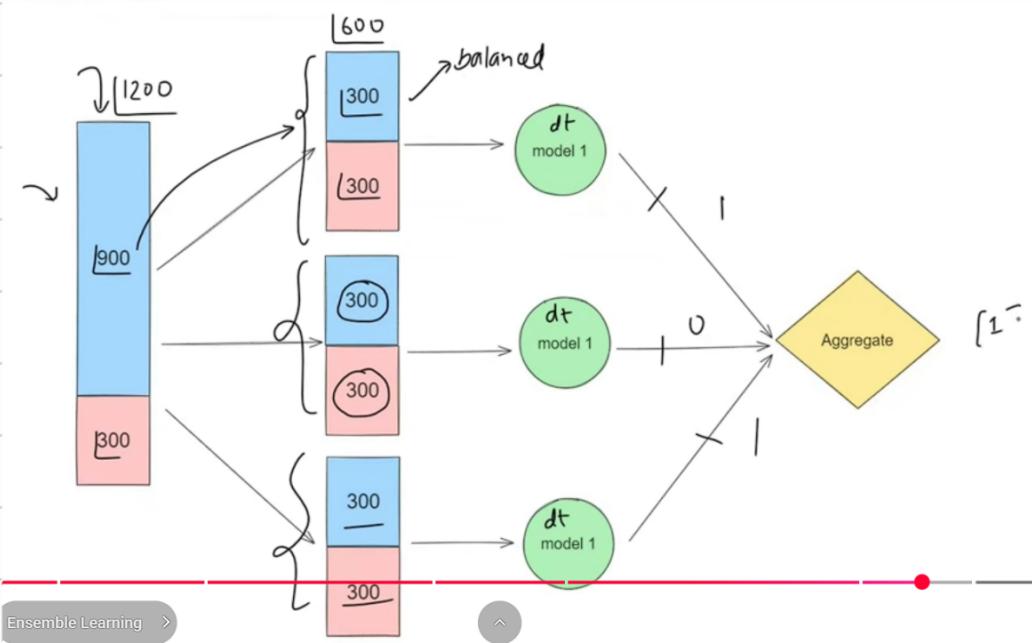
- 5 KNN
- Train a KNN on minority class observations - find each observation's 5 closest neighbours
- To create the new synthetic data:
- Select examples from the minority class at random
 - Select a neighbour of each example at random (for the interpolation)
 - Extract a random number between 0 and 1
 - Calculate the new examples as $(\text{original sample} - \text{factor} * (\text{original sample} - \text{neighbour}))$
 - The final dataset consists of the original dataset + the newly created examples

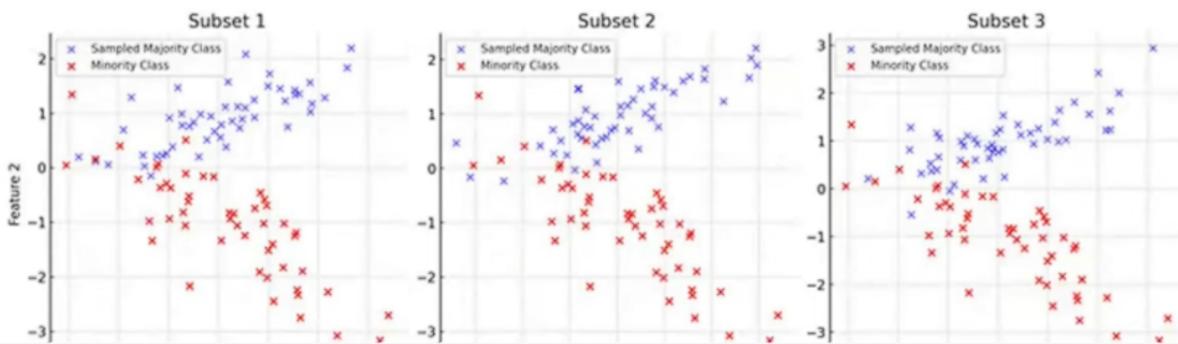
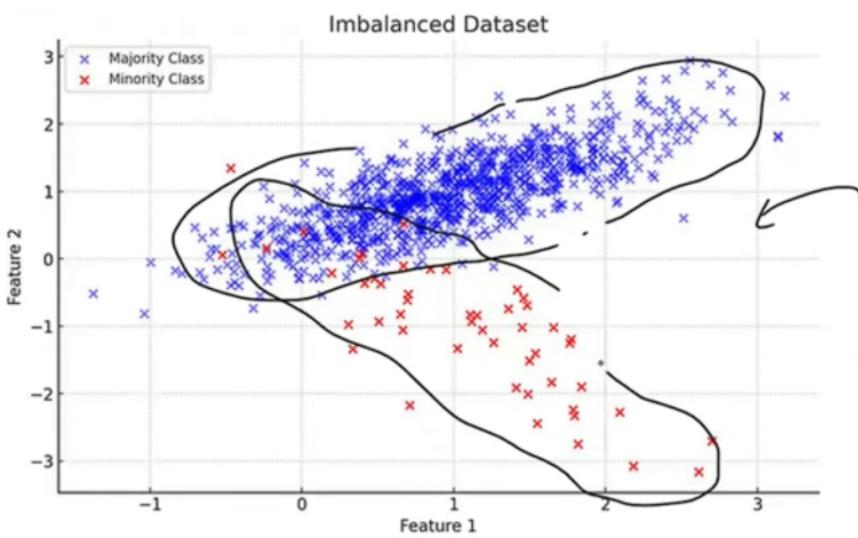
$$\begin{aligned} & \text{sample} - \underbrace{\text{factor}}_{\rightarrow 0.5} \times [\text{sample} - \text{neigh}] \\ & [1, 1] - 0.5 [(1, 1) - (2, 2)] \\ & (1, 1) - 0.5 [(-1, -1)] \\ & (1, 1) - (-0.5, -0.5) \\ & = 2 [1.5, 1.5] \end{aligned}$$

Disadvantages:-

- ① Do not handle Categorical data well
- ② Computational Complexity
- ③ Dependency on Choice of Neighbors
- ④ Sensitive to Outliers
- ⑤ Balanced dataset may not reflect True Nature

Ensemble Methods :-





⑤ Cost Sensitive Learning:-

Changes in Learning process:-

① Class Weights:-

$$\begin{matrix} 1 & 0 \\ 900 & 100 \end{matrix} \text{ weight} \rightarrow$$

$$\begin{bmatrix} 2 & 10 \end{bmatrix} =$$

Costs same for class 0
it will be weighted to
this the original cost
→ Model will take seriously.

⑥ Use Custom Loss Function:-

GB

Algo which give this feature:
↳ Xgboost
↳ Lightgbm