

Handling Imbalanced Dataset.

```
df["Virus Present"].value_counts()
```

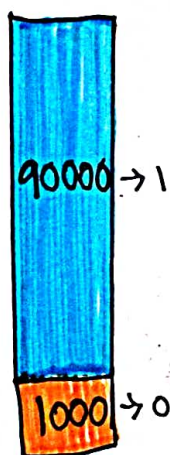
Out: 1 200
0 200

— Balanced Data Set

EX: 50-50%
(or)
45-55%

(+, - = 5%)

1. Under Sampling major class



EX: finance [cibil]

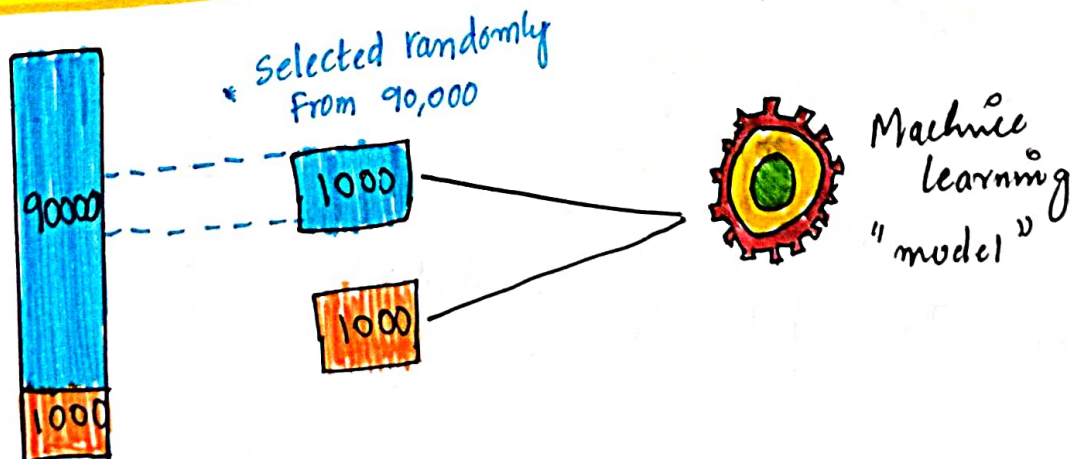
99% [pay]

1% [don't pay]

If we apply classification Technique

* it gives always answer as "1"

* Under Sampling

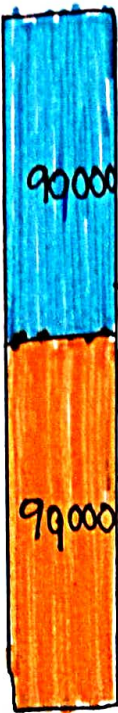


2. Over Sampling minority class by duplication

* Generate new Sample from Current Sample by Simply duplicating them



$1000 \times 99 =$
times
duplicates



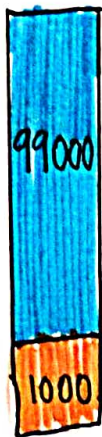
Duplication

same record we
duplicate

But, it gives wrong
answers (overfitting)

Machine
Learning "model"

3. Over Sampling minority class using SMOTE



machine
Learning
model.

* Generate Synthetic examples using "KNN"

* SMOTE : Synthetic minority Over Sampling
Technique.
Synthetic data

Ex:- Synthetic data.

age : 22

next record, we create

age : 22.1 (which is almost close to 22.
but not same data)

Ex: 2

22

21

24

22

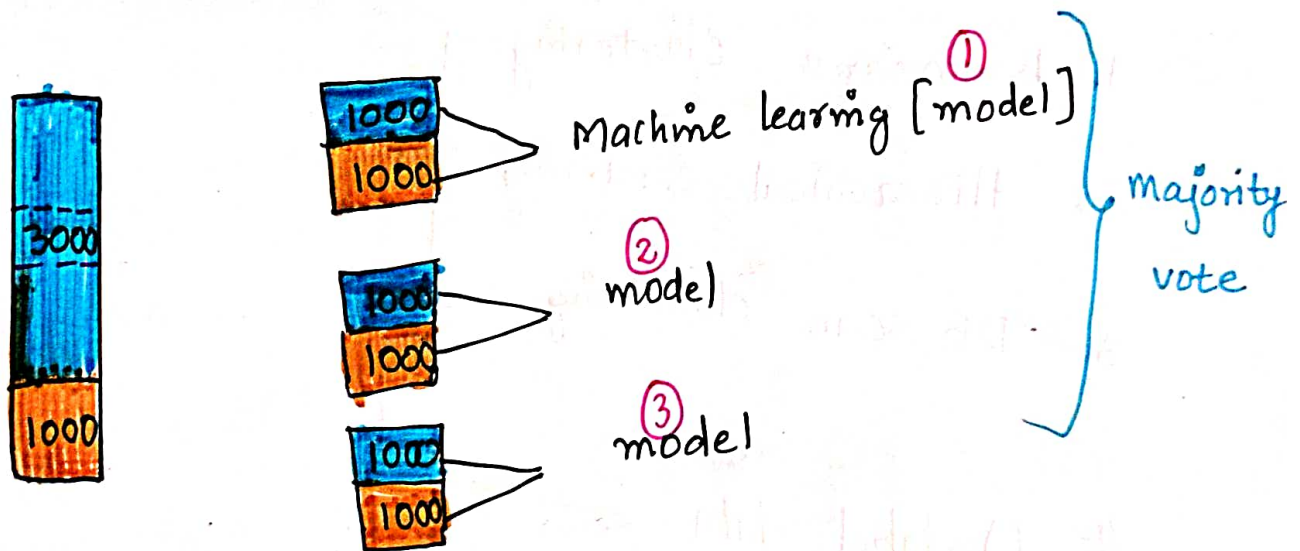
28

Avg: 23.4

If we want to create new data = 23.4
(Synthetic data)

4.

Ensemble Method



Ex:- 3100 records
1000 records } \rightarrow 1100 " 1000 " 1000 records $\pm 5\%$

- * 3000, divided into three parts.
 - * Then will make model, by Each part 1, 2, 3
 - * And we Take Majority vote.
 - * Applying of cross validation.
- divided by
"minority class"

30/4/22 4:00 PM