

Imbalanced data

Supervised Learning

Regression

y is continuous

Classification

y is discrete

$y \{ \begin{array}{l} \rightarrow \text{Pass/Fail} \\ \rightarrow \text{Diabetic/Not diabetic} \end{array}$

* 2 outcomes of target variable (y)
→ binary classification problem
(bi)
↳ two:

Sugar level (f_1)	Cholesterol (f_2)	y (diabetic/Not)
250	102	1 → diabetic
300	100	0 → non diabetic
400	200	0
—	—	—

90% — 0 (non-dia)
10% — 1 (diabetic)

* Assuming you don't build any model → you predict every thing
as class 0

⇒ 90% of accuracy =

* Class imbalance

↓

im + balanced class.

→ when one class has very high percentage as compared to other class, this is class imbalance.

90% class 0 \rightarrow majority class (non diabetic)
 10% class 1 \rightarrow minority class (diabetic)

f_1	f_2	f_3	f_4	y	y_{pred}
\rightarrow				0 \rightarrow non dia	0
\rightarrow				1 \rightarrow diabetic	0
\rightarrow				0	0
\rightarrow				0	0
				0	0

Overall — 90% accuracy

* you don't want your model to learn only the data pattern which is in majority \Rightarrow you need to deal with class imbalance.

Class Imbalance :- When the difference in the counts/Percentage of both the class is huge.

Class Imbalance || 80% \rightarrow Class 0
 20% \rightarrow Class 1

90% - 10%

95% - 5%

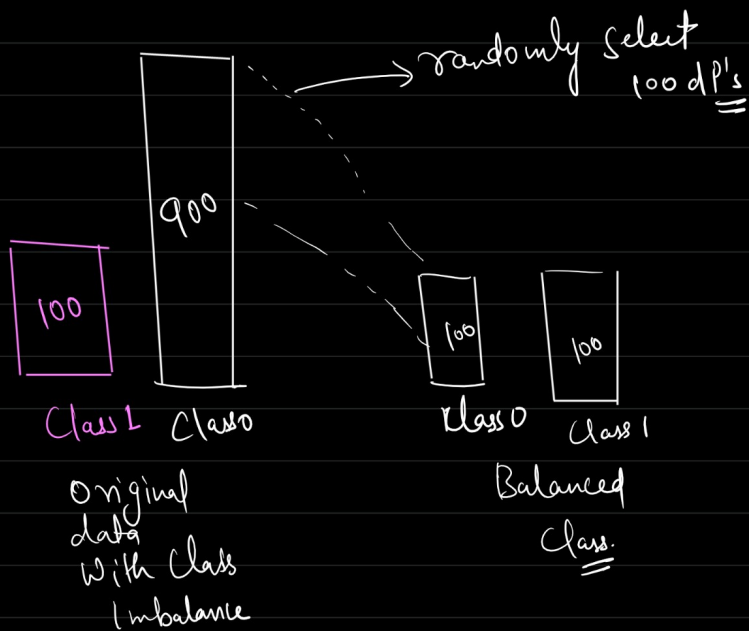
99% - 1%

* Solution to class imbalance

- ① Undersampling (down sampling)
- ② Oversampling (Upsampling)
- ③ SMOTE

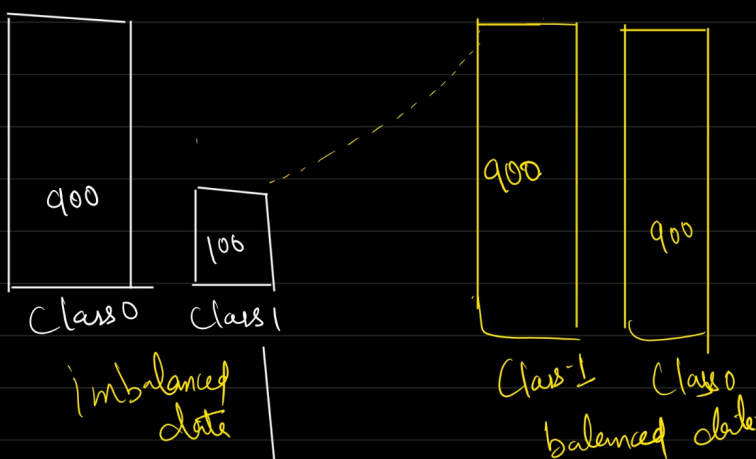
① Undersampling (Down Sample)

disadvantage — Loss of data



② Oversampling / Upsample

→ minority class is upsampled
→ repeat the minority class data.



disadvantage

* ML is about learning pattern in data

→ data is just repeating in case of oversampling

→ No pattern

→ Noise

Class - 100 dp's

$$100 \times 9 = 900.$$

	f_1	b_2	b_3	γ	
→	1	2	1	1	<div style="border-left: 1px solid black; border-bottom: 1px solid black; padding: 10px; display: inline-block;"> $(1, 2, 1)$ • </div>
→	1	2	1	1	
→	1	2	1	1	

② SMOTE (Synthetic Minority Oversampling Technique)

* Problem of No pattern is solved as seen in oversampling

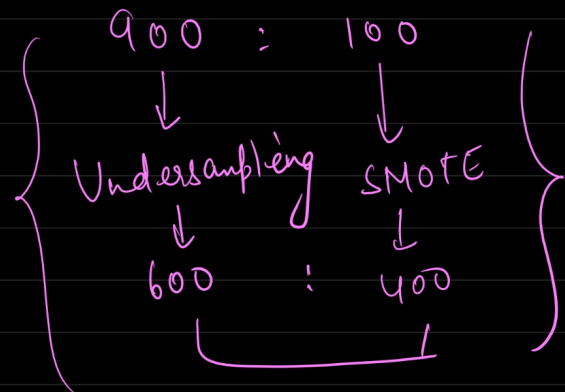


x → Class 0
o → Class 1

→ Like

datapoints (class 0) will always be together, this is the reason smote is effective.

How to use in industry



* In practical implementation of class imbalance.

Population.

1	2
3	4
5	

With replacement

$S_1 (1, 2)$ (data point can be repeated)
 $S_2 (1, 3)$
 $S_3 (2, 3)$
 $S_4 (4, 2)$

without replacement

$S_1 (1, 2)$
 $S_2 (3, 4)$
 $S_3 (5)$
 data will not be repeated.