# imbalanced data

Supervised learning

Classification    Regression

→ Pass / fail
→ Diabetes / Non diabetic

$1/0$
$1/0$

2 outcomes of target variable $(y)$
↳ binary classification

| Sugar level $(f_1)$ | Cholestrol $(f_2)$ | $y$ (diabetic / not) |
|---|---|---|
| 250 | 102 | 1 |
| 300 | 100 | 1 |
| 400 | 200 | 0 |
| — | — | — |

diabetic / Non diabetic

{ 90% → Non diabetic
  10% → diabetic

imbalanced class          class $(y)$

Im + balanced        Non diabetic    diabetic

Assuming diabetic is 1 (Class1)

non diabetic → 0 (Class 0)

non diabetic — 0 (class 0)

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

90% Class 0 ⟶ majority class (non-diabetic)

10% Class 1 ⟶ minority class (diabetic)

$f_1$   $f_2$   $f_3$   $f_4$   | $y$                      | $y_{pred}$

⟶                                    0  ⟶ non diabetic         0
⟶                                    0                          ⓪
⟶⟶                                  1  ⟶ diabetic              0
⟶⟶                                  0                          0
                                     0                          0
                                     0

⟶ — — — — —   ?         Overall ⟹ 90% accuracy

Class imbalance
├⟶ Undersampling
├⟶ Oversampling
└⟶ Smote

\* class imbalance ⟶ when the difference in the counts/percentage of the both the class is huge.

Class imbalance { 80% ⟶ Class 0 (non diabetic)
                 20% ⟶ Class 1 (diabetic)
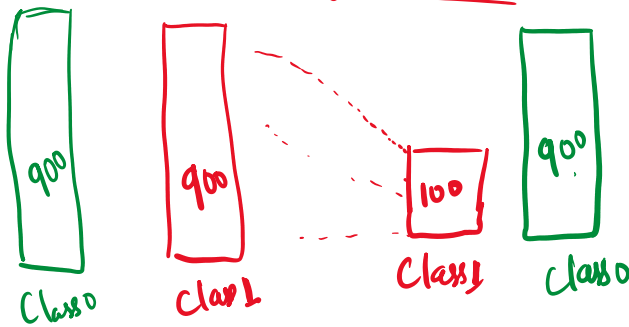
90 – 10%
95 – 5%
99 – 1%

99 - 1 %

* <u>Undersampling</u> (Downsample)

disadvantage
→ loss of
data

```
        ┌──┐
        │900│
┌──┐    │  │       ┌──┐  ┌──┐
│100│   │  │       │100│ │100│
└──┘    └──┘       └──┘  └──┘
Class 1  Class 0    class 0  Class 1
```

Original dataset
(Class imbalance)

* <u>Oversampling | Upsample</u>

```
┌──┐   ┌──┐           ┌──┐
│900│  │900│    ┌──┐   │900│
│  │   │  │     │100│  │  │
└──┘   └──┘     └──┘   └──┘
Class 0  Class 1  Class 1  Class 0
```

minority class is
                Upsampled

→ repeat the minority
   class data.

disadvantage
→ No pattern
→ Noise

Class 1 → 100 d $\beta$s
→ $100 \times 9 = \underline{\underline{900}}$

$f_1 \quad f_2 \quad f_3$

$$
\begin{array}{ccc}
- & - & 0 \\
1 & 2 & 1 \\
- & - & 0 \\
- & - & 0 \\
1 & 2 & 1 \\
1 & 2 & 1 \\
1 & 2 & 1
\end{array}
$$

(1,2) .

# * SMOTE

## Synthetic Minority Oversampling technique



$$4 \qquad 36$$
$$Cl-1 \qquad Cl-0 (X)$$
$$(\cdot)$$

## How used in industry

$$900 : 100$$

Undersampling → SMOTE

$$600 : 400$$

★

## with replacement

$$\begin{array}{ccc} 1 & 0^2 0 & 5_0 \\ 3 & 0 & 4_0 \end{array}$$

| with replacd | without replacd |
|---|---|
| $S_1 (1,2)$ | $S_1 = \{1,2\}$ |
| $S_2 (2,3)$ | $S_2 = \{3,4\}$ |
| $S_3 (1,3)$ | $S_3 = \{5\}$ |
| $S_4 (4,2)$ | |