

DT-1

Imbalanced data

Agenda

* Imbalanced Dataset

- Under Sampling
- Over Sampling
- SMOTE

* DT Intro

* Geometric Intuition

* Algorithm (mathematical)

* PURE NODE \rightarrow Entropy \rightarrow Information Gain

Imbalanced dataset

Binary Classification. $\rightarrow 1000$

One class \rightarrow dominated.

50% - 50% \rightarrow Balanced

60% - 40% \rightarrow Slightly balanced

70% - 30% \rightarrow Slightly Imbalanced

80% - 20% \rightarrow Imbalanced

90% - 10%
95% - 5% \rightarrow Extremely Imbalanced

950 \downarrow majority class
50 \downarrow minority class

① Classification - Model

Dumbest Model \rightarrow Normal
every trans is \rightarrow

Accuracy \Rightarrow 99%

② Model

1% - F

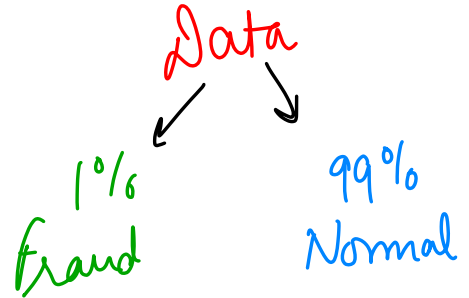
99% - NF

\hookrightarrow Tends to become biased
towards Majority Class.

Fraudulent Transaction

HDFC

10000 $\xrightarrow{1\%}$ 100 Fraud.
99% 9900 Normal

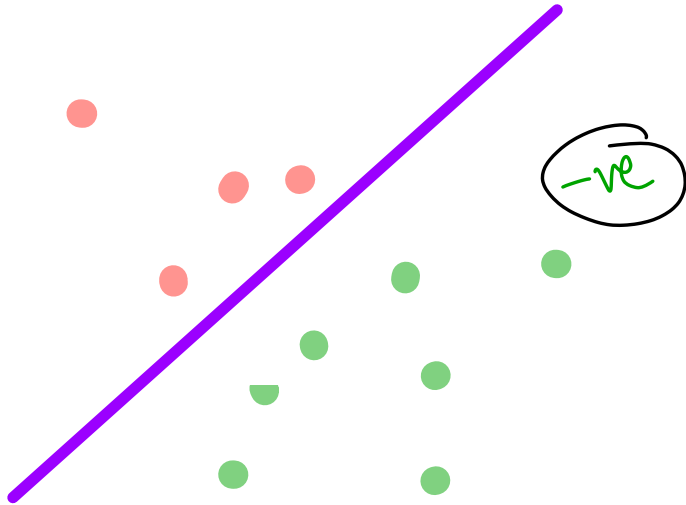


Impact of Imb data on log. Reg.

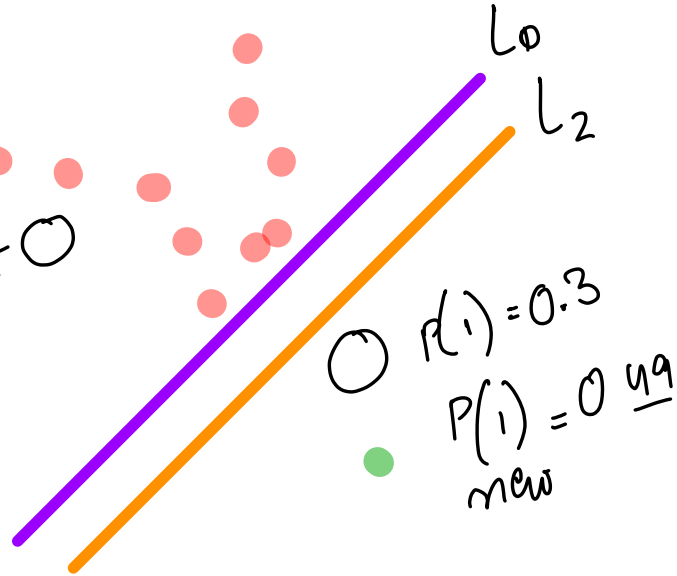
$$\log \text{loss} = \sum \log \text{loss}$$

$$y \log(\hat{y}) + (1-y) \log(1-\hat{y})$$

pos

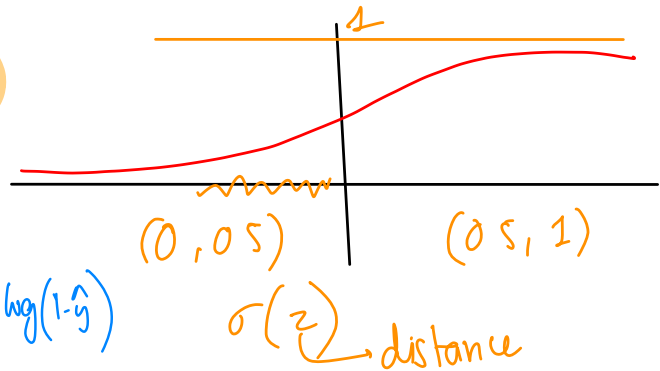


$$P(i) = 0.7 \bigcirc$$



$$P(i) = 0.3$$

$$P(i) = 0.49_{\text{new}}$$



Impact of Imbalanced data on KNN

$k=1,3$

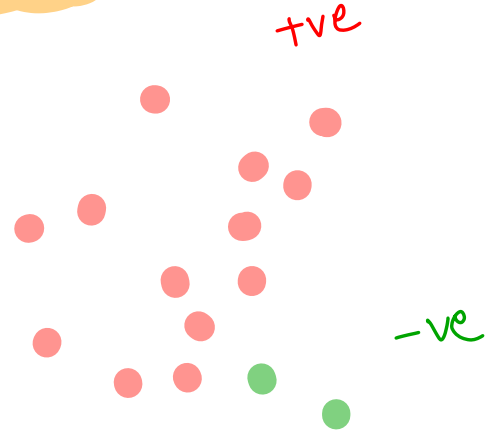
low

Impact of Imb. data
for low values of k
is comparatively low

$k=100$

high

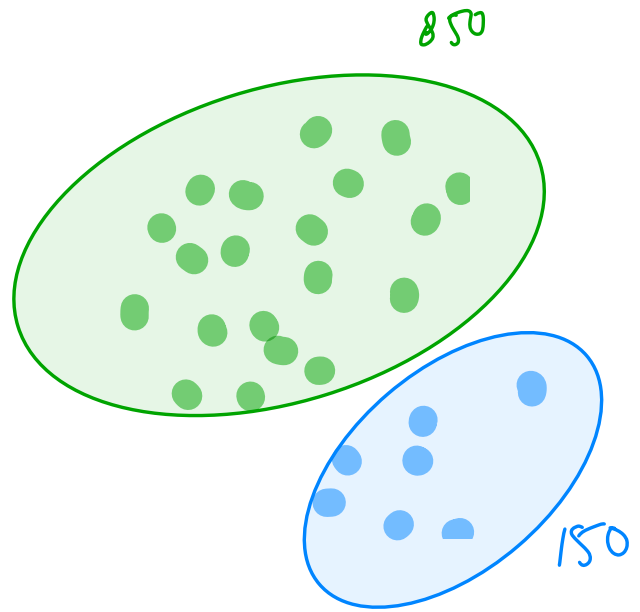
Impact of Imb. data
for high values of k
is HIGH.



Handle Imb. dataset.

① Class weights

1600 emails (Imbalanced)
850 Normal (majority)
150 Spam (minority)



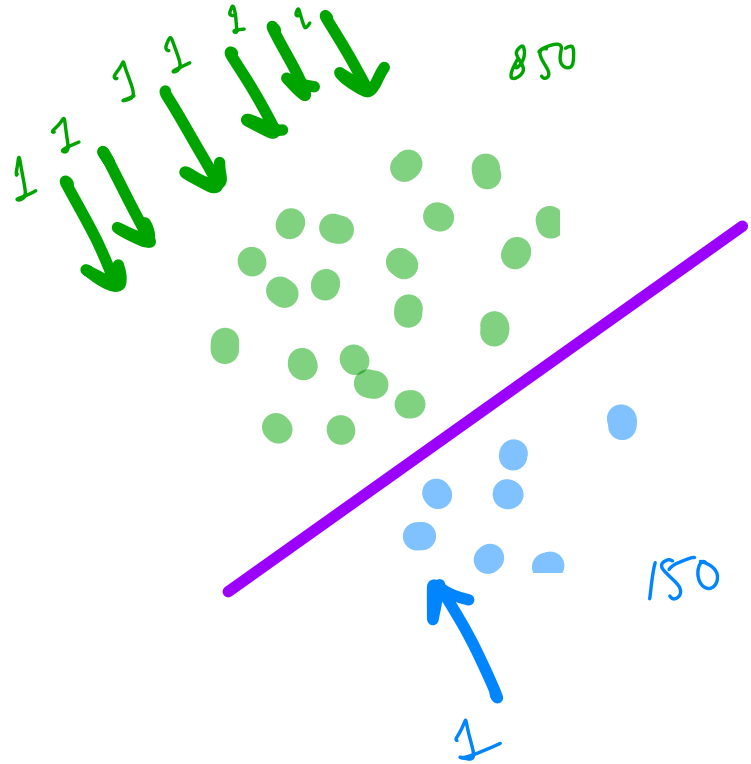
$$\frac{850 \text{ NS}}{150 \text{ S}} \approx 5.66 \approx 6$$

For every 1 Spam email, I have 6 NS emails.

$$w_i^0 \begin{matrix} (+ve) \\ (-ve) \end{matrix} = \begin{matrix} 1 \\ 6 \end{matrix} \quad \left| \quad y_i^0 = \begin{matrix} 1 (+ve) \\ 0 (-ve) \end{matrix}$$

$$\sum_{i=1}^n \log \text{loss} \cdot (w_i^0)$$

$$\{ \textcircled{0} \cdot 6, \textcircled{1} \cdot 1 \}$$



②

Over sampling

Advantage:

* No data loss

* Preferred over under sampling

Disadvantage

* Possibility of overfitting increases.



$$\begin{aligned} 850 + 850 \\ = 1700 \\ \uparrow \\ \text{Modelling} \end{aligned}$$

Imbalanced
↓

Balanced.

→ from where we will get this?

Random Sampling with Repetition
Duplicating minority class Randomly.

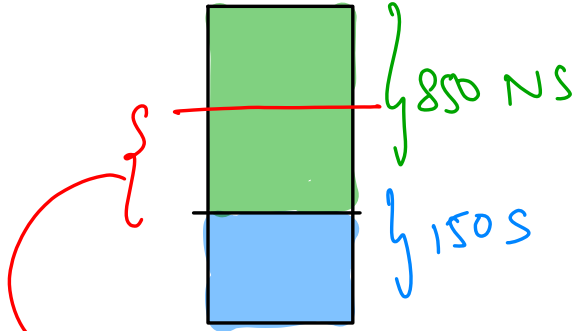
③ Under Sampling

Advantages:-

* Time & Space effective.

Disadvantages:-

- * loss of data {expensive}
- * lose out on potentially important data points.
- * Sample chosen might be biased.

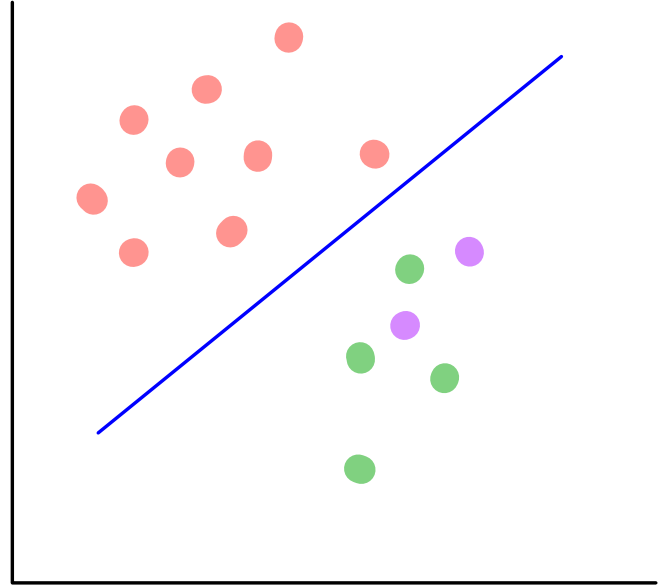
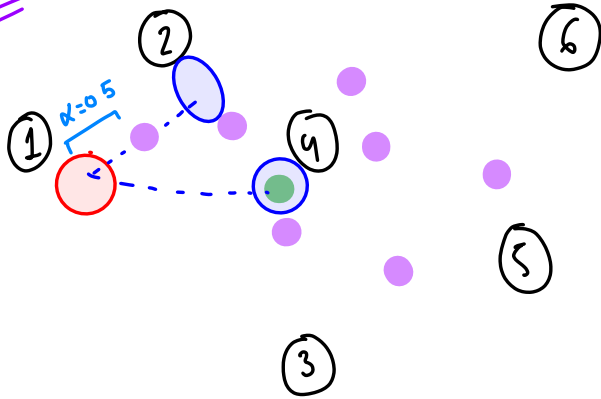


$$850 + 150 = 1000$$
$$\downarrow$$
$$150 + 150 = \boxed{300}$$

Randomly select equal
no. of data point
equal to minority
class.

④ SMOTE [Synthetic Minority Over sampling Technique]

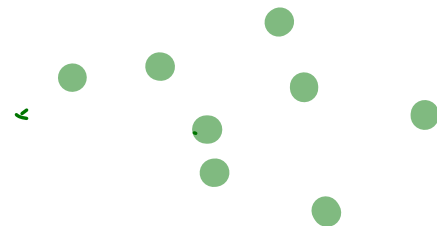
KNN



① $K = 2$

② Pick any value b/w $(0,1) = \alpha$

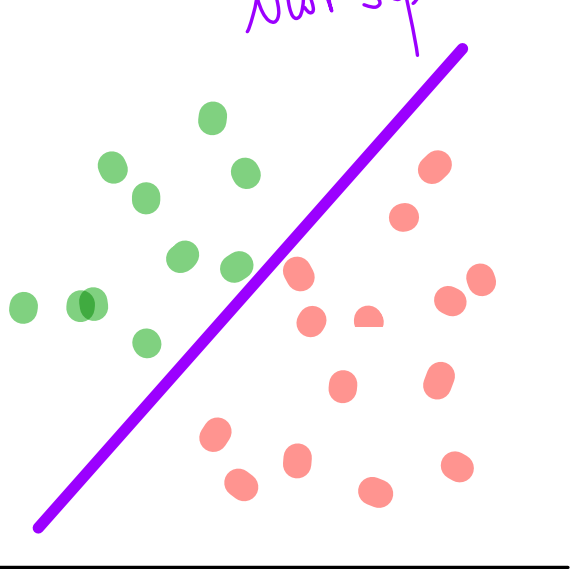
③ New point $x_{\text{New}} = x_1 + \alpha d(1,2)$



DOUBTS

F_1

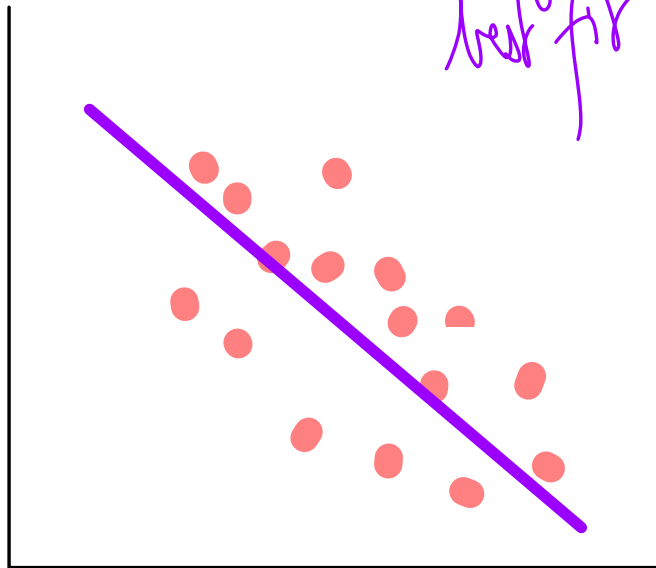
line of
best separation



F_2

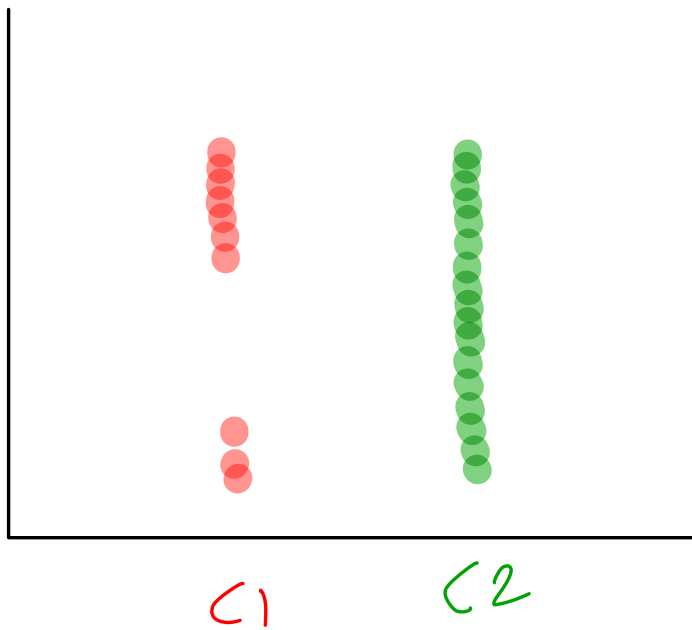
Y

line of
best fit



F

Cont



odds $\frac{1}{4:1}$ $\frac{0}{-ve}$

$$P(Y=1|X) = \frac{4}{5}$$

$$P(Y=0|X) = \frac{1}{5}$$

$$\frac{P(\text{Success})}{P(\text{Failure})} = \frac{P}{1-P}$$

σ

$$\log\left(\frac{P}{1-P}\right) =$$

$$w^T x + w_0$$

$$\log\left(\frac{P}{1-P}\right)$$

$$= w_1 x_1 + w_2 x_2 \dots w_n x_n + w_0$$