

- Introduction to Anomaly / Novelty / Outlier Detection
- Distribution Based Anomaly Detection
 - RANSAC (Random Sample Consensus)
 - Elliptic Envelope
- Sklearn Implementation
- Isolation Forest
- Disadvantages of Isolation Forest
- Sklearn Implementation


What is an Anomaly?

- Something that is not Normal
- Novel (new / never seen before)
- Unique
- Outliers

➤ IQR \Rightarrow mean and σ

- 1 Dimension at a time

mean and σ



outlier points

Distribution Based Outlier Detection

assumption: Data follow Gaussian
Distribution



$$N(\bar{x}, \sigma)$$

↓
Prob-score

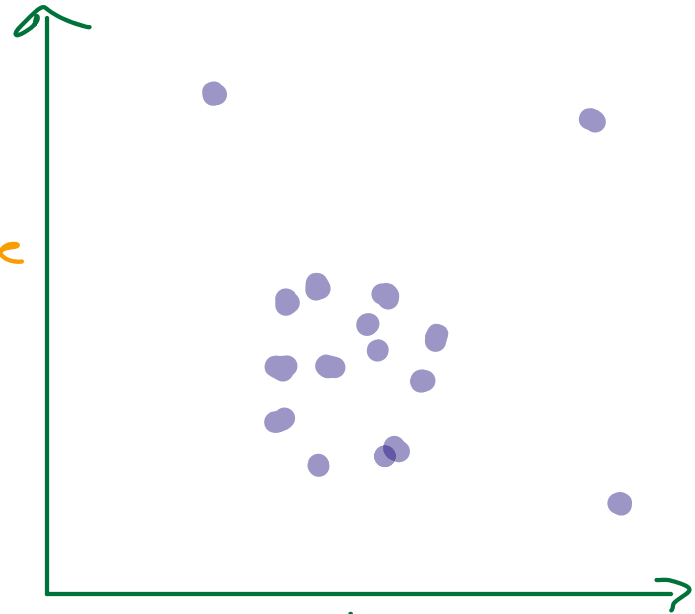
$\text{prob} < 0.05 \rightarrow \text{outlier}$

RANSAC (Random Sample Consensus)

1) Take Random Sample from Dataset

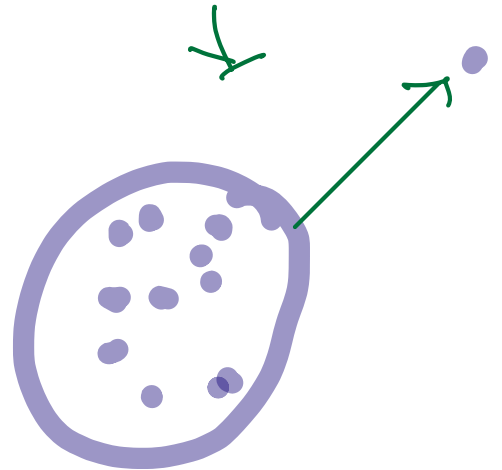
2) Calculate σ and mean of that Sample

3) Repeat the above two steps N Times



$[G_1, G_2, G_3, \dots, G_n]$

$[\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_n]$



mean of means

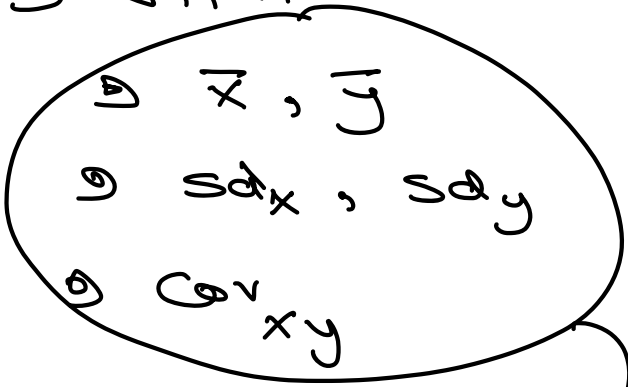
median of means

4) Robust Estimation of Distribution Parameters

5) By Doing Random Sampling Consensus we can reduce impact of outliers

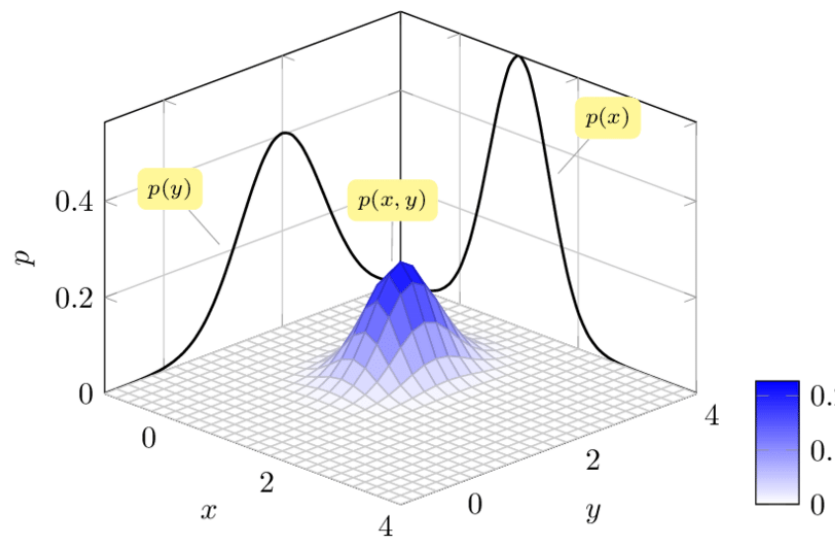
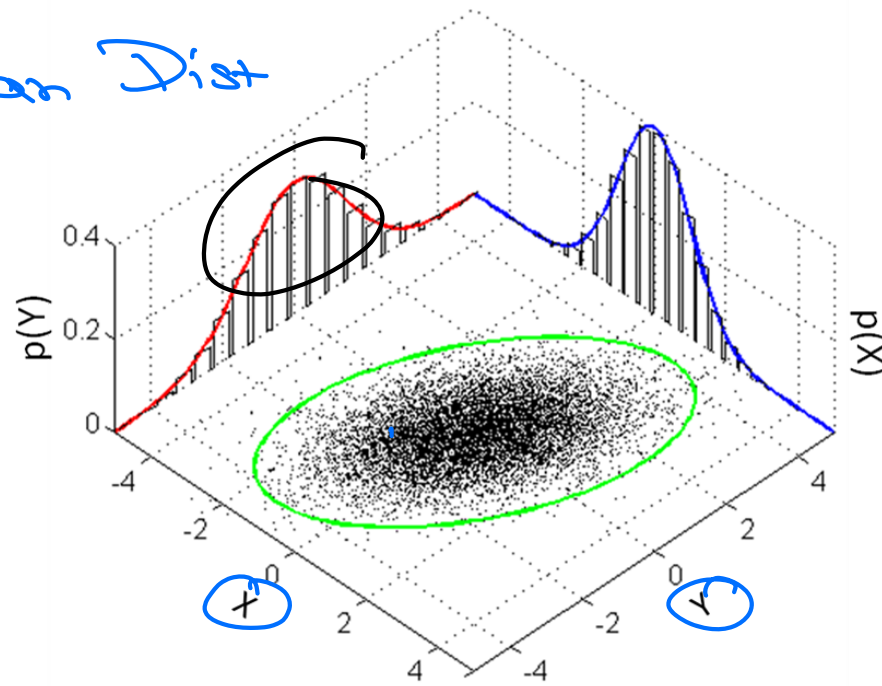
2d Gaussian Dist

① GMM



$$X \Rightarrow \begin{bmatrix} \bar{x} \\ \bar{y} \end{bmatrix} \Rightarrow$$

$$\text{Cov-mat} \Rightarrow \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_y^2 \end{bmatrix}$$



Elliptic Envelope

Estimating Gaussian Dist' parameters in multi-Dimension

$$\begin{matrix} 1 \times 3 \\ 3 \times 1 \end{matrix} \cdot \begin{matrix} \bar{x} \\ \bar{y} \\ \bar{z} \end{matrix}$$

$$\begin{bmatrix} \sigma_x^2 & \sigma_{xy} & \sigma_{xz} \\ \sigma_{yx} & \sigma_y^2 & \sigma_{yz} \\ \sigma_{zx} & \sigma_{zy} & \sigma_z^2 \end{bmatrix}$$

Elliptic Envelope \Rightarrow Estimate Params for Multi-Dim

① \Rightarrow Estimate Parameter



② \Rightarrow Calculate prob of every Single point in Dataset

③ How Do find threshold beyond which outliers are assumed?

$N=100$

$P \Rightarrow P_i$

P_{100}



Contamination Factor

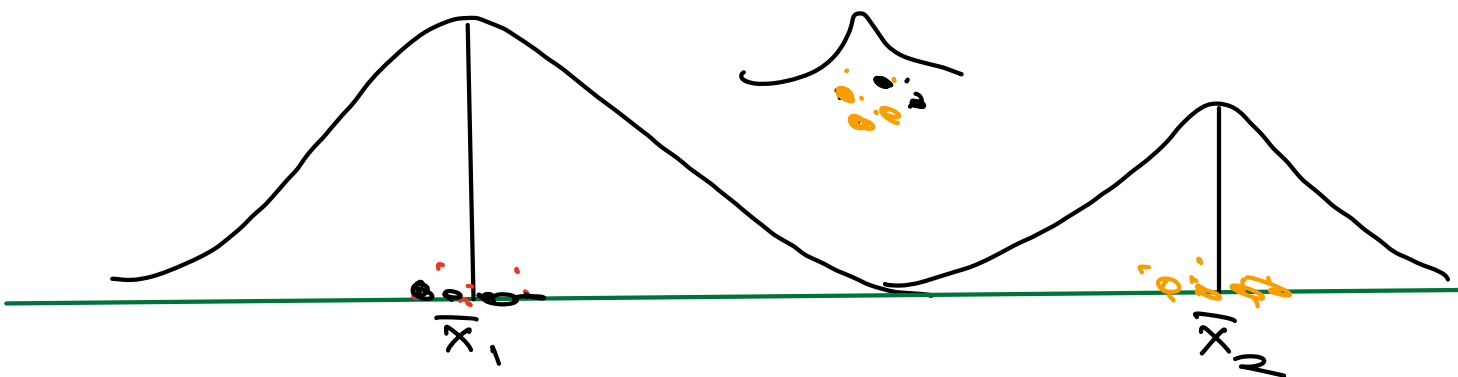
\downarrow
 $(0.05] \rightarrow 0.10$

Sort all point w.r.t prob-score and pick the lowest 10%.

- ① 1000 Data-points 1D.
- ② $G, \bar{x} \leftarrow \text{RANSAC}$
- ③ assign prob score to each point using PDF
- ④ C.F. ① 0.005 $\frac{5}{1000} \times 100 \%$
- ⑤ Filter 50 point with lowest prob-score
- ⑥ $\text{Sort}(\text{prob-score}, \text{point})[0:50]$
 \downarrow

⑥ Disadvantages

- ① Strict Assumption
Gaussian
- ② Doesn't work with multi-modal



Isolation Forest



① Simple and Elegant

② Works really well for Most use-cases

Forest :

Build Many Forests of multiple Trees

To Build Each tree in forest

Step 1 :

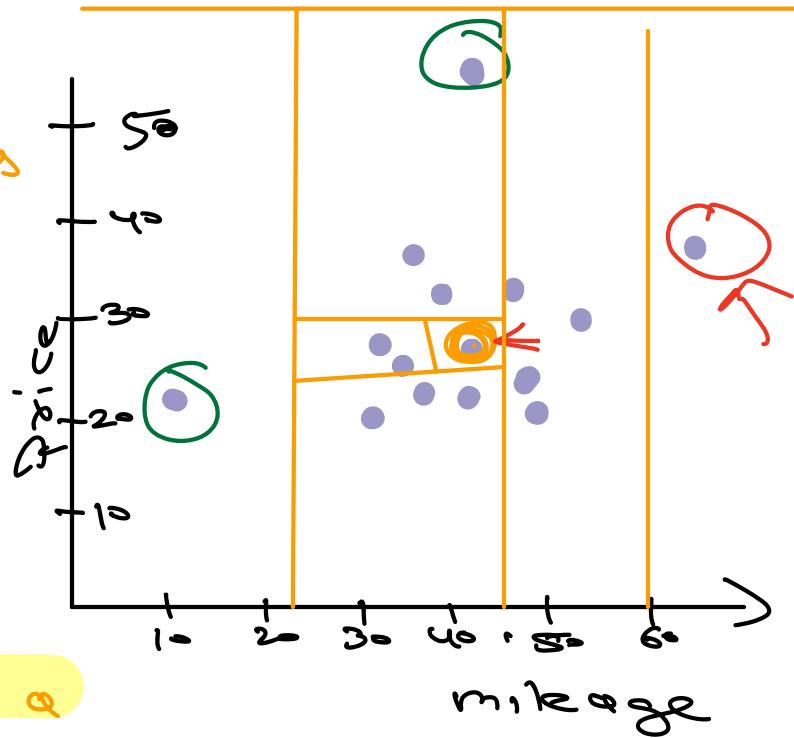
Randomly pick a

feature

Step 2 :

Randomly pick a threshold

③ Repeat Step 1 and Step 2 until Every single Data-point is isolated



③ Outliers points are isolated earlier than inlier points

④ Calculate Splits required to isolate every point across all trees

T - tree

$x_1 \Rightarrow [2, 1, 3, 2, 1, 4, 1]$

$x_2 \Rightarrow [5, 9, 8, 6, 3, 2, 1]$

$x_3 \Rightarrow$

x_n

|

|

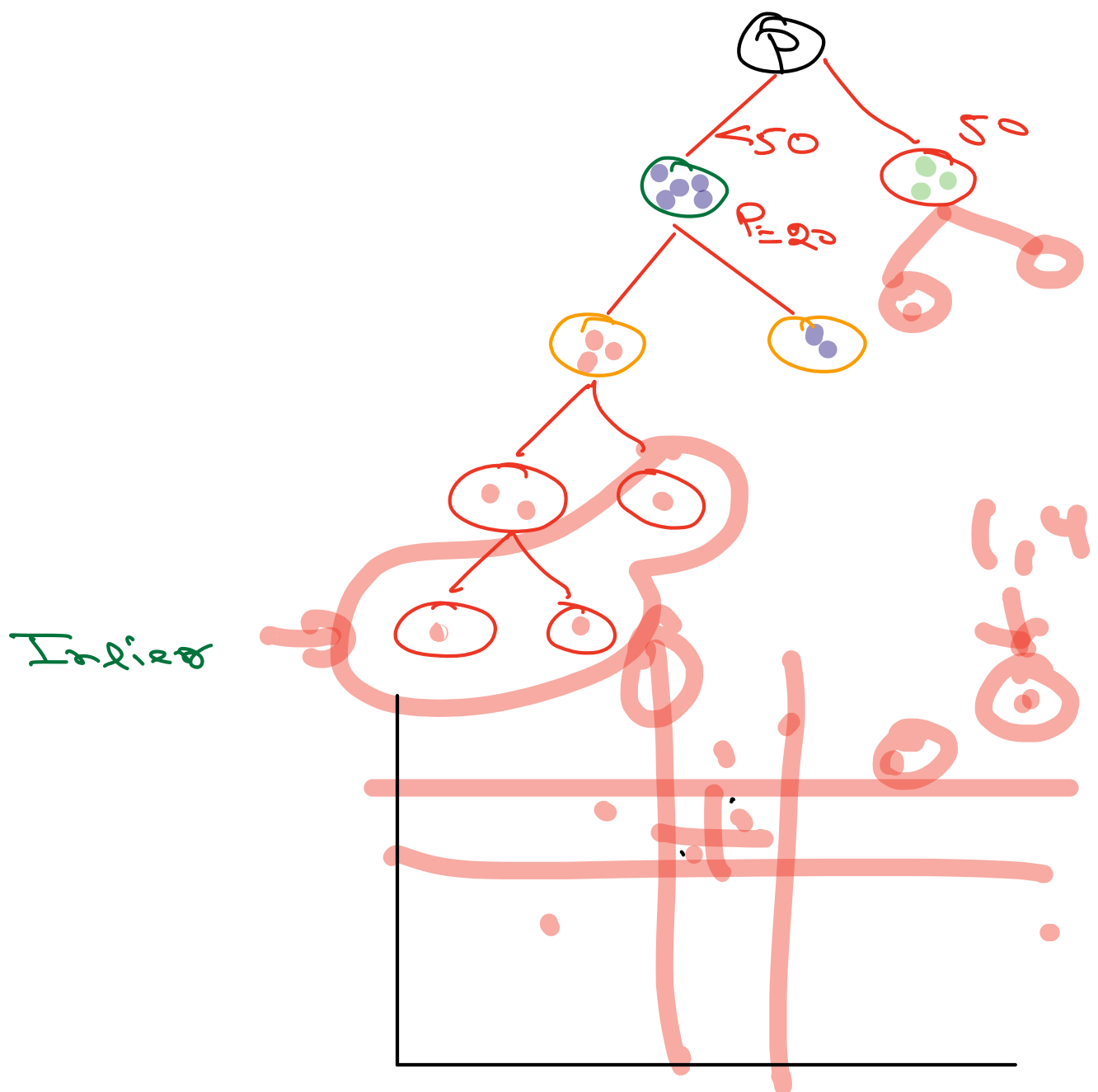
x_n

avg number of Splits across all trees

⑤ Define threshold to filter Data-point with less Splits

Contamination-Factor

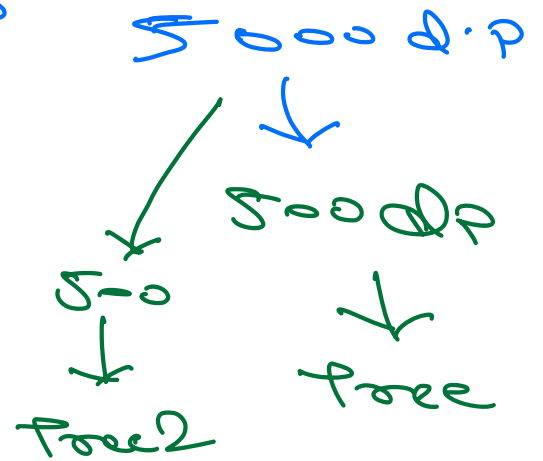
↓
Sort Avg Splits
and
Select min n%



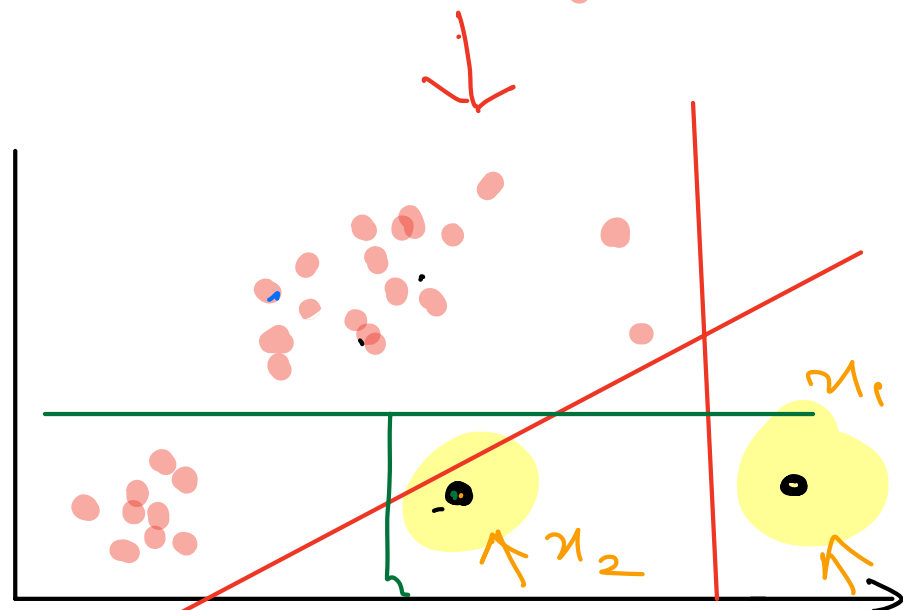
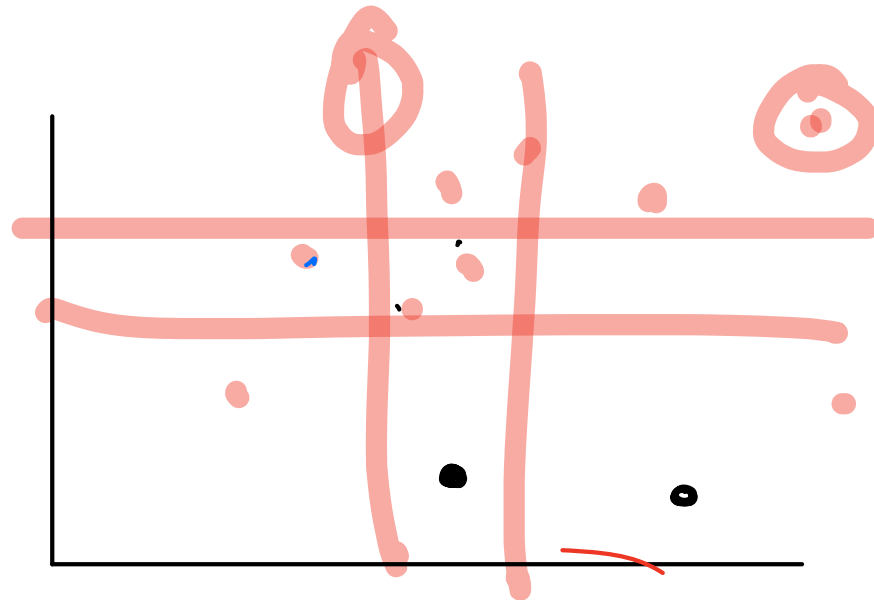
Disadvantage

① Isolation of Every D.P
is time consuming

Remedy: n° Sample $\ll n$



② Biased due to axis
parallel splits



✗ not possible with
iForest

③ One-Class Classifiers
SVM

④ LOF