



SNS COLLEGE OF TECHNOLOGY

Coimbatore-35
An Autonomous Institution

Accredited by NBA – AICTE and Accredited by NAAC – UGC with 'A++' Grade
Approved by AICTE, New Delhi & Affiliated to Anna University, Chennai



Build an Entrepreneurial Mindset Through Our Design Thinking FrameWork

DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING

23AMB201 - MACHINE LEARNING

II YEAR IV SEM

UNIT I – INTRODUCTION

TOPIC 4– Statistics for Machine Learning(Cont..)



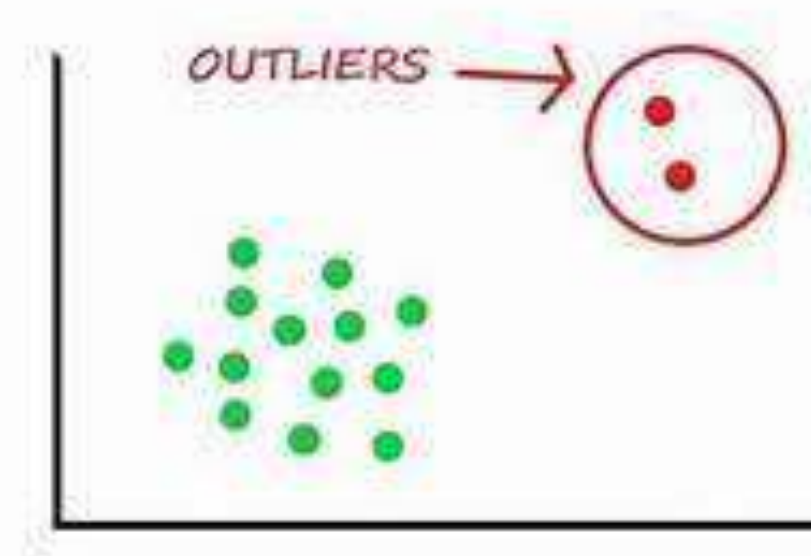
Outlier, Noise and Error



Noise: Irrelevant dataset.

| Student ID | Student Name | Age | GPA | Classification |
|------------|--------------|-----|-----|----------------|
| 100122014 | Joseph | 21 | 3.5 | Junior |
| 100232015 | Patrick | 200 | 3.2 | Sophomore |
| 100122012 | Seller | 24 | 3.0 | Senior |
| 100342013 | Roger | 23 | 234 | Senior |
| 100942012 | Davis | 2.8 | 3.7 | Sophomore |
| | Travis | 23 | 3.4 | Sr |
| 100982015 | Alex | 27 | | Sophomore |
| 100982013 | Trevor | -22 | 4.0 | Senior |
| AUC2016XC | Aman | 30 | 3.5 | Jr |

Outlier: relevant one, but so far from actual data



| Bill Amount | Tips |
|-------------|------|
| 500 | 25 |
| 500 | 35 |
| 600 | 40 |
| 300 | 20 |
| 500 | 20 |
| 500 | 140 |

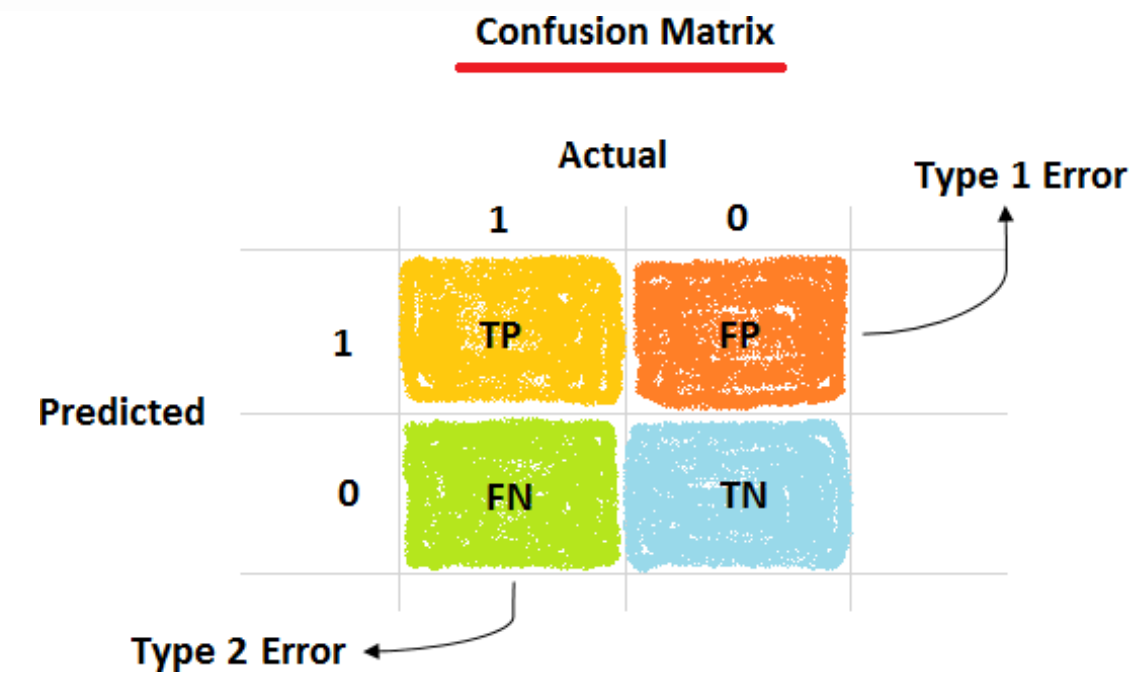
Missing Data

Inconsistent Data

Noisy Data

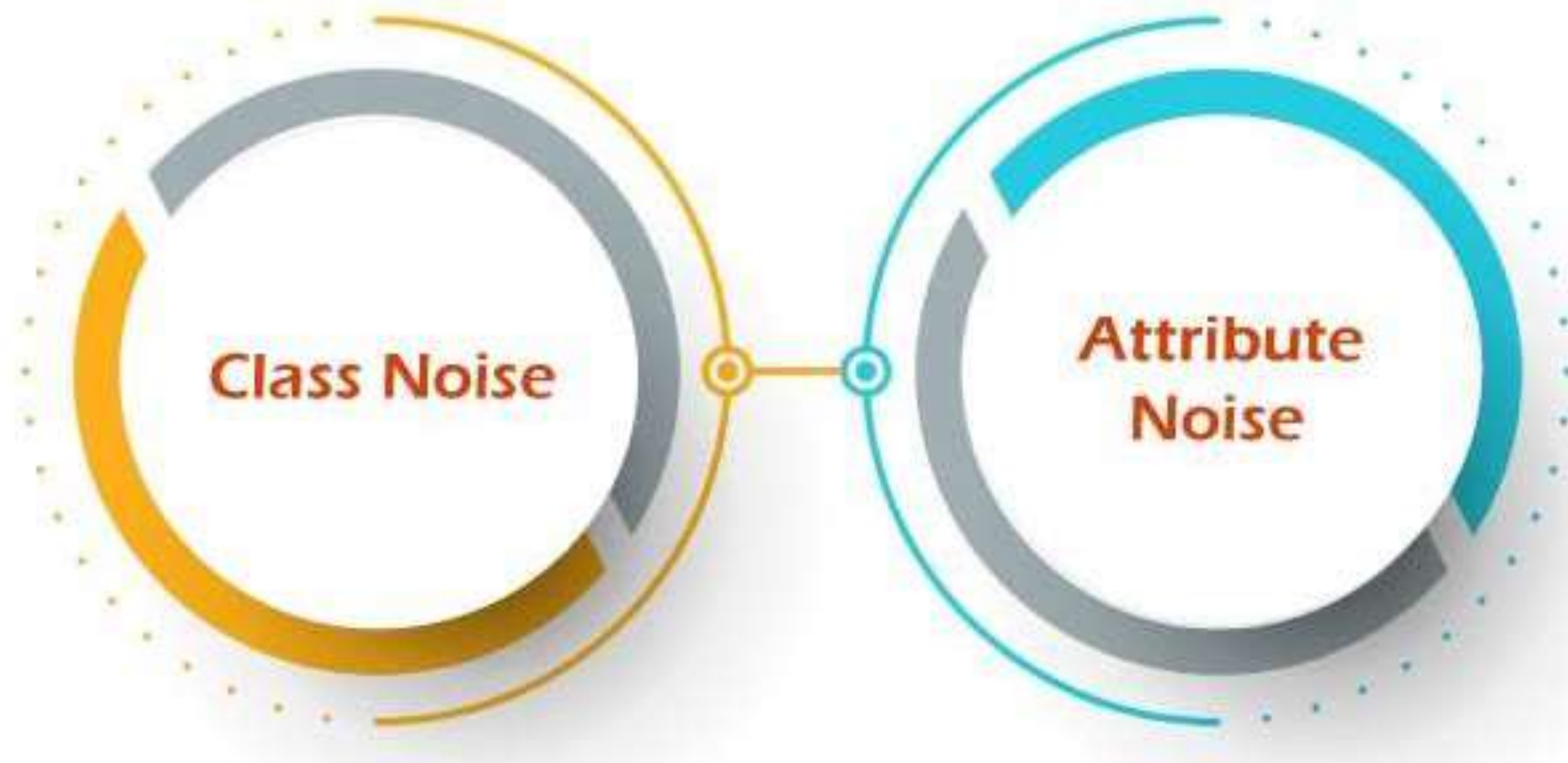
| | Name | Age | Gender | Marks |
|---|---------|-----|--------|-------|
| 0 | Jai | 17 | M | 90 |
| 1 | Princi | 17 | F | 76 |
| 2 | Gaurav | 18 | M | NaN |
| 3 | Anuj | 17 | M | 74 |
| 4 | Ravi | 18 | M | 65 |
| 5 | Natasha | 17 | F | NaN |
| 6 | Riya | 17 | F | 71 |

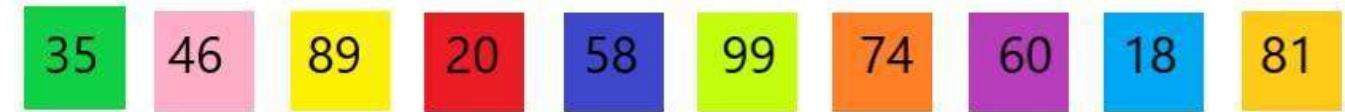
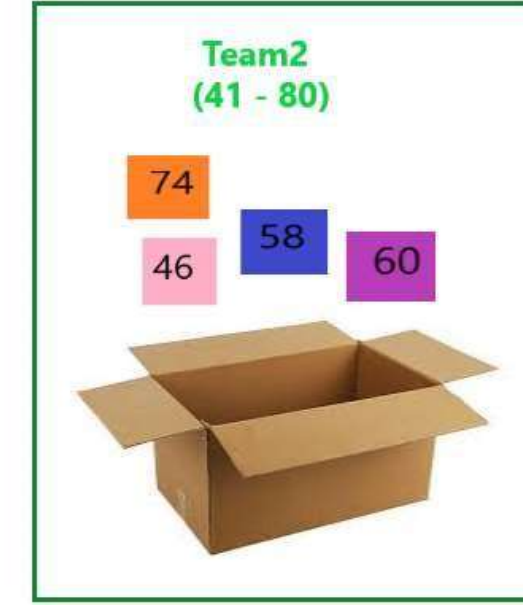
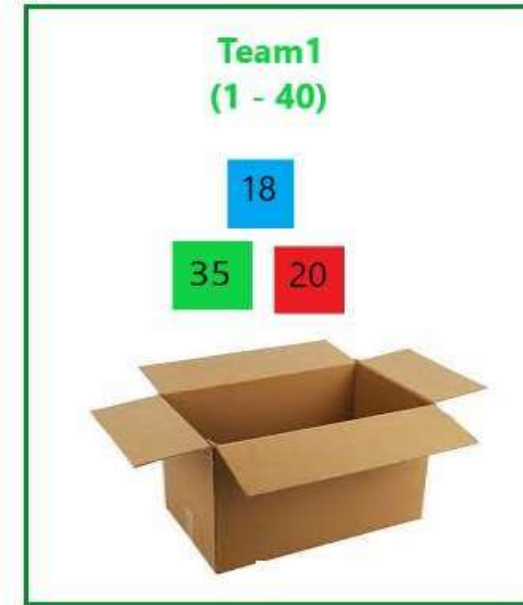
Error:



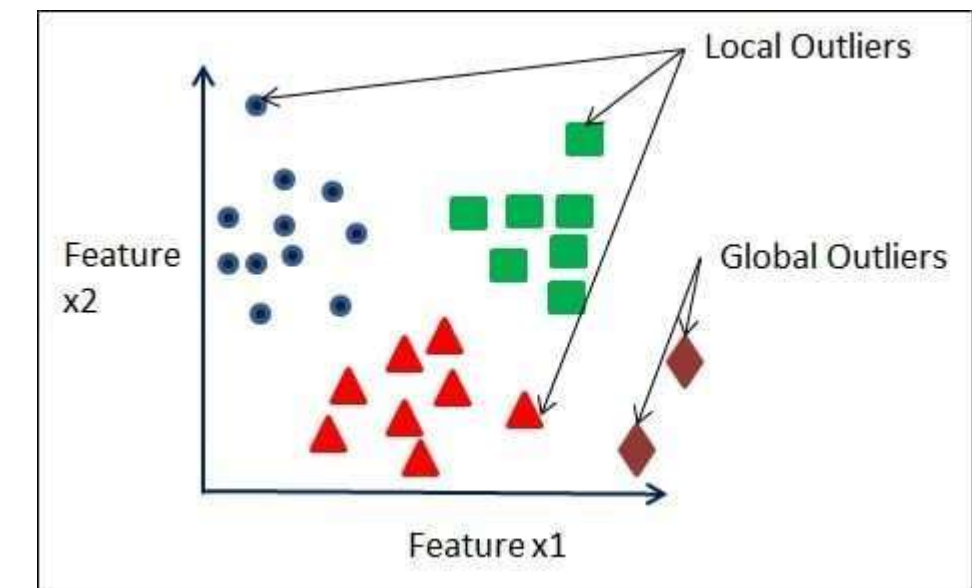
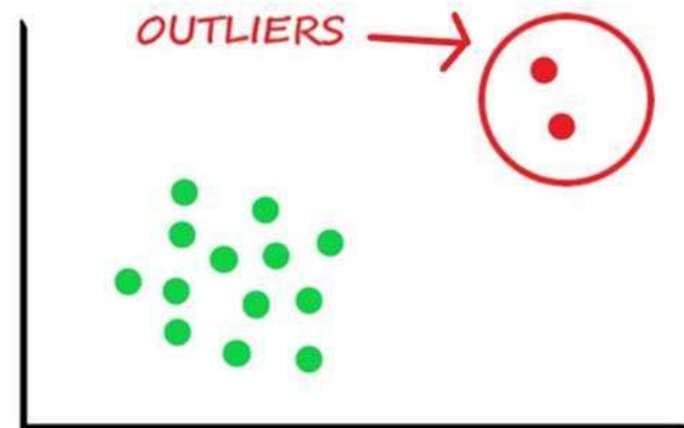
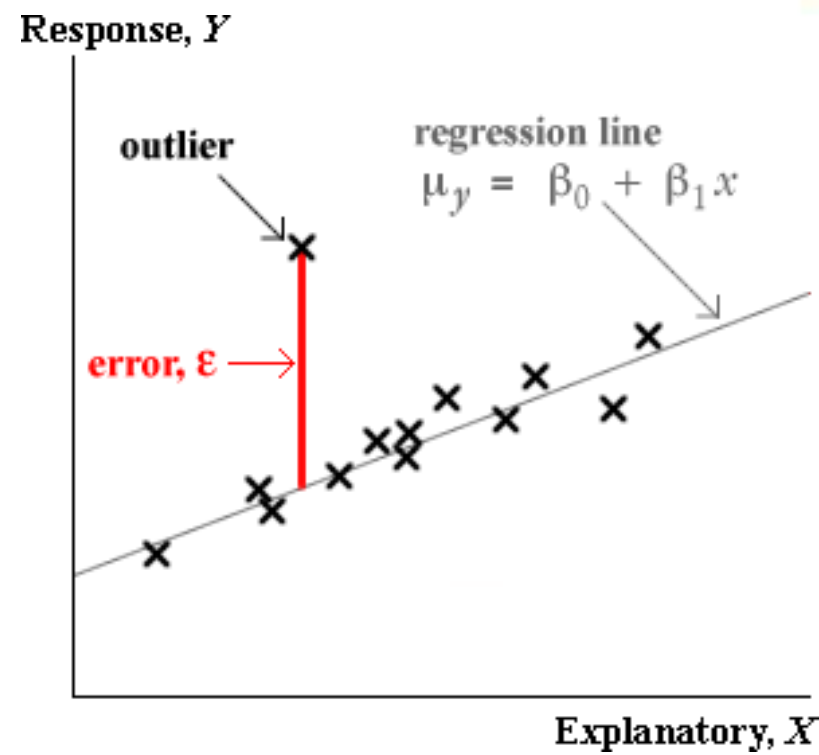


Types of Noise





Data





Data Wrangling



```
import pandas as pd

# Assign data
data = {'Name': ['Jai', 'Princi', 'Gaurav',
                'Anuj', 'Ravi', 'Natasha', 'Riya'],
        'Age': [17, 17, 18, 17, 18, 17, 17],
        'Gender': ['M', 'F', 'M', 'M', 'M', 'F', 'F'],
        'Marks': [90, 76, 'NaN', 74, 65, 'NaN', 71]}

# Convert into DataFrame
df = pd.DataFrame(data)

# Display data
df
```

| | Name | Age | Gender | Marks |
|---|---------|-----|--------|-------|
| 0 | Jai | 17 | M | 90 |
| 1 | Princi | 17 | F | 76 |
| 2 | Gaurav | 18 | M | NaN |
| 3 | Anuj | 17 | M | 74 |
| 4 | Ravi | 18 | M | 65 |
| 5 | Natasha | 17 | F | NaN |
| 6 | Riya | 17 | F | 71 |



Data Wrangling(Mean)



| | Name | Age | Gender | Marks |
|---|---------|-----|--------|-------|
| 0 | Jai | 17 | M | 90 |
| 1 | Princi | 17 | F | 76 |
| 2 | Gaurav | 18 | M | NaN |
| 3 | Anuj | 17 | M | 74 |
| 4 | Ravi | 18 | M | 65 |
| 5 | Natasha | 17 | F | NaN |
| 6 | Riya | 17 | F | 71 |

```
c = avg = 0
for ele in df['Marks']:
    if str(ele).isnumeric():
        c += 1
        avg += ele
avg /= c

# Replace missing values
df = df.replace(to_replace="NaN",
               value=avg)
```

```
# Display data
df
```

| | Name | Age | Gender | Marks |
|---|---------|-----|--------|-------|
| 0 | Jai | 17 | M | 90.0 |
| 1 | Princi | 17 | F | 76.0 |
| 2 | Gaurav | 18 | M | 75.2 |
| 3 | Anuj | 17 | M | 74.0 |
| 4 | Ravi | 18 | M | 65.0 |
| 5 | Natasha | 17 | F | 75.2 |
| 6 | Riya | 17 | F | 71.0 |



References

1. Aurélien Géron "Hands-On Machine Learning with Scikit-Learn and TensorFlow" Publisher(s): O'Reilly Media, Inc 2017.

Thank You