

## \* Outliers detection - Box plot

1. Median 25%
2. 75% population.
3. Inter Quartile
4. IQR

outlier is an observation i.e; unlike the other observation.

Caused by : measurement or i/p error

Data corruption

True outlier/observation.

May cause problem during model fitting.

noise: random behaviour , value scattered randomly in dataset

outlier: data follows trend, but some points do not.

Possible Solutions:

A model that can handle outliers, max margin classifiers SVM.

Seaborn feature.

import seaborn as sns

```
Sns.boxplot ( data=df, x=column, ax=axs [i, 1], color=color)
```

Step 1: Ascending

Step 2: Median

Q1

Q3

IQR

Step 3: Max, Min

## Imputing Missing Values

- ↳ Mean value (of a column)
  - ↳ Median value
  - ↳ Mode
  - ↳ Random Sample imputation

Naive Bayes model, Nearest neighbors.

\* Limited Quantity of Training Data - Problems



Test data: cat standing identifies as dog.  
(Picks up the noise).

\* Data Augmentation  $\rightarrow$  slightly modify the image.  
(perturbations)

original  $\Rightarrow$  learns more about distribution of <sup>data</sup> image.

Augment your dataset by creating additional training examples through rotation, flipping, cropping.

## \* Feature engineering

- ↳ which features are actually contributing to the best model.
- high quality features can help your model learn from the available data more effectively

## \* Reduction of parameters

## \* Alternate techniques

- ↳ regularization
- ↳ semi-supervised learning

## (data cleaning)

# Data Transformation

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

linear transformation;

$$X \alpha = y$$

↙      ↘  
matrix    vector

classification → not easy to find decision boundary

Transforms into a space where it can easily be distinguished.

Linear Transformation:

function that maps i/p vector  $\rightarrow$  o/p vector.

Rotation  $x = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$

Scaling

Shear → change the orientation of vector.

(If we can't find decision boundary in one space, it then makes/transforms into another space where it can be easily distinguished - decision boundary).

Support Vector

Neural networks: Complex datasets

how to transform data in another basis  
that could simplify the data.

Same info, but shown in diff way

we can learn about decision boundary, data properties

→ A line should remain a line after transformation.

→ And origin will be fixed.

→ Dist b/w the grid lines should remain equidistant.

Properties:  $T(\vec{u} + \vec{v}) = T(\vec{u}) + T(\vec{v})$

$$T(c\vec{u}) = cT(\vec{u})$$

$$T(0) = 0$$

## ① Data Normalization

- Some features dominating distances
- Models don't realize that datapoints may or may not be of same units.

Every feature must be treated with equal importance.

Train algo with normalization: Better accuracy  
with no normalization: neglecting most of the features

(may miss out on imp.)

why?

Transform data → similar distribution

→ dimensionless

→ same units

EQUAL IMP TO ALL FEATURES

\* Standardization, Z-score Normalization.

Transforming features in a way s.t

$$\text{Mean } (\mu) = 0$$

$$\text{Std dev } (\sigma) = 1$$

$$x_{\text{std}} = \frac{x_i - \mu}{\sigma}$$

- SVM, K-means clustering  $\rightarrow$  looking at <sup>NN</sup>  $\Rightarrow$  Belong to same cluster  
(usually scale our features)

$\Rightarrow$  If it is not of same units

4 features, 4 samples

$\therefore A \rightarrow$  very large may only take this into consideration.  
 $\therefore$  It may be not be the right neighbour.

from sklearn.preprocessing import StandardScaler.

scaler = StandardScaler()

x\_scaled = scaler.fit\_transform(x)

mod value  $\rightarrow$  similar line

(all datapoints lying in some range)

\* Min Max Normalization/ Feature Scaling

$\rightarrow$  converts feature values within a specific range (mostly (0-1))

$$x_{\text{norm}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

(NO OUTLIERS should be present)

$\therefore$  upper and lower bound may be influenced

Feature clipping (numpy.clip) → all values in an array  
 $0-100 \rightarrow$  as is  
 above 100 → clip to 100  
 ↳ defining lower and upper values

Ex: [10, 15, 200, 5, 25, 180]  
 ↓  
 clip  
 ↓  
 100

3 & limit

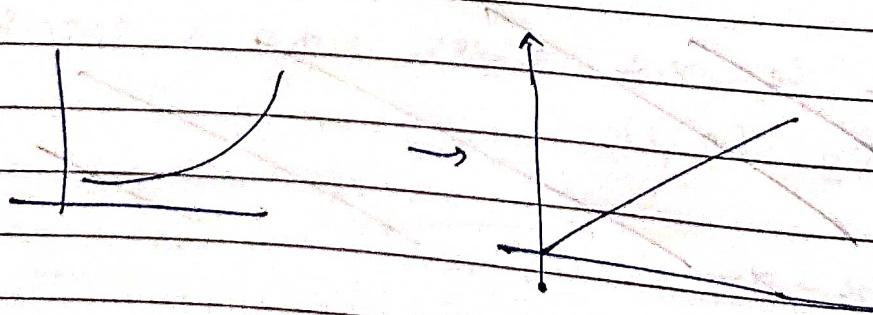
lower threshold = 0

higher threshold = 100

- ↳ Takes care of outliers, by limiting or bounding the extreme values of a feature with a specified range.
- ↳ determine clipping thresholds.
- ↳ Caps extreme values: The feature is constrained within a specific range.

#### \* log scaling

- ↳ When we deal with data that spans a wide range of values.
- ↳ It computes the log to compress a wide range to a narrow range, for a balanced visualization.



## SUMMARY

Normalization  
Technique

formula

when?

standardization

$$x_{\text{std}} = \frac{x_i - \mu}{\sigma}$$

when feature distribution does  
not contain extreme outliers.

Min-Max Scaling

$$x_{\text{norm}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$
 when the feature is one-

or less uniformly distributed

scaled to range in [0,1] across a fixed range.

Feature Clipping

If  $x > \text{max}$ ,  $x' = \text{max}$  when the feature

$x < \text{min}$ ,  $x' = \text{min}$  contains some extreme  
outliers

Log Scaling

$$x' = \log(x)$$