


Data Pre-Processing-I

(Introduction, Need, Data Cleaning)

Dr. JASMEET SINGH
ASSISTANT PROFESSOR, CSED
TIET, PATIALA

A solid orange horizontal bar spanning the width of the slide, located at the bottom.

Data Pre-Processing

- **Data Pre-processing:** It is that phase of any machine learning process, which transforms, or encodes, the data to bring it to such a state where it can be easily interpreted by the learning algorithm.

“Data pre-processing is not a single standalone entity but a collection of multiple interrelated tasks”

“Collectively data pre-processing constitutes majority of the effort in machine learning process (approx. 90 %)”

Data Pre-Processing

- **Data Pre-processing:** It is that phase of any Machine Learning process, which transforms, or Encodes, the data to bring it to such a state where it can be easily interpreted by the learning algorithm.

“Data pre-processing is not a single standalone entity but a collection of multiple interrelated tasks”

“Collectively data pre-processing constitutes majority of the effort in machine learning process (approx. 90 %)”

Need of Data Pre-Processing

- Data in the real world is “quite messy”
 - **incomplete**: missing feature values, absence of certain crucial feature, or containing only aggregate data.
 - e.g. Height=“ ”
 - **noisy**: containing errors or outliers
 - e.g. Weight=“5000” or “-60”
 - **inconsistent**: containing discrepancies in feature values.
 - e.g. Age=“20” and dob=“12 july 1990”
 - e.g. contradictions between duplicate records

Need for data Pre-processing

➤ **Unstructured Data (Text)**

- Lower Case
- Normalization (remove punctuation, special symbols, urls)
- Stopwords Removal (of, and, the,...)
- Stemming/Lemmatization (plays, playing, played → play)

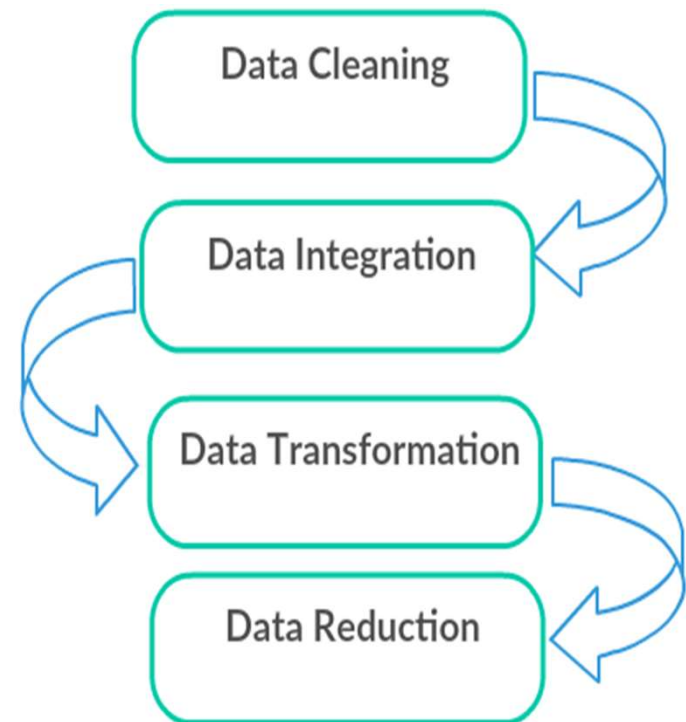
➤ **Unstructured Data (Images)**

- Read image
- Resize image
- Remove noise(Denoise)
- Segmentation
- Morphology(smoothing edges)

Pre- Processing in Structured Data

➤ Major data pre-processing tasks

- Data cleaning
- Data integration
- Data transformation
- Data reduction



Data Cleaning

- **Data cleaning:** It is a procedure to "clean" the data by filling in missing values, smoothening noisy data, identifying or removing outliers, and resolving data inconsistencies.
- Data cleaning tasks
 - Fill missing values
 - Noise smoothening and outlier detection
 - Resolving inconsistencies

Data Cleaning- Missing Values

Missing values: data values are not available.

i.e. many data entities have no data values corresponding to a certain feature like BMI value missing for some persons in a diabetes dataset.

- Probable reasons for missing values:
 - faulty measuring equipment
 - reluctance of person to share certain detail
 - negligence on part of data entry operator
 - feature unimportance at time of data collection

Data Cleaning- Missing Values Contd...

- Missing data handling techniques
 - Removing the data entity
 - Manually filling the values
 - Imputation (process used to determine and assign replacement values for missing, invalid, or inconsistent data)

“Technique selection is specific to user’s preference, dataset or feature type or problem set”

Data Cleaning- Missing Values Contd...


➤ Sample dataset related to forest fires

Month	FFMC	DC	temp	RH	wind
mar	86.2	94.3	8.2	51	6.7
oct	90.6	669.1	18	33	0.9
oct	90.6	686.9	NaN	33	NaN
mar	NaN	77.5	8.3	97	4
mar	89.3	102.2	11.4	99	1.8
aug	92.3	NaN	22.2	NaN	NaN
aug	NaN	495.6	24.1	27	NaN
aug	91.5	608.2	8	86	2.2
sep	91	692.6	NaN	63	5.4
sep	92.5	698.6	22.8	40	4

Data Cleaning- Missing Values Contd...

Removing the data entity: Most easiest way directly to clean the data, but this is usually discouraged as it leads to loss of data, as you are removing the data entity or feature values that can add value to data set as well.

Month	FFMC	DC	temp	RH	wind
mar	86.2	94.3	8.2	51	6.7
oct	90.6	669.1	18	33	0.9
oct	90.6	686.9	NaN	33	NaN
mar	NaN	77.5	8.3	97	4
mar	89.3	102.2	11.4	99	1.8
aug	92.3	NaN	22.2	NaN	NaN
aug	NaN	495.6	24.1	27	NaN
aug	91.5	608.2	8	86	2.2
sep	91	692.6	NaN	63	5.4
sep	92.5	698.6	22.8	40	4




Month	FFMC	DC	temp	RH	wind
mar	86.2	94.3	8.2	51	6.7
oct	90.6	669.1	18	33	0.9
mar	89.3	102.2	11.4	99	1.8
aug	91.5	608.2	8	86	2.2
sep	92.5	698.6	22.8	40	4

Data Cleaning- Missing Values Contd...

- **Manually filing up of values** : This approach is time consuming, and not recommended for huge data sets.

Month	FFMC	DC	temp	RH	wind
mar	86.2	94.3	8.2	51	6.7
oct	90.6	669.1	18	33	0.9
oct	90.6	686.9	NaN	33	NaN
mar	NaN	77.5	8.3	97	4
mar	89.3	102.2	11.4	99	1.8
aug	92.3	NaN	22.2	NaN	NaN
aug	NaN	495.6	24.1	27	NaN
aug	91.5	608.2	8	86	2.2
sep	91	692.6	NaN	63	5.4
sep	92.5	698.6	22.8	40	4



Month	FFMC	DC	temp	RH	wind
mar	86.2	94.3	8.2	51	6.7
oct	90.6	669.1	18	33	0.9
oct	90.6	686.9	17	33	0.8
mar	91.6	77.5	8.3	97	4
mar	89.3	102.2	11.4	99	1.8
aug	92.3	380	22.2	92	1.8
aug	90	495.6	24.1	27	2
aug	91.5	608.2	8	86	2.2
sep	91	692.6	22	63	5.4
sep	92.5	698.6	22.8	40	4

Data Cleaning- Missing Values Contd...

- **Imputation** : process used to determine and assign replacement values for missing, invalid, or inconsistent data. Various imputation methods include:
 - Central Tendency Imputation
 - Hot Deck Imputation
 - Cold Deck Imputation
 - Model Based Imputation
 - Nearest Neighbor Imputation
 - Tree-Based Imputation

Data Cleaning (Missing values)- Contd...

Central tendency Imputation : Replacing the missing value by central tendency (mean, median, mode) for a feature vector or belonging to same class of feature vector.

Mean $\bar{x} = \frac{\sum_{i=1}^N x_i}{N}$

Median


$$md = x_{\frac{(n-1)}{2}} \text{ for } n \text{ is odd}$$
$$md = \frac{1}{2} \left(x_{\frac{n}{2}} + x_{\frac{n}{2}+1} \right) \text{ for } n \text{ is even}$$

Mode : Mode is the most frequent value corresponding to a certain feature in a given data set

Data Cleaning- Missing Values Contd...

- Replacing Mean Value:

Month	FFMC	DC	temp	RH	wind
mar	86.2	94.3	8.2	51	6.7
oct	90.6	669.1	18	33	0.9
oct	90.6	686.9	NaN	33	NaN
mar	NaN	77.5	8.3	97	4
mar	89.3	102.2	11.4	99	1.8
aug	92.3	NaN	22.2	NaN	NaN
aug	NaN	495.6	24.1	27	NaN
aug	91.5	608.2	8	86	2.2
sep	91	692.6	NaN	63	5.4
sep	92.5	698.6	22.8	40	4



Month	FFMC	DC	temp	RH	wind
mar	86.2	94.3	8.2	51	6.7
oct	90.6	669.1	18	33	0.9
oct	90.6	686.9	15.3	33	3.57
mar	90.5	77.5	8.3	97	4
mar	89.3	102.2	11.4	99	1.8
aug	92.3	458.3	22.2	58.7	3.57
aug	90.5	495.6	24.1	27	3.57
aug	91.5	608.2	8	86	2.2
sep	91	692.6	15.3	63	5.4
sep	92.5	698.6	22.8	40	4

Data Cleaning- Missing Values Contd...

- “ Replacing by mean value: Not a suitable method if data set has many outliers”
- For example: weighs of humans 67, 78, 900, -56, 389, -1 etc. Outlier
Mean is 229.5
- Can be replaced with median in such cases.
- “Mode is a good option for missing values in case of categorical variables”

Data Cleaning- Missing Values Contd...

Hot Deck Imputation

- Computes how many number of features (other than feature with missing data) have same values in the entire training examples and choose it for replacement.
- Used mostly in categorical data.

Cold Deck Imputation

- Similar to hot deck imputation.
- In it missing observations are replaced by values from a source unrelated to the data set under consideration.

Data Cleaning- Missing Values Contd...

Nearest Neighbor- Based Imputation

- Rely on distance metrics
- evaluate the distance between recipients and donors.
- Used after converting all features to numerical (quantitative)

Handling Missing Values using Random Forests

Random Forests Classifiers can also be used to handle the missing values in the dataset.

- Whenever, there is some missing values in training data, random forest applies following general technique,” It makes an initial guess and then gradually refines the guess until it is hopefully good.
- Initial bad guess can be the mode for categorical feature and mean/median for a numerical feature.
- The guesses are refined by finding samples which are similar to one with missing data.

Handling Missing Values using Random Forests (Contd...)

- In order to determine similarity following steps are followed.
 - i. Build a random forest.
 - ii. Run all the data down all the trees.
 - iii. Note the samples that end up at same leaf nodes. Keep track of similar samples using a Proximity Matrix.
 - iv. Then divide each entry in proximity matrix with total number of trees.
 - v. Use proximity values to make better guesses about the missing data i.e. calculated weighted frequency of each value using proximity values as weight.
 - vi. Repeat steps 1 to 5 until the missing values converge.

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	???	???	No

However, for patient #4, we've got some missing data.

Since No is the most commonly occurring value (it occurs 2 out of 3 times)

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	???	No

So "No" is our initial guess.

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	180	No

In this case, the median value is 180

Step 2: Run all of the data down all of the trees.

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	180	No

Step 1: Build a random forest...

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	180	No

etc. etc. etc.

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	180	No

Now we want to refine these guesses.

We do this by first determining which samples are similar to the one with missing data.

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	180	No

Ultimately, we run the data down all the trees and the proximity matrix fills in.

	1	2	3	4
1		2	1	1
2	2		1	1
3	1	1		8
4	1	1	8	

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	180	No

	1	2	3	4
1		0.2	0.1	0.1
2	0.2		0.1	0.1
3	0.1	0.1		0.8
4	0.1	0.1	0.8	

Then we divide each proximity value by the total number of trees. In this example, assume we had 10 trees.

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	???	???	No

For Blocked Arteries, we calculate the weighted frequency of "Yes" and "No", using proximity values as the weights.

	1	2	3	4
1		0.2	0.1	0.1
2	0.2		0.1	0.1
3	0.1	0.1		0.8
4	0.1	0.1	0.8	

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	No	???	No

The weighted frequency for "No" is...

Yes = $\frac{1}{3} \times 0.1 = 0.03$

No = $\frac{2}{3} \times 0.9 = 0.6$

Yes = 1/3

No = 2/3

"No" has a way higher weighted frequency, so we'll go with it.

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	???	???	No

The weighted frequency for "Yes" is...

Yes = $\frac{1}{3} \times 0.1 = 0.03$

The weighted frequency for "Yes"

Yes = 1/3

No = 2/3

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	???	???	No

The weighted frequency for "Yes" is...

Yes = $\frac{1}{3} \times$ The weight for "Yes"

Yes = 1/3

No = 2/3

The weight for "Yes" = $\frac{\text{Proximity of "Yes"}}{\text{All Proximities}}$

$$\text{Weighted average} = (125 \times 0.1) + (180 \times 0.1) + (210 \times 0.8)$$

$$= 198.5$$

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	NO	???	No

	1	2	3	4
1		0.2	0.1	0.1
2	0.2		0.1	0.1
3	0.1	0.1		0.8
4	0.1	0.1	0.8	

Filled-in Missing Values

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	NO	198.5	No

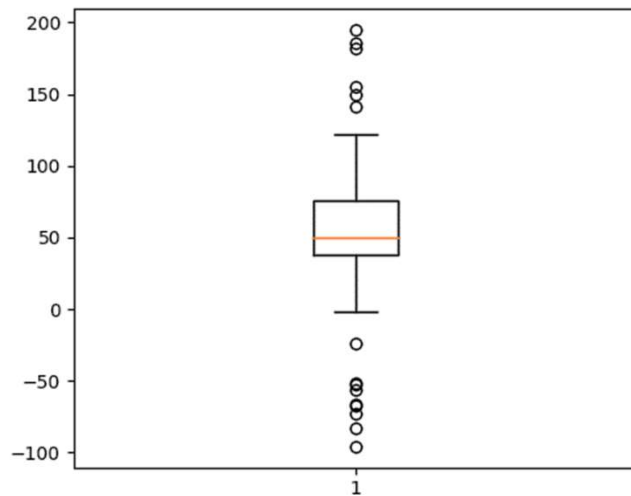
Now that we've revised our guesses a little bit, we do the whole thing over again...

We build a random forest, run the data through the trees, recalculate the proximities and recalculate the missing values.

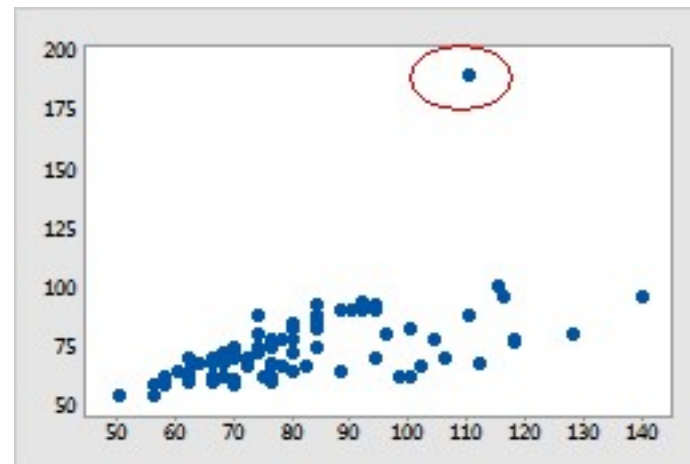
We do this 6 or 7 times until the missing values converge (i.e. no longer change each time we recalculate).

Data Cleaning- Noisy Data

- **Noise** is defined as a random variance in a measured variable.
- For numeric values, boxplots and scatter plots can be used to identify outliers.



Boxplot



Scatter plot

Data Cleaning- Noisy Data

➤ Major reasons of random variations in data are:

- Malfunctioning of collection instruments.
- Data entry lags.
- Data transmission problems

To deal with these anomalous values, data smoothing techniques are applied, some of the popular ones are

- Binning method
- Outlier analysis

Binning Method for Noisy Data

Binning method : performs the task of data smoothening.

Steps to be followed under binning method are:

Step 1: Sort the data into ascending order.

Step 2: Calculate the bin size (i.e. number of bins)

Step 3: Partition or distribute the data equally among the bins starting with first element of sorted data.

Step 4: perform data smoothening using **bin means, bin boundaries, and bin median**.

Last bin can have one less or more element!!

Binning Method for Noisy Data

Example : 9, 21, 29, 28, 4, 21, 8, 24, 26

Step1: sorted the data 4, 8, 9, 21, 21, 24, 26, 28, 29

Step 2 : Bin size calculation

$$\begin{aligned}\text{Bin size} &= \frac{\text{Max value} - \text{Min value}}{\text{data size}} \\ &= \frac{29-4}{9} = 2.777\end{aligned}$$

But we need to take ceiling value, so bin size is 3 here

Binning Method for Noisy Data

- Step 3 : Bin partitioning (equi-size bins)

Bin 1: 4, 8, 9

Bin 2: 21, 21, 24

Bin 3: 26, 28, 29

Step 4 : data smoothening

- Using mean value : replace the bin values by bin average

Bin 1: 7, 7, 7

Bin 2: 22, 22, 22

Bin 3: 27, 27, 27

Binning Method for Noisy Data

- ➤ Using boundary values : replace the bin value by a closest boundary value of the corresponding bin.

Bin 1: 4, 9, 9

Bin 2: 21, 21, 24

Bin 3: 26, 29, 29

“Boundary values remain unchanged in boundary method”

- Using median values : replace the bin value by a bin median.

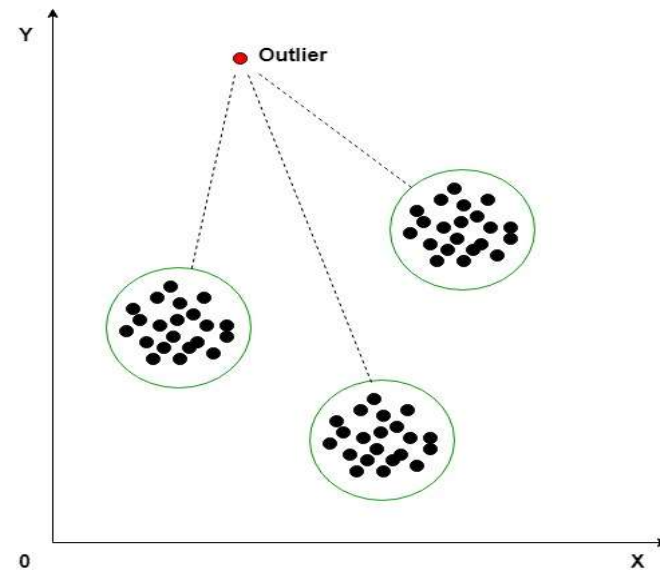
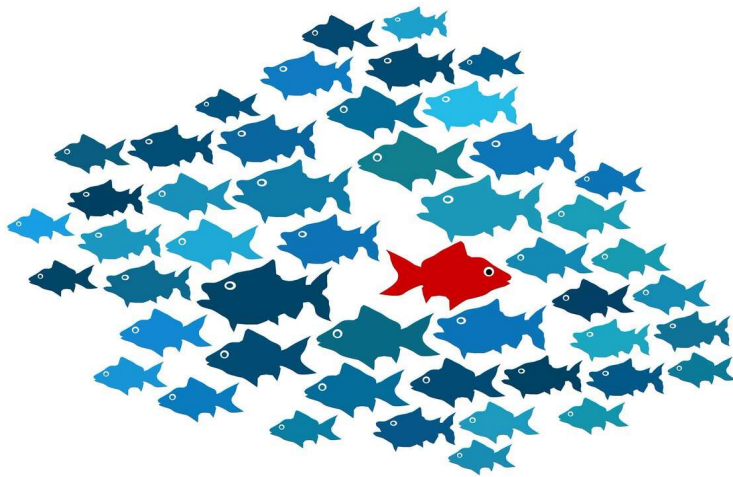
Bin 1: 8, 8, 8

Bin 2: 21, 21, 21

Bin 3: 28, 28, 28

Outlier Analysis

- An **outlier** is an object that deviates significantly from the rest of the objects.
- They can be caused by measurement or execution error.
- The analysis of outlier data is referred to as outlier analysis or outlier mining.



Outlier Analysis

- Types of Outliers:
 - Univariate: A univariate outlier is a data point that consists of an extreme value on one variable.
 - Multivariate: A multivariate outlier is a combination of unusual scores on at least two variables.
- Outlier Detection and Handling Methods
 - Extreme Value Analysis
 - Linear Models
 - Proximity-based Methods
 - Information Theoretic Methods
 - Isolation Forests Based Methods

Extreme Value Analysis- Outlier Analysis

■ Numeric Outlier

- This is the simplest, nonparametric outlier detection method in a one dimensional feature space.
- Here outliers are calculated by means of the *IQR* (InterQuartile Range).
- The first and the third quartile ($Q1$, $Q3$) are calculated.
- An outlier is then a data point x_i that lies outside the interquartile range. That is:

$$x_i > Q_3 + k(IQR) \text{ and } x_i < Q_1 - k(IQR) \\ \text{where } IQR = Q_3 - Q_1 \text{ and } k \geq 0$$

Assume the data 6, 2, 1, 5, 4, 3, 50. If these values represent the number of chapatis eaten in lunch, then 50 is clearly an outlier.

Sorted Values: 1, 2, 3, 4, 5, 6, 50

Q1 25 percentile of the given data is, 2

Q2 50 percentile of the given data is, 4.0

Q3 75 percentile of the given data is, 6

$IQR = 6 - 2 = 4$, $k = 1.5$

Range: $6 - 1.5 \times 4 = 0$ and $6 + 1.5 \times 4 = 12$

50 is Outlier

Extreme Value Analysis- Outlier Analysis

- **Z-score** is a parametric outlier detection method in a one or low dimensional feature space.
- This technique assumes a Gaussian distribution of the data.
- The outliers are the data points that are in the tails of the distribution and therefore far from the mean.
- How far depends on a set threshold z_{thr} for the normalized data points z_i calculated with the formula:

$$z_i = \frac{x_i - \mu}{\sigma}$$

where x_i is a data point, μ is the mean of all x_i and σ is the standard deviation of all x_i .

An outlier is then a normalized data point which has an absolute value greater than z_{thr} .

$$|z_i| > z_{thr}$$

Commonly used z_{thr} values are 2.5, 3.0 and 3.5.

Outlier Analysis

- **Linear Models:**

- Projection methods that model the data into lower dimensions using linear correlations.
- For example, principle component analysis and data with large residual errors may be outliers.

- **Proximity-based Models:**

- Data instances that are isolated from the mass of the data as determined by cluster, density or nearest neighbor analysis.

- **Information Theoretic Models:**

- Outliers are detected as data instances that increase the complexity (minimum code length) of the dataset.

Outlier Analysis-Isolation Forest Based Methods

- Isolation Forests(IF), similar to Random Forests, are build based on decision trees. And since there are no pre-defined labels here, it is an unsupervised model.
- In an Isolation Forest, randomly sub-sampled data is processed in a tree structure based on randomly selected features.
- The samples that travel deeper into the tree are less likely to be anomalies as they required more cuts to isolate them.
- Similarly, the samples which end up in shorter branches indicate anomalies as it was easier for the tree to separate them from other observations.

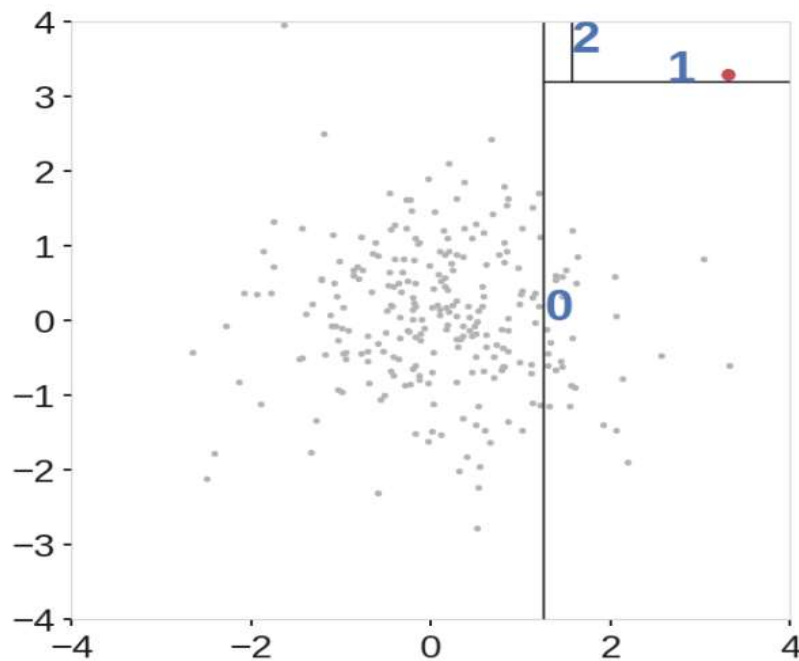
Outlier Analysis-Isolation Forest Based Methods (Contd...)

- The algorithm starts with the training of the data, by generating Isolation Trees:
 1. When given a dataset, a random sub-sample of the data is selected and assigned to a binary tree.
 2. Branching of the tree starts by selecting a random feature (from the set of all N features) first. And then branching is done on a random threshold (any value in the range of minimum and maximum values of the selected feature).
 3. If the value of a data point is less than the selected threshold, it goes to the left branch else to the right. And thus a node is split into left and right branches.
 4. This process from step 2 is continued recursively till each data point is completely isolated or till max depth(if defined) is reached.
 5. The above steps are repeated to construct random binary trees.

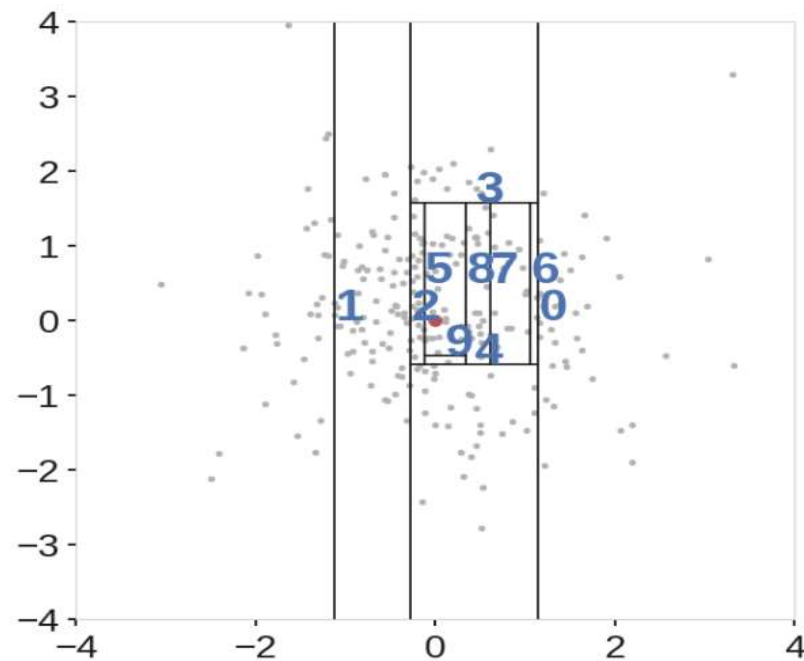
Outlier Analysis-Isolation Forest Based Methods (Contd...)

- After an ensemble of iTrees(Isolation Forest) is created, model training is complete.
- During scoring, a data point is traversed through all the trees which were trained earlier.
- Now, an 'anomaly score' is assigned to each of the data points based on the depth of the tree required to arrive at that point. This score is an aggregation of the depth obtained from each of the iTrees.

Outlier Analysis-Isolation Forest Based Methods (Contd...)



(a) Anomaly point



(b) Nominal point

Data Cleaning – Inconsistent Data

Inconsistent Data: discrepancies between different data items.

e.g. the “Address” field contains the “Phone number”

To resolve inconsistencies

- Manual correction using external references
- Semi-automatic tools
 - To detect violation of known functional dependencies and data constraints
 - To correct redundant data

To avoid inconsistencies, perform data assessment like

knowing what the data type of the features should be and whether it is the same for all the data objects.”

Data Cleaning-Summary

