

# **AI in Anomaly Detection and Image Compression**

*A Project Report Submitted  
in Partial Fulfillment of the Requirements  
for the Degree of*

**Bachelor of Technology**

*by*

**Kaushal Kishore**  
(111601008)

*under the guidance of*

**Dr. Chandra Shekar**



INDIAN INSTITUTE  
OF TECHNOLOGY  
**PALAKKAD**

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

# CERTIFICATE

*This is to certify that the work contained in this thesis entitled “**AI in Anomaly Detection and Image Compression**” is a bonafide work of **Kaushal Kishore (Roll No. 111601008)**, carried out in the Department of Computer Science and Engineering, Indian Institute of Technology Palakkad under my supervision and that it has not been submitted elsewhere for a degree.*

**Dr. Chandra Shekar**

Assistant Professor

Department of Computer Science & Engineering

Indian Institute of Technology Palakkad

# Acknowledgements

Apart from the efforts of myself, the success of any project depends largely on the encouragement and guidelines of many others. I take this opportunity to express my gratitude to the people who have been instrumental in the successful completion of this project. I would like to show my greatest appreciation to Dr. Chandra Shekar. I can't say thank you enough for his tremendous support and help. I feel motivated and encouraged every time I attend his meeting. Without his encouragement and guidance, this project would not have materialized.

# Contents

<b>List of Figures</b>	<b>iv</b>
<b>List of Tables</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Anomaly Detection . . . . .	1
1.2 Use Cases . . . . .	2
1.3 Challenges . . . . .	3
1.3.1 Masking and Swamping . . . . .	3
1.3.2 Concept drift . . . . .	3
1.3.3 Miscellaneous . . . . .	3
1.4 Organization of The Report . . . . .	4
<b>2 Isolation Forest</b>	<b>7</b>
2.1 Isolation based anomaly detection . . . . .	7
2.2 Training Stage . . . . .	9
2.3 Evaluation Stage . . . . .	10
2.4 Anomaly Score . . . . .	11
2.5 Isolation vs. Distance and Density . . . . .	12
2.6 Conclusion . . . . .	14

<b>3</b>	<b>PIDForest</b>	<b>15</b>
3.1	Issues with Isolation Forest . . . . .	15
3.2	Partial Identification . . . . .	16
3.2.1	Boolean Setting . . . . .	16
<b>4</b>	<b>Contributions</b>	<b>19</b>
4.1	Construction . . . . .	19
4.2	Improved Method . . . . .	19
4.3	Conclusion . . . . .	19
<b>5</b>	<b>AI in Image Compression</b>	<b>21</b>
5.1	Construction . . . . .	21
5.2	Improved Method . . . . .	21
5.3	Conclusion . . . . .	21
<b>6</b>	<b>Conclusion and Future Work</b>	<b>23</b>
	<b>References</b>	<b>25</b>

# List of Figures

2.1	<b>Left:</b> a normal point $x_i$ requires twelve random partitions to be isolated; <b>Right:</b> an anomaly $x_o$ requires only four partitions to be isolated. . . . .	8
2.2	Averaged path lengths of $x_i$ and $x_o$ converge when the number of trees increases. . . . .	8
2.3	Isolation Forest . . . . .	12
2.4	High density and short distance do not always imply normal instances. . .	13
2.5	Low density and long distance do not always imply anomalies. . . . .	13

# List of Tables

# Chapter 1

## Introduction

This chapter discusses anomaly detection, its use cases and some major challenges.

### 1.1 Anomaly Detection

Anomaly detection (also outlier detection) is the identification of rare items, events or observations which raise suspicions by differing significantly from the majority of the data. Typically, the anomalous items will translate to some kind of problem such as bank fraud, a structural defect, medical problems or errors in a text. Anomalies also known as outliers, novelties, noise, deviations and exceptions.

Three broad categories of anomaly detection techniques exist:

1. Unsupervised anomaly detection techniques detect anomalies in an unlabeled test data set under the assumption that the majority of the instances in the data set are normal by looking for instances that seem to fit least to the remainder of the data set.
2. Supervised anomaly detection techniques require a data set that has been labeled as "normal" and "abnormal" and involves training a classifier (the key difference to



many other statistical classification problems is the inherently unbalanced nature of outlier detection).

3. Semi-supervised anomaly detection techniques construct a model representing normal behavior from a given normal training data set, and then test the likelihood of a test instance to be generated by the learnt model.

We will restrict ourselves to unsupervised anomaly detection and semi-supervised anomaly detection problem.

## 1.2 Use Cases

The ability to detect anomalies has significant relevance, and anomalies often provides critical and actionable information in various application domains.

Identification of potential outliers is important for the following reasons: [1]

1. An outlier may indicate bad data. For example, the data may have been coded incorrectly, or an experiment may not have been run correctly. If it can be determined that an outlying point is in fact erroneous, then the outlying value should be deleted from the analysis (or corrected if possible).
2. In some cases, it may not be possible to determine if an outlying point is bad data. Outliers may be due to random variation or may indicate something scientifically interesting. In any event, we typically do not want to simply delete the outlying observation.

For example, anomalies in credit card transactions could signify fraudulent use of credit cards. An anomalous spot in an astronomy image could indicate the discovery of a new star. An unusual computer network traffic pattern could stand for unauthorised access. These applications demand anomaly detection algorithms with high detection accuracy and fast execution.

## 1.3 Challenges

### 1.3.1 Masking and Swamping

Masking and swamping is the biggest problem affecting any anomaly detection algorithm.

Masking is the existence of too many anomalies concealing their own presence. It happens when anomaly clusters become large and dense. For example, if we are testing for a single outlier when there are in fact more outliers, these additional outliers may influence the value of the test statistic enough so that no points are declared as outliers.

On the other hand, swamping refers to situations where normal instances are wrongly identifying as anomalies. It happens when the number of normal instances increases, or they become more scattered. For example, if we are testing for two or more outliers when there is in fact only a single outlier, both points may be declared outliers.

Masking is one reason that trying to apply a single outlier test sequentially can fail. For example, if there are multiple outliers, masking may cause the outlier test for the first outlier to return a conclusion of no outliers. So the testing is not performed for any additional outliers.

### 1.3.2 Concept drift

In the case of streaming data, the anomaly context can change over time. For example, consider a user's behaviour change from one system to another. The anomaly detection algorithm should adapt to this change in the behaviour of the external agent. This deviation of the normal behaviour time to time is called concept drift. Any online anomaly detection algorithm must have a way to deal with this.

### 1.3.3 Miscellaneous

Apart from the above conceptual challenges, here are some general challenges pointed out by Vatsal et al. [2]

**High dimensional, heterogeneous data:** The data collected could contains measurements of metrics like cpu usage, memory, bandwidth, temperature, in addition to categorical data such as day of the week, geographic location, OS type. This makes finding an accurate generative model for the data challenging. The metrics might be captured in different units, hence algorithms that are unit-agnostic are preferable. The algorithm needs to scale to high dimensional data.

**Scarce labels:** Most of the data are unlabeled. Generating labels is time and effort intensive and requires domain knowledge. Hence supervised methods are a non-starter, and even tuning too many hyper-parameters of unsupervised algorithms could be challenging.

**Irrelevant attributes:** Often an anomaly manifests itself in a relatively small number of attributes among the large number being monitored. For instance, a single machine in a large datacenter might be compromised and behave abnormally.

## 1.4 Organization of The Report

This chapter [1] provides a background for the topics covered in this report. We provided a description of anomaly detection problem and discussed some use cases. Then we discussed some challenges to anomaly detection problem: masking, swamping and concept drift. In the next chapter [2] we will discuss a very efficient ensemble method Isolation Forest for anomaly detection. In chapter [3] we will discuss another ensemble method PIDForest which has been recently developed. The major drawback of the above mentioned algorithms is that they are used in offline setting without dealing with concept drift. Most of the anomaly detection algorithm is offline and fail to address the problem of concept drift. In chapter [4] we will present some methods to address these issues. In chapter [5] we will review our work did till mid-sem. And finally in chapter [6], we conclude with some future works.

**Code Repository:** <https://github.com/KishoreKaushal/AnomalyDetection>

**Report Repository:** <https://github.com/KishoreKaushal/btp-report>



# Chapter 2

## Isolation Forest

### 2.1 Isolation based anomaly detection

Isolation is the process or fact of isolating or being isolated. The authors of [3] proposed an isolation based anomaly detection which takes advantage of two quantitative properties of anomalies:

1. They are the minority consisting of few instances.
2. They have attribute-values that are very different from those of normal instances.

Hence, anomalies are 'few and different' which make them more susceptible to a mechanism we called Isolation. Isolation can be implemented by any means that separates instances. Lui et al. [3] proposed to use a binary tree structure called isolation tree (iTree) which can be constructed effectively to isolate instances. Because of the susceptibility to isolation, anomalies are more likely to be isolated closer to the root of an iTree; whereas normal points are more likely to be isolated at the deeper end of an iTree.

The proposed method, called Isolation Forest (iForest) builds an ensemble builds an ensemble of iTrees for a given data set. Anomalies are those instances which have short average path lengths on the iTrees. There are two training parameters and one evaluation

parameter in this method: the training parameters are the number of trees to build and subsampling size. The evaluation parameter is the tree height limit during evaluation.



**Fig. 2.1:** Left: a normal point  $x_i$  requires twelve random partitions to be isolated; Right: an anomaly  $x_o$  requires only four partitions to be isolated.

**Source:** Lui et al. [3]



**Fig. 2.2:** Averaged path lengths of  $x_i$  and  $x_o$  converge when the number of trees increases.

**Source:** Lui et al. [3]

## 2.2 Training Stage

Formally, isolation tree is defined as follows:

**Definition: 1** Let  $T$  be a node of an isolation tree.  $T$  is either an external-node with no child, or an internal-node with one test and exactly two daughter nodes  $(T_l, T_r)$ . A test at node  $T$  consists of an attribute  $q$  and a split value  $p$  such that the test  $q < p$  determines the traversal of a data point to either  $T_l$  or  $T_r$ . Let  $X = \{x_1, \dots, x_n\}$  be the given data set of a  $d$ -variate distribution. A sample of instances  $X' \subset X$  is used to build an isolation tree<sup>[2]</sup>. We recursively divide  $X'$  by randomly selecting an attribute  $q$  and a split value  $p$ , until either:  
*i) Node has only one instance ii) Or, all data at the node have the same values.*

**Definition: 2** Isolation forest is defined as 4-tuple  $(X, t, \psi, S)$  where

- $X$  is input data,
- $t$  is number of trees,
- $\psi$  is subsampling size and
- $S$  is the set of isolation trees.

The elements of set  $S$  is constructed<sup>[1]</sup> by sampling  $\psi$  instances from  $X$  without replacement.

---

### Algorithm 1: $iForest(X, t, \psi)$

---

**Complexity:** Time -  $O(t\psi^2)$ , Space -  $O(t\psi)$

**Input:**  $X$  - input data,  $t$  - number of trees,  $\psi$  - subsampling size

**Output:** List of  $iTrees$

```

1 Forest  $\leftarrow$  EmptyList
2 for  $i = 1$  to  $t$  do
3    $X' \leftarrow sampleWithoutReplacement(X, \psi)$ 
4   Forest  $\leftarrow$  Forest  $\cup iTree(X')$ 
5 end
6 return Forest
```

---



---

**Algorithm 2:** *iTree*( $X$ )

---

**Complexity:** Time -  $O(\psi^2)$ , Space -  $O(\psi)$

**Input:**  $X$  - input data

**Output:** an *iTree*

```
1  $q \leftarrow \text{RandomChoice}(X.\text{attributes})$ 
2  $p \leftarrow \text{RandomNumber}(X[\text{splitAttr}].\text{min}(), X[\text{splitAttr}].\text{max}())$ 
3  $\text{tree} \leftarrow \text{Node} \{ \text{left} \leftarrow \text{None}, \text{right} \leftarrow \text{None}, \text{size} \leftarrow X.\text{size}$ 
4        $\text{splitAttr} \leftarrow q, \text{splitVal} \leftarrow p \}$ 
5 if  $X.\text{size} > 1$  and  $X[\text{splitAttr}].\text{numUnique}() > 1$  then
6   |  $X_l \leftarrow X.\text{where}(q < p)$ 
7   |  $X_r \leftarrow X.\text{where}(q \geq p)$ 
8   |  $\text{tree.left} \leftarrow \text{iTree}(X_l)$ 
9   |  $\text{tree.right} \leftarrow \text{iTree}(X_r)$ 
10 end
11 return  $\text{tree}$ 
```

---

## 2.3 Evaluation Stage

---

**Algorithm 3:** *PathLength*( $x$ )

---

**Complexity:** Time -  $O(t\psi)$ , Space -  $O(1)$

**Input:**  $x$  - input instance,  $T$  - an *iTree*,  $hlim$  - height limit,  $e$  - current path length to be initialized to zero when called first time

**Output:** path length of  $x$

```
1 if ( $T.\text{right}$  is  $\text{None}$ ) and ( $T.\text{left}$  is  $\text{none}$ ) and ( $e \geq hlim$ ) then
2   | return  $e + c(T.\text{size})$  //  $c(\dots)$  is defined in Equation 2.1
3 end
4  $a \leftarrow T.\text{splitAttr}$ 
5 if  $x[a] < T.\text{splitVal}$  then
6   | return  $\text{PathLength}(x, T.\text{left}, hlim, e + 1)$ 
7 end
8 else
9   | return  $\text{PathLength}(x, T.\text{right}, hlim, e + 1)$ 
10 end
```

---

In the evaluation stage<sup>[3]</sup>, a single path length  $h(x)$  is derived by counting the number of edges  $e$  from the root node to an external node as instance  $x$  traverses through an iTree. When the traversal reaches a predefined height limit  $hlim$ , the return value is  $e$  plus an adjustment  $c(size)$ . This adjustment accounts for estimating an average path length of a random sub-tree which could be constructed using data of  $size$  beyond the tree height limit. When  $h(x)$  is obtained for each tree of the ensemble, an anomaly score is computed. The anomaly score and the adjustment  $c(size)$  is defined in the next section.

## 2.4 Anomaly Score

The difficulty in deriving an anomaly score from  $h(x)$  is that while the maximum possible height of iTree grows in the order of  $\psi$ , the average height grows in the order of  $\log \psi$ . When required to visualize or compare path lengths from models of different subsampling sizes, normalization of  $h(x)$  by any of the above terms either is not bounded or cannot be directly compared. Thus, a normalized anomaly score is needed for the aforementioned purposes.

Since iTrees have an equivalent structure to Binary Search Tree, the estimation of average  $h(x)$  for external node terminations is the same as that of the unsuccessful searches in BST.

Section 10.3.3 of [4] gives the average path length of unsuccessful searches in BST as:

$$c(\psi) = \begin{cases} 2H(\psi - 1) - 2(\psi - 1)/n & \text{for } \psi > 2, \\ 1 & \text{for } \psi = 2, \\ 0 & \text{otherwise.} \end{cases} \quad (2.1)$$

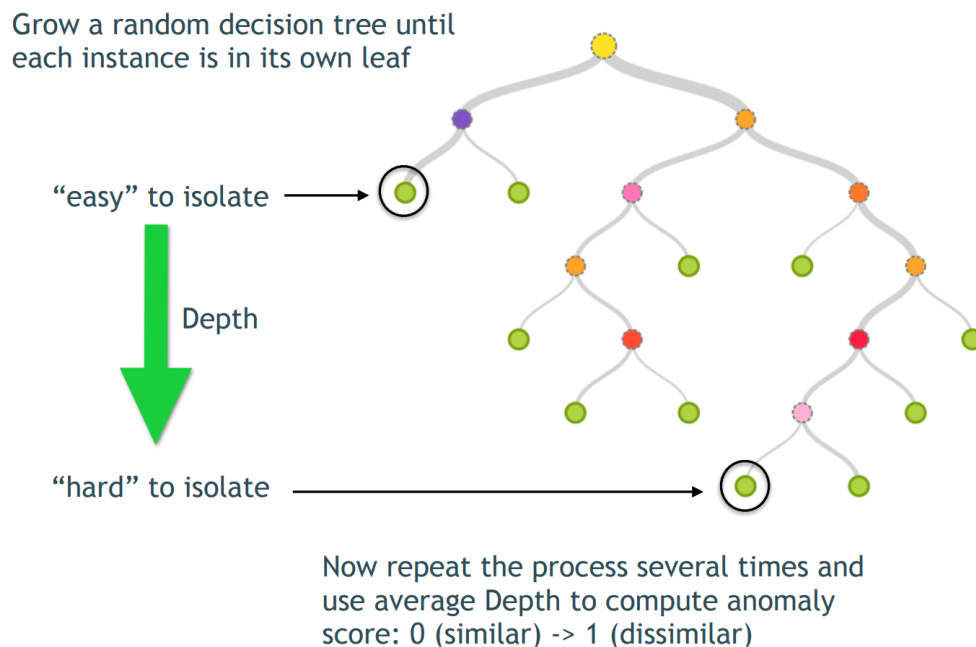
where  $H(i) \approx \ln(i) + 0.5772156649$  (euler's constant), is the harmonic number. As  $c(\psi)$  is the average of  $h(x)$  given  $\psi$ , we use it to normalise  $h(x)$ . The anomaly score  $s$  of an instance  $x$  is defined as:

$$s(x, \psi) = 2^{\frac{-E[h(x)]}{c(\psi)}} \quad (2.2)$$

where  $E[h(x)]$  is the average of  $h(x)$  from a collection of iTrees.

Anomaly score  $s(x, \psi)$  is interpreted as follows:

1. if  $s \approx 1$ , instance is abnormal
2. if  $s \approx 0$ , instance is nominal
3. if  $s \approx 0.5$ , no distinct anomaly



**Fig. 2.3:** Isolation Forest

Source: slideshare.com

In practical use cases, one has to set a threshold deciding the result. And finding a good threshold is a very difficult task. Most of the times the anomalous points' score overlaps with the nominal points' score.

## 2.5 Isolation vs. Distance and Density

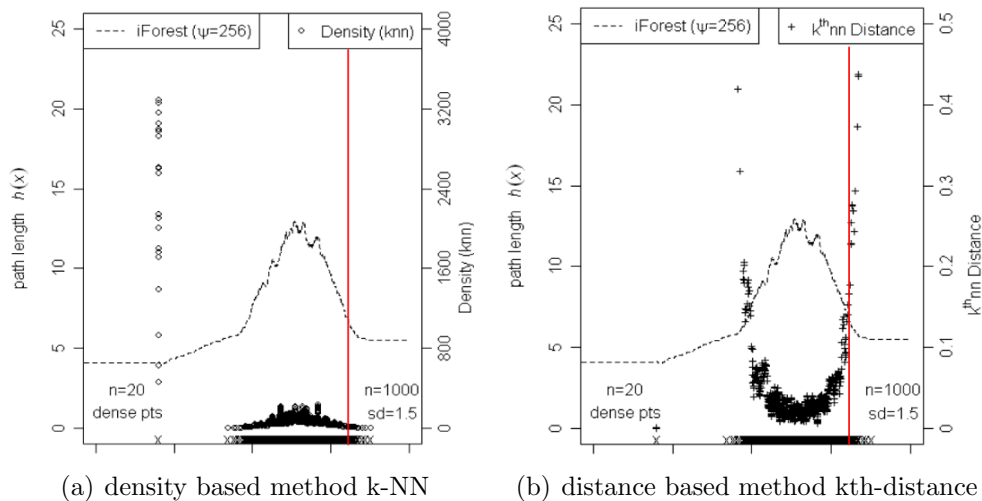
Using basic density measures, the assumption is that ‘Normal points occur in dense regions, while anomalies occur in sparse regions’. Using basic distance measures, the basic

assumption is that ‘Normal point is close to its neighbours and anomaly is far from its neighbours’.

There are violations of these assumptions: i) high density and short distance do not always imply normal instances ii) low density and long distance do not always imply anomalies.



**Fig. 2.4:** High density and short distance do not always imply normal instances.  
**Source:** Lui et al. [3]



**Fig. 2.5:** Low density and long distance do not always imply anomalies.  
**Source:** Lui et al. [3]

Isolation Forest compares favourably to distance and density based methods in terms

of accuracy and processing time. Distance and density based suffers immensely in terms of accuracy and processing time because of curse of dimensionality.

Distance based methods also suffers from masking and swamping effect. In isolation forest, masking and swamping effects can be managed by adjusting the *hlim* parameter during evaluation. Refer section 4.5, 5.3 and 5.4 of Lui et al. [3].

Codes for Isolation Forest and Isolation Tree can be found at this repository.

## 2.6 Conclusion

To conclude, isolation forest is one of the most efficient and accurate methods for anomaly detection. It outperforms various density and distance based methods like ORCA, DOLPHIN, LOF, ROF, etc. and their variants in terms of accuracy and performance. Still there are some shortcomings of the isolation forest that we will discuss in next chapter.

# Chapter 3

## PIDForest

In chapter 1 we discussed some of the challenges that an anomaly detection algorithm face. In the previous chapter we discussed an isolation based ensemble method which is a very good method in terms of complexity and accuracy. Isolation forest tried to address the problems related to masking and swamping. In this chapter we will point out some major issues with isolation forest and discuss another ensemble method, which improves upon those issues.

### 3.1 Issues with Isolation Forest

**Random Split:** iForest repeatedly samples a set  $X'$  of  $\psi$  points from  $X$  and builds a random tree with those points as leaves. The tree is built by choosing a random coordinate  $q$ , and a random value  $p$  in its range about which to split. Since iForest chooses which coordinate we split on as well as the breakpoint at random. Thus to be isolated at small depth frequently, picking splits at random must have a good chance of isolating an anomalous point. Although, there are other variants like extended isolation forest which improves on this, but there is no significant improvements.

**High Dimensional:** As the number of co-ordinates or attributes increases, probability of choosing a sequence of attributes for split which gives rise to most of the anomalies will be very less. Hence, it is very likely that anomalous points won't be isolated near the root and false negative cases will increase. (In anomaly detection problem, anomaly is the true class.)

**Presence of non-ordinal categorical attributes:** A very big limitation of iForest is that it only works for those datasets where all the features are real-values or ordinal.

## 3.2 Partial Identification

We will briefly discuss some concepts present in section 2 of Vatsal et al. [2].

**Notations:** Let  $T$  denote a dataset of  $n$  points in  $d$  dimensions. Given indices  $S \subseteq [d]$  and  $x \in R$ , let  $x_S$  denote the projection of  $x$  onto coordinates in  $S$ . All the logarithms are to base 2.

### 3.2.1 Boolean Setting

In Boolean setting  $T \subseteq \{0, 1\}^d$  and assume that  $T$  has no duplicates.

**Definition: 3** (*IDs for a point*) We say that  $S \subseteq [d]$  is an ID for  $x \in T$  if  $x_S \neq y_S, \forall y \in T \setminus x$ . Let  $ID(x, T)$  be the smallest ID for  $x$  breaking ties arbitrarily. Let  $idLength(x, T) = |ID(x, T)|$ .

**Definition: 4** (*Imposters*) Given  $x \in T$  and  $S \subseteq [d]$ , the set of impostors of  $x$  in  $T$  are all points that equal  $x$  on all coordinates in  $S$ , i.e.,  $Imp(x, T, S) = \{y \in T \mid x_S = y_S\}$ .

**Definition: 5** (*Partial ID*)  $PID(x, T) = \arg \min_{S \subseteq [d]} (|S| + \log_2(|Imp(x, T, S)|))$ ,  
and  $pidLength(x, T) = \min_{S \subseteq [d]} (|S| + \log_2(|Imp(x, T, S)|))$ .

**Geometric view of `pidLength`:** A subcube  $C$  of  $\{0, 1\}^d$  is the set of points obtained by fixing some subset  $S \subseteq [d]$  coordinates to values in  $\{0, 1\}$ . The sparsity of  $T$  in a subcube  $C$  is  $\rho_{0,1}(T, C) = \frac{|C|}{|C \cap T|}$ . The notation  $C \ni x$  means that  $C$  contains  $x$ , hence  $\min_{C \ni x}$  is the minimum over all  $C$  that contain  $x$ . Anomalies are points that lie in relatively sparse subcubes. Low scores come with a natural witness: a sparse subcube  $PID(x, T)$  containing relatively few points from  $T$ .





# Chapter 4

## Contributions

The algorithm presented in previous chapter has  $O(n)$  time complexity. We further propose another distributed algorithm in this chapter based on xyz which has linear time complexity.

### 4.1 Construction

Write ...

### 4.2 Improved Method

Write...

### 4.3 Conclusion

In this chapter, we proposed another distributed algorithm for XYZ. This algorithm has both time complexity of  $O(n)$  where  $n$  is the total number of nodes. In next chapter, we conclude and discuss some of the future aspects.



# Chapter 5

## AI in Image Compression

This chapter contains information about work done before COVID-19 pandemic lockdown.

### 5.1 Construction

Write ...

### 5.2 Improved Method

Write...

### 5.3 Conclusion

In this chapter, we proposed another distributed algorithm for XYZ. This algorithm has both time complexity of  $O(n)$  where  $n$  is the total number of nodes. In next chapter, we conclude and discuss some of the future aspects.



# Chapter 6

## Conclusion and Future Work

write results of your thesis and future work.



# References

- [1] Pham H, *Handbook of Engineering Statistics*. Springer, London, 2006. [Online]. Available: <https://doi.org/10.1007/978-1-84628-288-1>
- [2] P. Gopalan, V. Sharan, and U. Wieder, “Pidforest: Anomaly detection via partial identification,” in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 15 809–15 819. [Online]. Available: <http://papers.nips.cc/paper/9710-pidforest-anomaly-detection-via-partial-identification.pdf>
- [3] F. T. Liu, K. M. Ting, and Z.-H. Zhou, “Isolation-based anomaly detection,” *ACM Trans. Knowl. Discov. Data*, vol. 6, no. 1, Mar. 2012. [Online]. Available: <https://doi.org/10.1145/2133360.2133363>
- [4] B. R. Preiss, *Data Structures and Algorithms with Object-Oriented Design Patterns in C++*. USA: John Wiley and Sons, Inc, 1998.