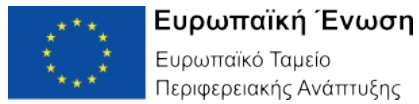




Partially supported by the Cyprus Research and Innovation Foundation under the grant agreement POST-DOC/0916/0090 and by EU H2020 Innovation Associate grant PREFAIL (957149)



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Ταμείο
Περιφερειακής Ανάπτυξης



Κυπριακή Δημοκρατία

ML Techniques for Detecting Storage Failures

Marko Krstic, Innovation Associate, Algolysis Ltd





Challenges for detection of Soon-to-Fail Drives

- High level of noise
- Extremely class imbalanced distributions
- Concept drift
- Large number and heterogeneity of hard drives in a dataset



Data processing

- Backblaze Dataset Format
- Feature selection
- Test/Train/Validation Dataset Creation
- Labeling



Backblaze Dataset Format

- **365 daily snapshots** with following columns:
- **Date:** The date of the file in yyyy-mm-dd format
- **Serial Number:** The manufacturer-assigned serial number of the drive.
- **Model:** The manufacturer-assigned model number of the drive.
- **Capacity:** The drive capacity in bytes
- **Failure:** Contains a “0” if the drive is healthy. Contains a “1” if this is the last day the drive was operational before failing.
- **Attributes:** 90 columns of S.M.A.R.T attributes and their normalized values, each associated with an identifier.

Feature Selection

Table 1. List of S.M.A.R.T attributes used.

ID	Name	Description
1	Read Error Rate	Frequency of errors during read operations.
3	Spin-Up Time	Time required a spindle to spin up to operational speed.
4	Start/Stop Count	Raw value holds the actual number of spin-up/spin-down cycles.
5	Reallocated Sectors Count	The number of the unused spare sectors. When encountering a read/write/check error, a device remaps a bad sector to a "healthy" one taken from a special reserve pool.
9	Power-On Hours Count	The Raw value shows the actual powered-on time, usually in hours.
12	Power Cycle Count	The Raw value holds the actual number of power cycles.
194	Temperature	Temperature, monitored by a sensor somewhere inside the drive. Raw value typically holds the actual temperature (hexadecimal) in its rightmost two digits.



Test/Train/Validation Dataset Creation

- We generate the testing dataset by selecting 30 healthy and 30 failed drives from each year at random. The remaining data is divided into training and validation dataset in proportion of 70% and 30%, respectively.
- With the sequence length parameter it is possible to control how many previous days S.M.A.R.T measurements are added to the current to describe operational state of the drive. In this paper we adopted sequence length of 30 days as long enough to model process of disk failing.



Labeling process

- **H**: a drive record is labeled as healthy at a date D , if the drive is not marked as failed in our dataset in any day d s.t., $D < d \leq D + 29$ or we have reached the end of the dataset (whichever is first).
- **F**: a drive record is labeled as soon-to-fail (or failed) at date D , if the drive is marked as failed in the dataset in any day d s.t., $D < d \leq D + 29$.



Prediction Models

- **LSTM** (Long Short Term Memory)
- Combination of **CNN**(Convolutional Neural Network) and **LSTM**
- **ResNet** (Residual Neural Network)
- **Simple CNN** model

LSTM and CNN-LSTM

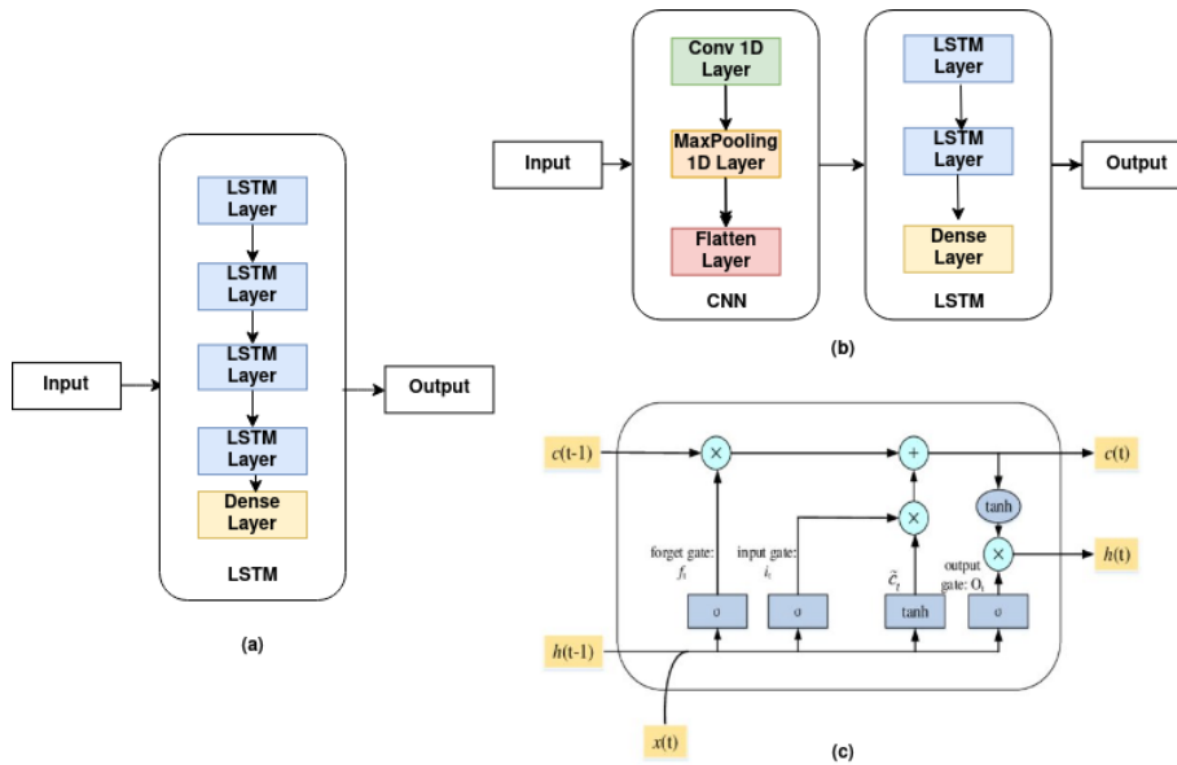


Fig. 2. a) LSTM model [9], (b) CNN LSTM model [9], and (c) LSTM Unit[30]

ResNet

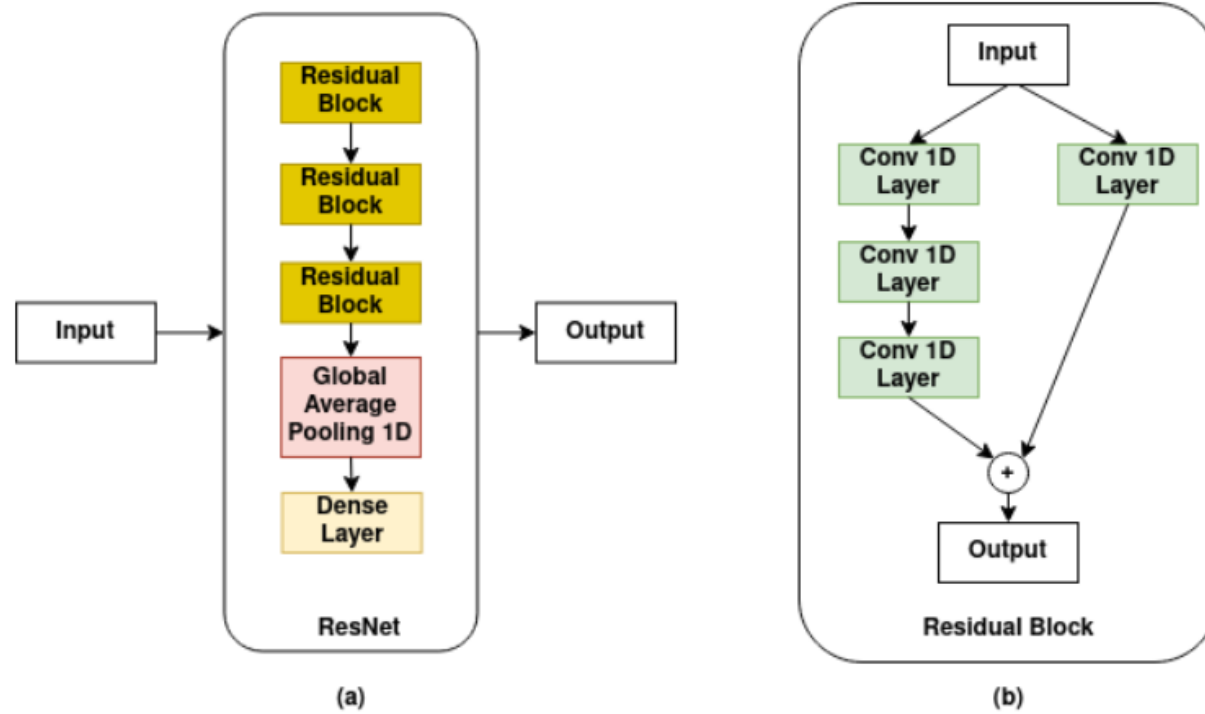


Fig. 3. a) Structure of ResNet model [38] b) Structure of residual block [1]

Simple CNN

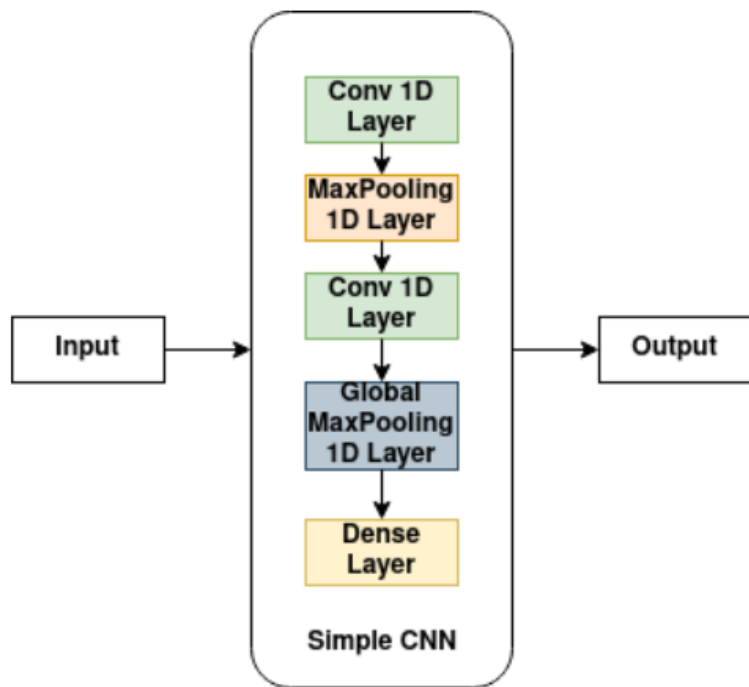


Fig. 4. Structure of Simple CNN model



Training Process

- Dealing with large amounts of data
- Dealing with class imbalance



- **Dealing with large amounts of data**

- We examined two approaches where each batch is formed by:
 - data collected for a single drive
 - combination of data from failed and healthy drives



• Dealing with class imbalance

- In the case where a batch is formed from data collected for a single drive, we apply both undersampling and cost-sensitive learning methods
- On the other hand, in the case when batch is formed from both healthy and failed drives data, we did not use undersampling. The reason for this decision was to support learning of differences between healthy and soon-to-fail operating characteristics in the same batch with as many examples. In addition we tried advanced method to fight with class imbalance in ML models with sequential inputs are also considered. Recent research showed that if cost/weights in training process is adjusted dynamically to class imbalance ratio in each batch, instead to global class imbalance ratio.



Experimental Results

- Comparison of different ways to generate data in training batches
- Performance of advanced class imbalance technique for time series



- **Comparison of different ways to generate data in training batches**

Table 2. Performance of models in case when batches are formed from one disk data.

Model	Precision	Recall	F1 measure	Accuracy
LSTM	0.45	1	0.62	0.55
CNN LSTM	0.45	1	0.62	0.55
ResNet	0.57	0.23	0.32	0.57
Simple CNN	0.56	0.8	0.66	0.6

Table 3. Performance of models in case when batches are formed from mixture of failed and non-failed drives.

Model	Precision	Recall	F1 measure	Accuracy
LSTM	0.57	0.88	0.69	0.61
CNN LSTM	0.57	0.88	0.69	0.61
ResNet	0.36	0.6	0.45	0.55
Simple CNN	0.32	0.66	0.43	0.52



Performance of advanced class imbalance technique for time series

Table 4. Performance of models in case when different class imbalance weights/costs are incorporated at batch level.

Model	Precision	Recall	F1 measure	Accuracy
LSTM	0.51	0.79	0.61	0.55
CNN LSTM	0.51	0.79	0.61	0.55
ResNet	0.34	0.54	0.37	0.51
Simple CNN	0.3	0.58	0.4	0.53



Conclusions about sequential modeling

- The proposed detection system is designed in such a way to overcome challenges that arise from data characteristics, with the main focus on Big data and class imbalance implication. It is experimentally shown that **LSTM** neural network, trained in batch-by-batch manner, where each batch are formed from both information about healthy and failed drives, is capable to recognize 88% of drives that are in soon-to-fail conditions. However, this comes at price of 43% of false alarms, thus further research is needed to improve prediction precision.



Additional experiments with non-sequential modeling

ID	Name	Description
1	Read Error Rate	Frequency of errors during read operations.
5	Reallocated Sectors Count	The number of the unused spare sectors. When encountering a read/write/check error, a device remaps a bad sector to a "healthy" one taken from a special reserve pool.
7	Seek Error Rate	Frequency of the errors during disk head positioning.
12	Power Cycle Count	The Raw value holds the actual number of power cycles.
192	Power-Off Retract Cycles	The number of unexpected power outages when the power was lost before a command to turn off the disk is received. On a hard drive, the lifetime with respect to such shutdowns is much less than in case of normal shutdown. On an SSD, there is a risk of losing the internal state table when an unexpected shutdown occurs.
197	Current Pending Sectors	The number of unstable sectors which are waiting to be re-tested and possibly remapped.
198	Off-line Uncorrectable	The number of bad sectors which were detected during offline scan of a disk.]

SMART attributes value for current day
+
difference between current and previous day

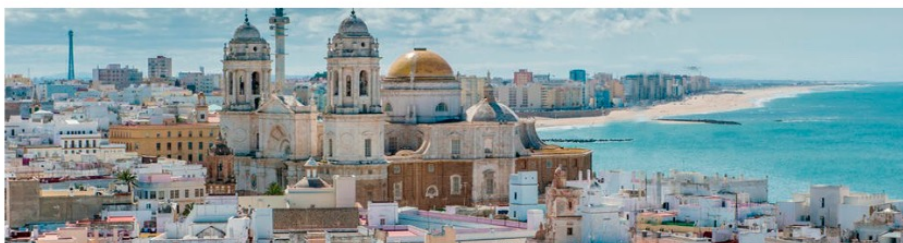
Model	Precision	Recall	F1 measure	Accuracy
Feed forward neural network	0.24	0.23	0.24	0.62
Bayesian neural network	0.29	0.61	0.39	0.52

Potential of AIOps from Digital Forensics Perspective

*7th WORKING GROUP MEETING
DIGFORASP COST ACTION 17124
Centro de Transferencia “El Olivillo”*

Cádiz, Spain

October 21st – 23rd, 2021



12:10–13:30	Workshop. Session 2	
12:10–12:30		<i>Detecting stylistic features to identify hate-speech spreaders on social media</i> Antonio Pascucci
12:30–12:50		<i>Potential of AIOps from Digital Forensics Perspective</i> Marko Krstic, Nicolas Nicolaou, Efstathios Stavrakis
12:50–13:10		<i>Secret-Key Agreement by Asynchronous EEG over Authenticated Public Channels</i> M. Galis, M. Milosavljević, A. Jevremović, Z. Banjac, A. Makarov, J. Radomirović.
13:10–13:30		<i>Artificial Intelligence and Law Enforcement: Main applications and Impact on Fundamental Rights</i> José María Castillo-Secilla

DIGFORASP



Thank you for your attention!