# Using AI to Predict Server **Hard Drive Failure**

Project by Samy Djemaï
LOG6309E
Teacher: Heng Li

# What's the **Problem?**

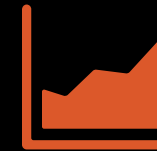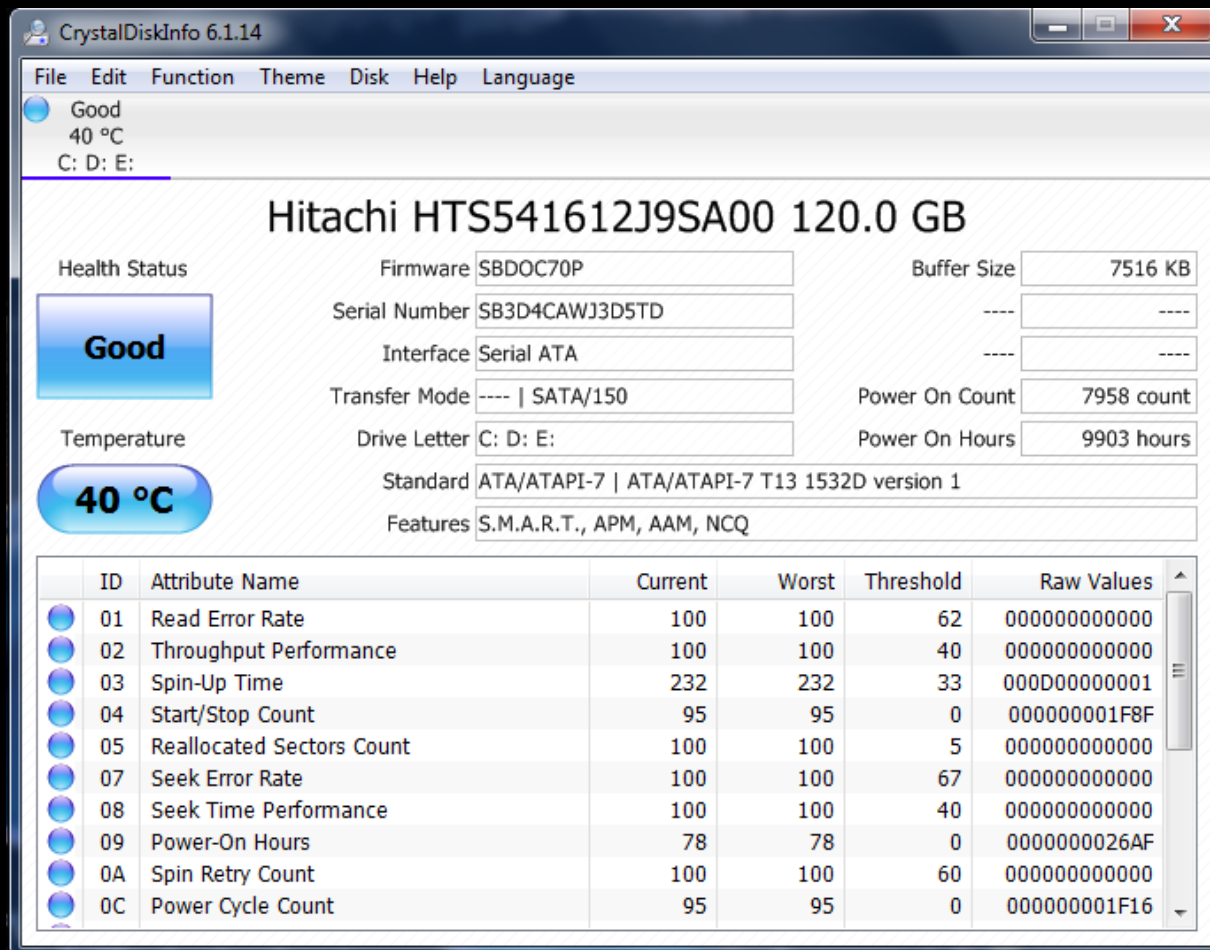- A datacenter contains servers which contain hard drives, used to store application data

- For 98% of companies, one hour of downtime costs **over $150,000**

  Source: ITIC 2017-2018 Global Server Hardware, Server OS Reliability Survey

- Server operators rely on arbitrary rules or wait for failure to replace drives

- **What if we could predict hard drive failures?**

# A S.M.A.R.T.er Solution



Modern hard drives provide S.M.A.R.T. stats

Self-monitoring

Train and test predictive models

# So Many Questions...

1 What are the most important features to consider?

2 How accurate can classifiers get?

3 Can models be ported to other drives?

# Backblaze's Dataset

Backup company provides S.M.A.R.T. and failure stats for its hard drives, from Q1 2019 to Q3 2021

Daily snapshots of >175,000 disks' stats, 131 columns of data

| date | serial_number | model | capacity_bytes | failure | smart_1_normalized | smart_1_raw | ... |
|---|---|---|---|---|---|---|---|
| 2021-04-01 | ZHZ65F2W | ST1200NM0008 | 12,000,138,625,024 | 0 | 82 | 159,565,280 | ... |
| 2021-04-01 | ZLW0EGC7 | ST12000NM001G | 12,000,138,625,024 | 0 | 74 | 22,618,672 | ... |
| 2021-04-01 | ZA1FLE1P | ST8000NM0055 | 8,001,563,222,016 | 0 | 82 | 167,665,584 | ... |
| ... | ... | ... | ... | ... | ... | ... | ... |

# The **Stuff** We Used

**12-core Intel Core i7-9750H, 2.60GHz**

16 GB RAM

Pop_OS! 21.04

**Anaconda 4.10.3**

Python 3.8.12

NumPy, Pandas, Ruptures, Scikit-learn, Multiprocessing

**Visual Studio Code**

Jupyter Notebook

# Which Models Should We Study?

Q1 2020 Dataset

| | count | unique | fail_count | failure_rate | missing_stats |
|---|---|---|---|---|---|
| ST8000DM004 | 209 | 2 | 1 | 0.4785% | 4 |
| TOSHIBA MQ01ABF050 | 39,102 | 413 | 39 | 0.0997% | 9 |
| ... | ... | ... | ... | ... | ... |
| **ST4000DM000** | 1,744,529 | 19,142 | 68 | 0.0039% | 5 |
| ST12000NM0008 | 750,681 | 10,876 | 29 | 0.0039% | 6 |
| **ST12000NM0007** | 3,368,588 | 36,997 | 126 | 0.0037% | 6 |

# ST4000DM000 It Is!

- One of the **most used drives**

- **High failure rate**

- Reports **nearly all** **S.M.A.R.T. stats**

# No Stats For You!

`df.isnull().sum().sort_values(ascending=False).head(72)`



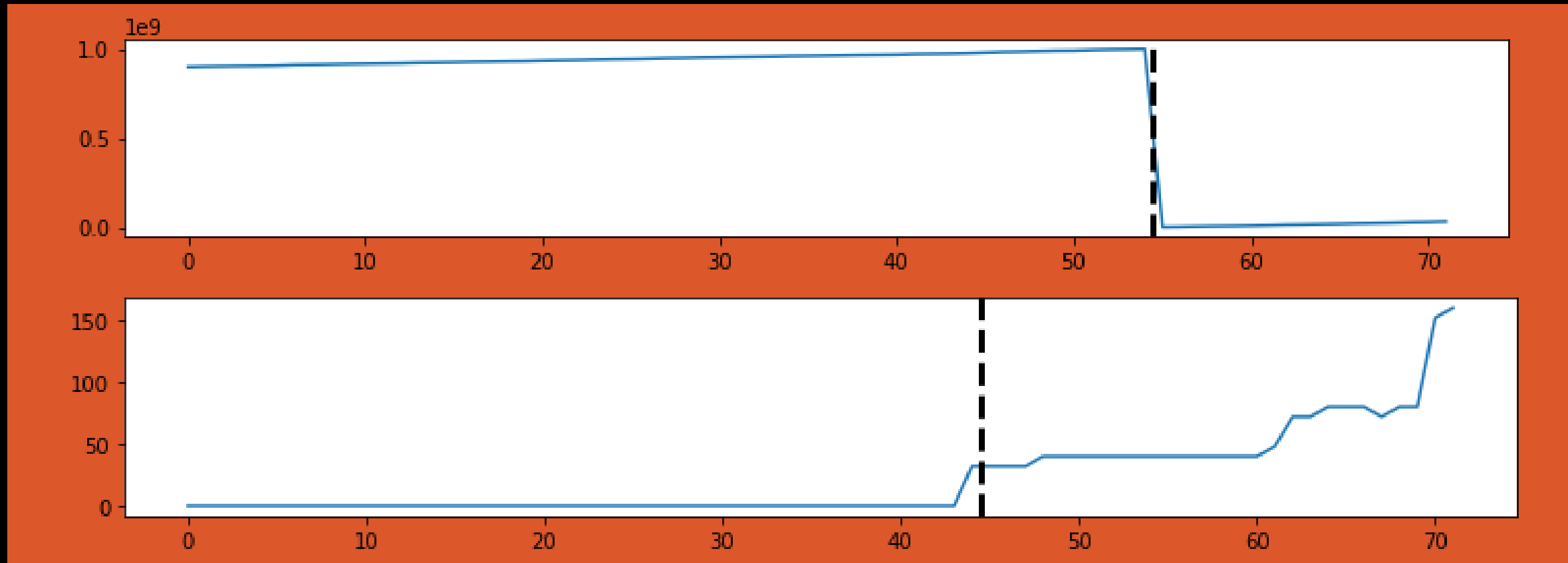| | |
|---|---|
| smart_255_raw | 132,339 |
| smart_250_normalized | 132,339 |
| smart_15_raw | 132,339 |
| smart_15_normalized | 132,339 |
| ... | ... |
| smart_183_raw | 112,270 |
| smart_8_raw | 94,583 |
| smart_8_normalized | 94,583 |



Q1_2020 Dataset (132,339 rows)

# The Quest for **Relevancy**

**RQ1**: Which S.M.A.R.T. stats must be considered?

- Dataset: **Q1 2019** to **Q4 2020**

- **615** failed ST4000DM000 drives

- Up to **60 days** of data before failure

# Change Is Now

*ruptures: change point detection in Python*, Truong et al., 2018



Seek Error Rate (SMART_7_RAW), Off-line Uncorrectable (SMART_198_RAW)

# The Results (Number 4 Will Shock You)

Analysis ran on 615 failed drives, from Q1 2019 to Q4 2020 (last 60 days)

| Name | Description | Frequency |
|---|---|---|
| smart_242_raw | Total LBAs Read | 43.74% |
| smart_9_normalized | Power-On Hours Count (Norm.) | 42.44% |
| smart_241_raw | Total LBAs Written | 42.28% |
| **smart_7_raw** | Seek Error Rate | 33.66% |
| **smart_7_normalized** | Seek Error Rate (Norm.) | 25.69% |
| smart_193_raw | Load/Unload Cycles | 21.46% |
| **smart_187_{raw,normalized}** | Reported Uncorrectable Errors | 15.12% |
| **smart_197_raw** | Current Pending Sectors | 15.12% |
| **smart_198_raw** | Off-line Uncorrectable | 15.12% |
| ...and 19 more | including Reallocated Sectors | >1.00% |

# Be Accur8 M8
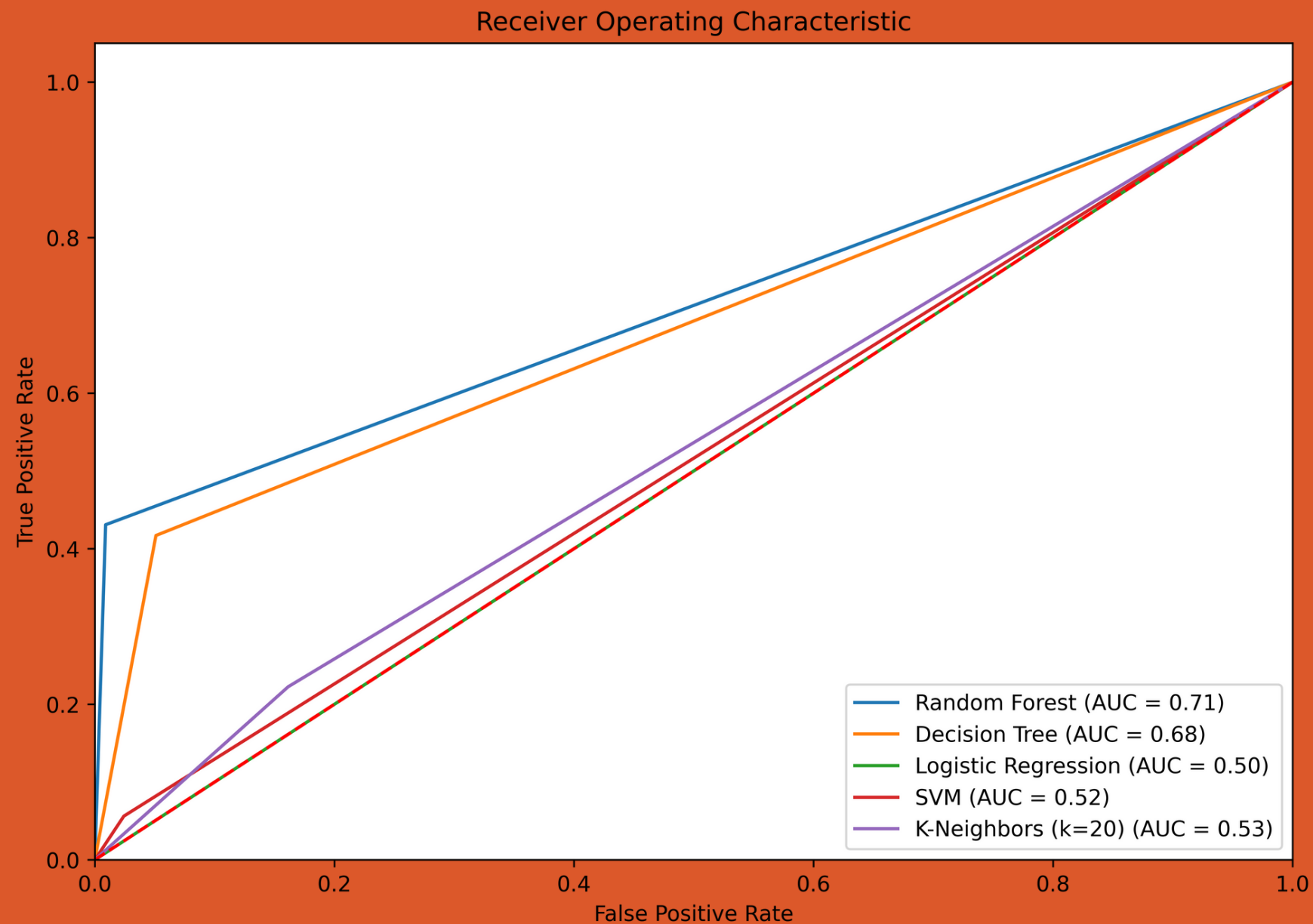
RQ2: How accurate can models get?

- Training set: Q1 2019 to Q4 2020

- Testing set: Q1 & Q2 2021

- Dataset processing: **backtracking**
  failure = 1 up to 15 days before actual failure

- Low failed drive count: **sampling**
  1/3 failed drives, 2/3 healthy drives

# The **Accuracy** of this... Classifier

**RQ2**: How accurate can models get?

| | Accuracy | Recall | F1-Score |
|---|---|---|---|
| **Random Forest** | **0.8138** | **0.4309** | **0.5942** |
| Decision Tree | 0.7807 | 0.4170 | 0.5461 |
| Logistic Regression | 0.6836 | 0.0000 | 0.0000 |
| SVM | 0.6848 | 0.0562 | 0.1014 |
| K-Neighbors (k = 20) | 0.6436 | 0.2225 | 0.2832 |
| K-Neighbors (k = 2) | 0.6513 | 0.1114 | 0.1682 |

# Takeaways

Thank you for listening :)

## Discussions

- Possible to predict disk failures using common classifiers
- Many factors are at play
- Slight clustering of failed disks, as shown by K-Neighbors
- Low recall and F1-Score overall

## Limitations

- Analysis limited by chosen S.M.A.R.T. attributes
- Single drive model
- No real time-series approach

## Next steps

- Answer RQ3: how do these classifiers fare on other drives?
- Try more classifiers/models
- Try other drives
- Expand dataset to Q3 2021