

Predicting Hard Drive Failure

by Sulman Khan

4/30/2020

Abstract

Data is an integral component of our society. From the simple caloric deficits collected in your apple watch to the user history in your Netflix account, data is used in a myriad of applications. With such an abundance of data being used daily, how is it stored? The solution is computer backup or cloud storage services. Furthermore, Backblaze is a world leader in computer backup and storage. Since 2013, Backblaze has published statistics and insights based on the hard drives in their data center [1]. In this study, we'll explore various features in a hard drive dataset to predict hard drive failure.

Dataset

The original dataset is separated into four folders depicting the four quarters of the year. Each quarter contains daily CSV files, which has information on all the hard drives in operation. When a hard drive fails, it is removed from the proceeding day. There are roughly 130 features in the dataset, and most are SMART stats, which are industry known statistics used by various hard drive manufacturers for quality assurance. Despite having access to all SMART stats, these can result in an abundant amount of null entries due to company bias - a company may record only SMART stats, they deem necessary. Preprocessing is split into two distinct phases: EDA and Prediction.

1. Create EDA dataframes
 - a. Merge daily CSV's into quarterly CSV's (4 CSV)
 - b. Merge quarterly CSV's into a yearly CSV (1 CSV)
 - c. Track the frequency of hard drive failures (1 CSV)
2. Create Prediction dataframe
 - a. Preprocess the yearly CSV file (1 CSV)
 - i. Remove null entries
 - ii. Scale columns by feature

If you want more information on the functions used in the creation of these CSV files, check out the Backblaze-parser script and Jupyter Notebook.

Exploratory Data Analysis

The EDA explored features such as drive size, manufacturer, and age because these are SMART statistics that most manufacturers record.

Quarterly Frequencies

The daily and quarterly hard drive failure frequency was explored in Figure 1. Across the four quarters, there is an increasing trend in both the number of hard drives and failures. The maximum number of hard drives and failures is 125k and 678, respectively. The minimum number of hard drives and failures is 108K and 444, respectively. There is an outlier visible in November - many hard drives were taken down for maintenance and were not recorded.

Daily and Quarterly Hard Drive and Failure Frequency

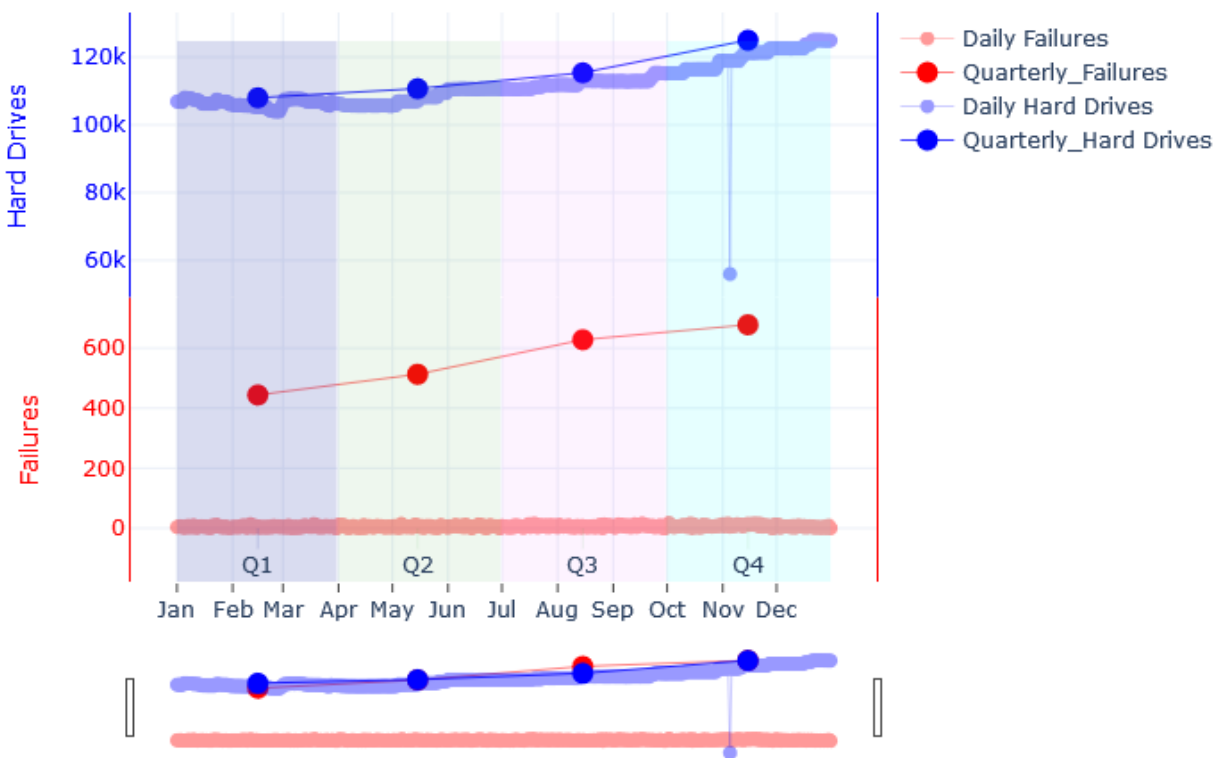
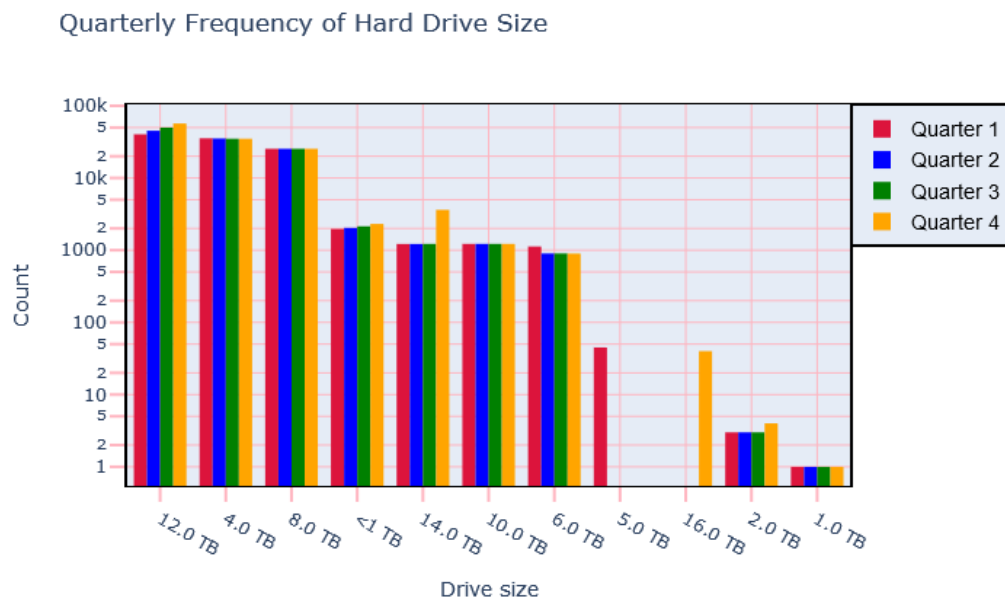


Figure 1: The daily and quarterly hard drive failure frequency.

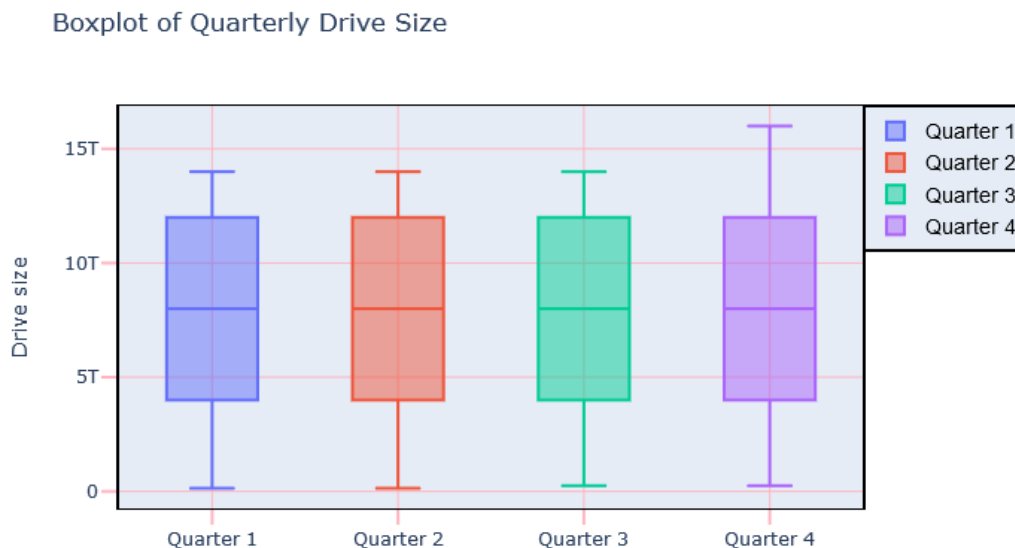
Drive Size

Various data visualizations were used to explore the hard drive size in Figure 2. Across the four quarters, there is a trend in using larger hard drive sizes such as 12 TB, 14 TB, and 16 TB. The size of the hard drive's ranges from < 1 TB to 16 TB. As of quarter 4, 16 TB hard drives were added in datacenters. The hard drive size with the largest failure rate is the <1 TB hard drives - a failure rate of 7.668% and a count of 190 hard drives. The hard drive size with the greatest failure frequency is the 12 TB hard drives - a failure rate of 2.09% and a count of 1200 hard drives.

(a) Quarterly Frequency of Hard Drive Size



(b) Boxplot of Quarterly Hard Drive Size



(c) Frequency of Hard Drive Size Failures

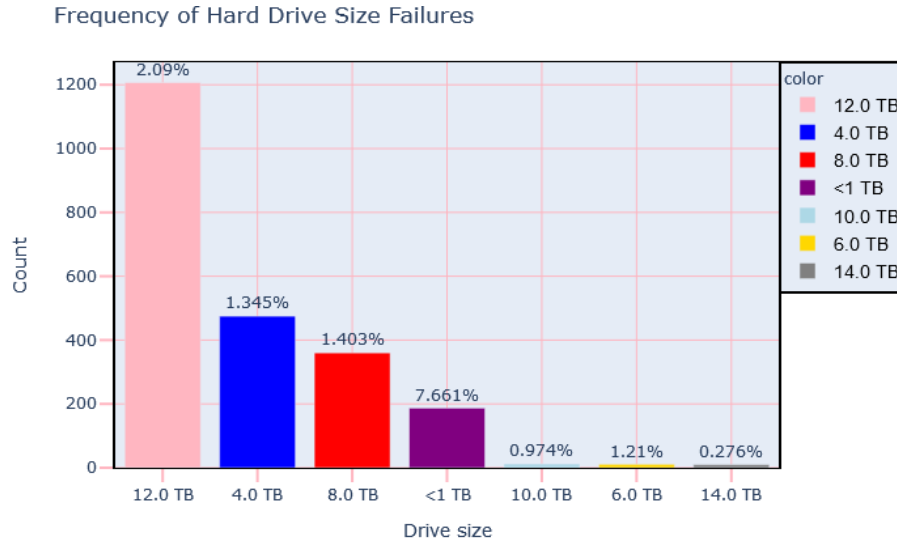
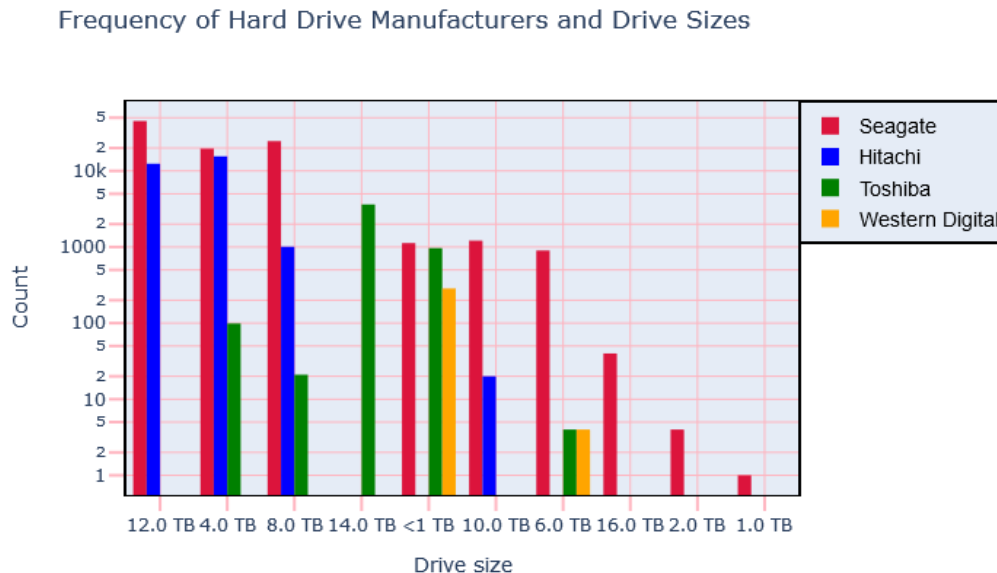


Figure 2: Feature Exploration of Hard Drive Size. (a) Quarterly Frequency of Hard Drive Size. (b) Boxplot of Quarterly Hard Drive Size. (c) Frequency of Hard Drive Size Failures.

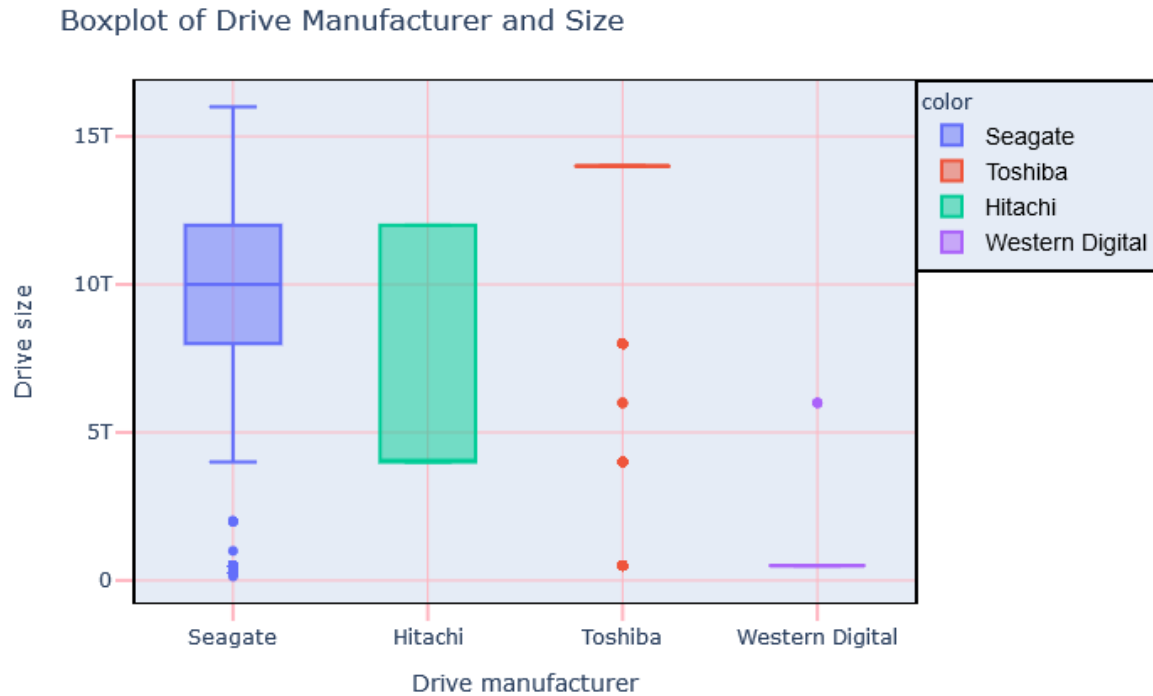
Drive Manufacturers

Various data visualizations were used to explore the hard drive manufacturers in Figure 3. Across the four manufacturers, the frequency of hard drives from highest to lowest: Seagate, Hitachi, Toshiba, and Western Digital. Seagate is the only drive manufacturer that has a complete boxplot. The greatest failure rate is 4.884% from Western Digital with a count of ~20 hard drives. The largest failure count is the Seagate hard drives with a failure rate of 2.172% and a count of 2000 hard drives.

(a) Frequency of Hard Drive Sizes with Manufacturers



(b) Boxplot of Hard Drive Manufacturer and Size



(c) Frequency of Hard Drive Manufacturer Failures

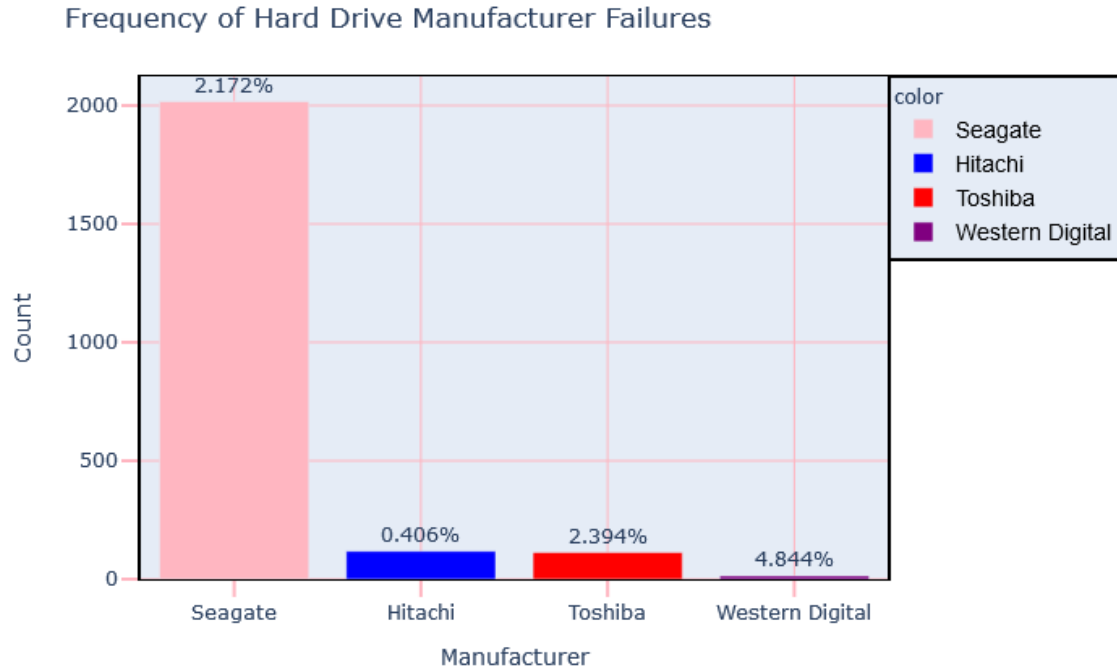
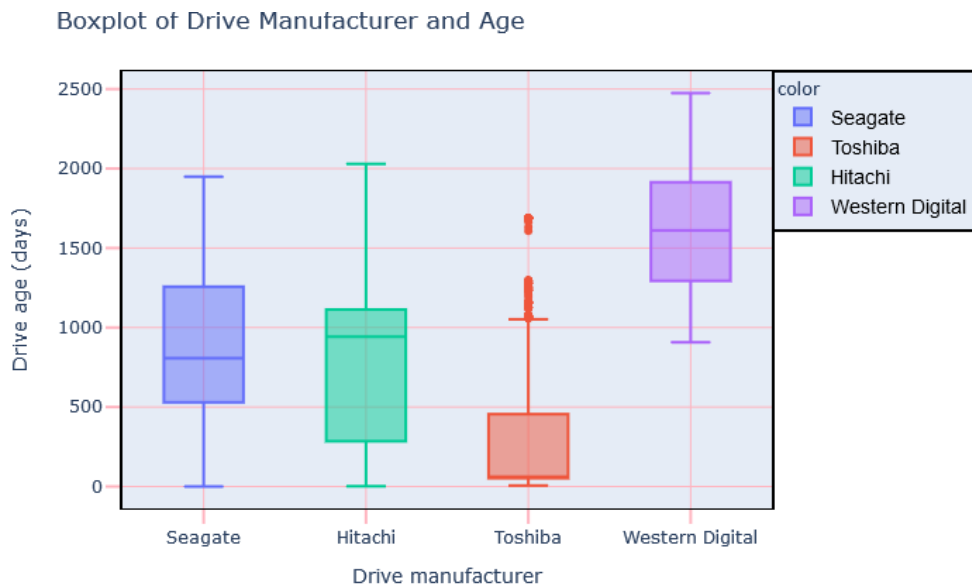


Figure 3: Feature Exploration of Hard Drive Manufacturer. (a) Frequency of Hard Drive Sizes with Manufacturers. (b) Boxplot of Hard Drive Manufacturers and Sizes. (c) Frequency of Hard Drive Manufacturer Failures.

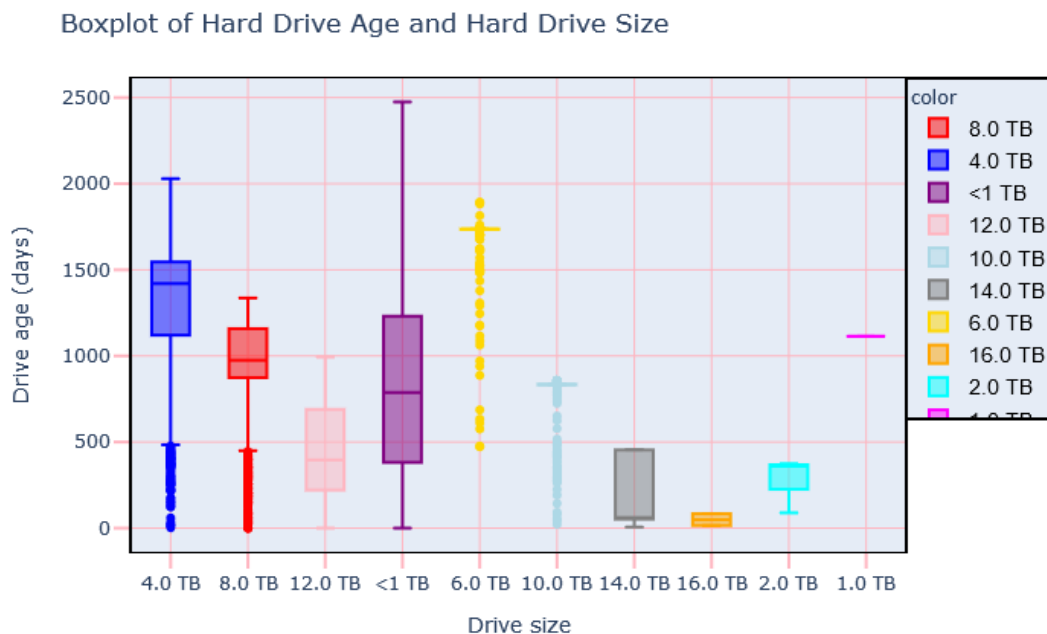
Drive Age

Various data visualizations were used to explore the hard drive age feature in Figure 4. The hard drive age ranges from 0 – 2500 days. Most of the older hard drives are Western Digital, and the recent hard drives are Toshiba. The greatest failure rate is 6.249%, given by the 400-599 days interval. The highest failure count is the 1000+ days interval with a failure rate of 1.169%, and a count of 550. As a drive manufacturer, older drive ages or lifetimes before failure are ideal. Smaller lifetimes in certain models could be attributed to defects or poor design.

(a) Boxplot of Drive Manufacturer and Age



(b) Boxplot of Hard Drive Age and Hard Drive Sizes



(c) Frequency of Hard Drive Age Failures

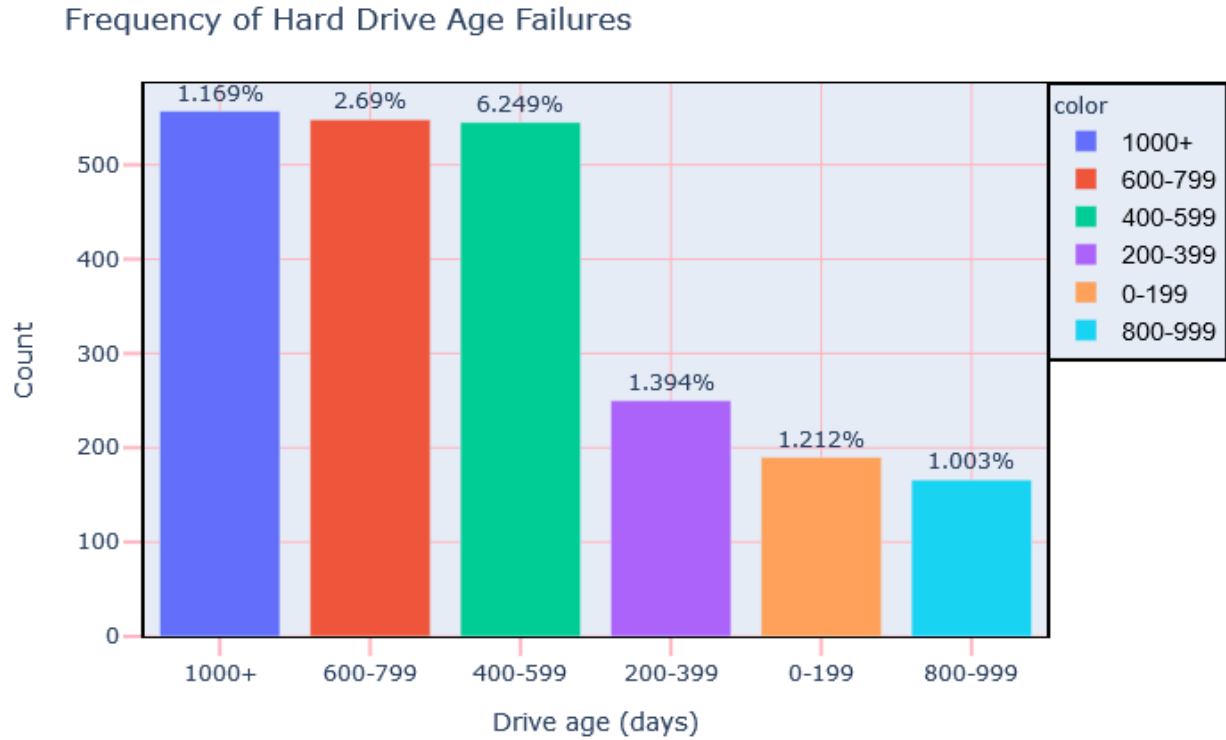


Figure 4: Feature Exploration of Hard Drive Age. (a) Boxplot of Drive Manufacturers and Age. (b) Boxplot of Hard Drive Age and Hard Drive Sizes. (c) Frequency of Hard Drive Age Failures.

Prediction

In the prediction stage, our problem is identifying additional features to predict failures in hard drives. Backblaze has selected five features that industry experts have stated correlates to hard drive failure. Therefore, we have a baseline model that includes all these five features and we'll compare that to our final model.

In order to characterize the performance, we'll be using the Receiver Operating Characteristics, mainly the AUC as a metric. Our baseline algorithm for our predictions is Logistic Regression. Our process will be selecting pertinent SMART statistics provided in the dataset by researching their functions and then reporting their performance in our model [2].

Procedure

The samples are split into an 80/20 ratio between the training and test samples, respectively. There was k-fold cross-validation performed on the training set for hyperparameter tuning. In k-fold cross-validation, the data is evenly split into k non-overlapping pieces and each piece contains a validation set. Then, average the selected scoring metric across the k trials. Cross-validation was used to test the generalizability of the model. As CV checks in how it is performing on new unseen data with a limited amount of training samples. Since the problem is binary classification, stratification can be used to make sure the same number of proportion of classes are in the validation and training folds.

Smart Statistics

The following section explains the various features selected in the model, the last two features Smart 199 and Smart 242 displayed promising results in predicting hard drive failure [3].

Smart 5: Reallocated Sector Count

When the hard drive finds a read/write/verification error, it marks this sector as "reallocated" and transfers data to a special reserved area (spare area). This process is also known as remapping and "reallocated" sectors are called remaps. Therefore, on a modern hard disk, you will not see "bad blocks" while testing the surface - all bad blocks are hidden in reallocated sectors.

Smart 187: Reported Uncorrectable Errors

The number of errors that could not be recovered using hardware ECC (error-correcting code).

Smart 188: Command Timeout

The number of aborted operations due to hard disk timeout.

Smart 197: Current Pending Sector Count

The current count of unstable sectors (waiting for remapping). The raw value of this attribute indicates the total number of sectors waiting for remapping. Later, when some of these sectors are read successfully, the value is decreased. If errors still occur when reading some sector, the hard drive will try to restore the data, transfer it to the reserved disk area (spare area) and mark this sector as remapped.

Smart 198: Offline Uncorrectable Errors

The number of errors that the drive has attempted to correct itself but failed.

Smart 199: Ultra Direct Memory Access Cyclic redundancy check (UDMA CRC) Errors

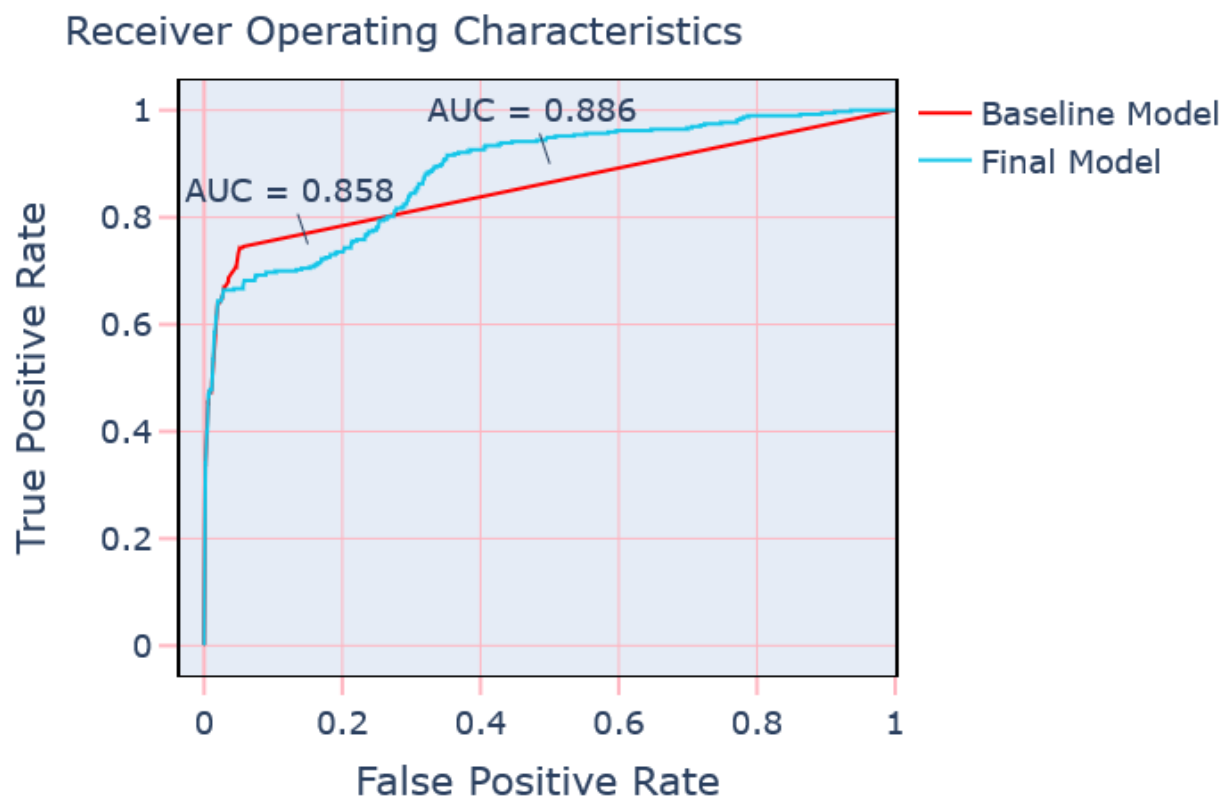
The total quantity of CRC errors during UltraDMA mode. The raw value of this attribute indicates the number of errors found during data transfer in UltraDMA mode by ICRC (Interface CRC).

Smart 242: Total Logical Block Addressing (LBA) Read

The total number of LBAs read.

ROC Curve

Lastly, the ROC curve was plotted for the baseline and final models. The AUC score increased from 85.8% to 88.6% with the added features.



Conclusion

The project was a success, as we have investigated two additional features in predicting hard drive failure. The AUC score increased from 85.8% to 88.6% with the added features. As most of the difficulty in this project can be contributed by the preprocessing steps – the data was very dirty and needed creative solutions to create many of the data visualizations. Possible future projects include survival analysis on hard drives (predicting the probability a hard drive will survive with a given lifetime) and a regression problem involving predicting the lifetime of a hard drive given a set of SMART statistics.

References

- [1] "Backblaze Hard Drive Stats", *Backblaze.com*, 2020. [Online]. Available: <https://www.backblaze.com/b2/hard-drive-test-data.html>. [Accessed: 27- Apr- 2020].
- [2] "SMART Drive and Failure Rates", *Backblaze.com*, 2020. [Online]. Available: <https://www.backblaze.com/blog-smart-stats-2014-8.html>. [Accessed: 29- Apr- 2020].
- [3] "Backup Software & Data Protection Solutions - Acronis", *Acronis.com*, 2020. [Online]. Available: <https://www.acronis.com/en-us/>. [Accessed: 01- May- 2020].