

UNI VARIATE LINEAR REGRESSION ANALYSIS FOR SONATYPE NEXUS REPOSITORY SPACE MANAGEMENT**Renukadevi Chuppala¹, Dr. B.Purnachandra Rao²**¹Western Union Financial Services, CA, USA renu.chuppala@gmail.com²Sr. Solutions Architect, HCL Technologies, Bangalore, Karnataka, India.
pcr.bobbepalli@gmail.com**Abstract**

Sonatype Nexus is a powerful repository manager widely used in DevOps and continuous integration/continuous delivery (CI/CD) pipelines. It is designed to manage, store, and retrieve binary artifacts efficiently. Nexus plays a crucial role in software development by providing centralized storage and management for build artifacts, dependencies, and containers, enabling teams to collaborate more effectively and maintain control over their software supply chain. As part of the business activity on daily basis, lots of artifacts will be uploaded to number of repositories to each nexus instance. The space will be consumed proportional to the volume of artifacts. It is administrator's responsibility to clear the space as and when it reaches to beyond the thresh hold limit. If the admin misses to clear the space it leads to the situation to shuts down the server, which will effect the entire business. Admin can use the available automated tasks from the nexus admin tasks, but they are having some limitations on the deletion of artifacts, we need to opt only unwanted artifacts. But in the case of high usage of wanted artifacts the space will get decreased and we can't know when the server will go down. This paper resolves this issue by providing the solution using machine learning algorithm (Linear Regression Analysis) on the usage of space by users. It will findout the regression equation for the given data, so that we can find out prediction value for each actual value. Using the regression equation we can predict the value.

Keywords: *Sonatype Nexus Repository Manager, NXRM, Release Repository, Snapshot Repository, Docker registry, npm repository, Maven, Nuget, LDAP, Univariate Linear Regression, Feature Engineering, Linear Regression Analysis.*

INTRODUCTION

Nexus Repository is a centralized platform used to store and manage software artifacts [1], such as libraries, packages, and dependencies, in various formats. It acts as a universal repository manager that supports formats like Maven, npm, NuGet, Docker [2], and more, allowing teams to store, retrieve, and share these artifacts during the software development lifecycle. Nexus [3] stores build artifacts like JAR files, Docker images, etc., making them accessible to development and deployment pipelines. It can proxy external repositories (e.g., Maven Central) and cache dependencies locally, which reduces build times and network usage. Allows teams to store multiple versions of artifacts and manage them through consistent versioning policies. Nexus can act as a private Docker registry [4] for storing and managing container images, ensuring secure access to containerized applications.

Nexus integrates seamlessly with continuous integration and continuous deployment tools like Jenkins [5] and Bamboo, automating the retrieval and publication of build artifacts. It offers fine-grained role-based access control (RBAC) and can scan repositories for vulnerabilities, ensuring compliance with security policies. Nexus is widely used in DevOps pipelines to manage dependencies [6], store build outputs, and ensure a secure software supply chain. Each business unit will have one nexus instance and the users will start uploading artifacts to corresponding repository. There are snapshot and release repositories. All non prod environment artifacts will be created as snapshot artifacts and the prod artifacts will be created as non snapshot artifacts called release artifacts. ReactJS and Node JS artifacts will be stored in npm hosted repositories.

LITERATURE REVIEW**Sonatype Nexus Repository:**

Nexus instances are connected to load balancer as shown in the Fig 1. There are different types of clients available to work with nexus repository manager (NXRM). Based on the dev team request NXRM admin [7] will create the hosted

repositories for the respective type of deliveries. For Java development hosted repositories SNAPSHOT , RELEASE will be created, for nodes js npm repositories will be created. To upload the docker images docker registries will be created.

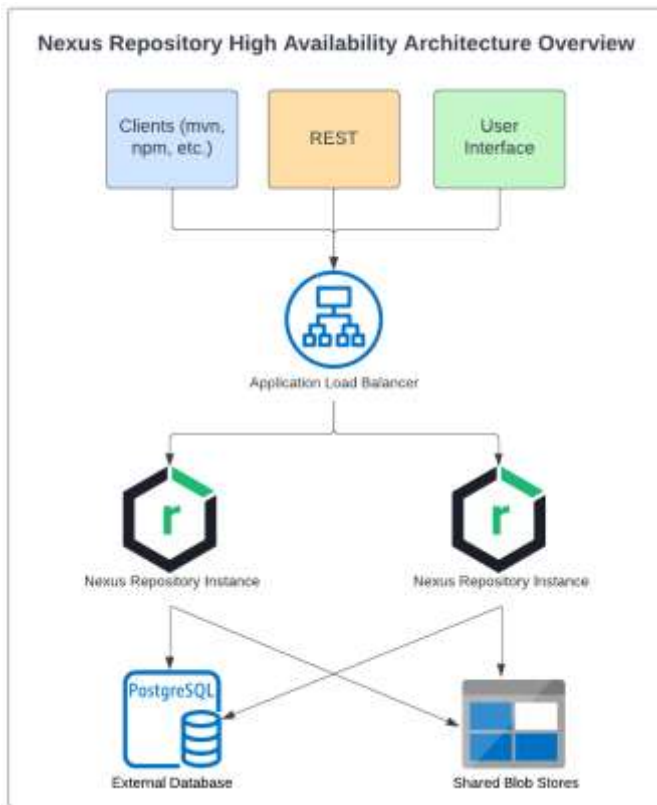


Fig: 1 Sonatype Nexus Repository Management Architecture diagram

The Fig: 1 shows the two instances backed to one load balancer and maintaining the consistency of the data using blob store and database postgres SQL [8]. NXRM admin will provide access to respective users. If there are number of users for the same repository he will create Active Directory (AD) Groups and add the number of users to the same AD Groups so that he no need to add each user to the repository for access. The AD group [9] will be given with the corresponding permissions. Java development team will create artifact type as jar file. This will be uploaded to snapshot or release repository based on the version of the artifact. If it is fixed version it will be uploaded to release repository , if it is snapshot version it will be uploaded to snapshot repository. For web development npm packages will be created and these will be uploaded to npm hosted repository. Dot Net development packages will be uploaded to Nuget repository. Docker images will be uploaded to Docker registry.

Artifact Storage:

On daily basis artifacts from the different developers will be uploaded to corresponding repositories. The access has been given by the NXRM admin team to work with the uploads. We can not re upload if it is NON SNAPSHOT artifact. If it is SNAPSHOT version it will be re uploaded.

Snapshot Repository:

In the development environment the code needs to be corrected frequently , so we can not finalize the artifact in the initial cycles. For each delivery or the completion of cycle the artifact will be uploaded. So in this case we need to use SNAPSHOT (Ex: 1.0.0-SNAPSHOT) version so that for each upload of the same artifact the time stamp will be appended to the artifact. For example if we are creating the artifact on the date 10/17/2024 21:03 then the artifact name is abc-1.0.0-171020242103. In this case abc is the name of the artifact. If we do the same build again after 5 minutes then the artifact

name is 1.0.0-171020242603. I have assumed the same minutes for easy explanation. In this case for each build it will create new artifact name. So the rebuild is happening without any issue. Snapshot repository allows rebuild operation. But this is technically wrong statement, why because we are not uploading same artifact each time. Every time we are uploading different artifact.

Release Repository:

Once all cycles are completed in dev environment i.e If we feel that the build is working fine and we are ready to create the artifact , we can change the version of the artifact as fixed version. In this case it doesn't matter about the date of working. The name of the artifact will be abc-1.0.0. No time stamp will be added to artifact name. In this case re upload is not allowed unless and until we need to enable the option to re upload the same artifact. Next release the artifact name will be abc-1.0.1. This will come from the snapshot artifact abc-1.0.1-SNAPSHOT [10].

Repository Proxying:

Different teams are having different repositories in the Nexus instance. Suppose team A wants to use the artifacts of team B then they can use the repositories inside Maven settings.xml [11] file so that these repositories will be scanned based on the dependencies mentioned at pom.xml file. We need to mention groupId , artifactId and version at pom.xml if we want to use any artifact. If the artifact is internal to organization then it will be downloaded from the repos mentioned at corresponding tags in settings.xml file. If the artifact is available at outside of the organization then it will throw the error if the corresponding location is not mentioned at settings.xml file. We need to mention the external urls as proxies inside settings.xml file. Maven will look into local cache for artifact. If it is not available at local cache It will connect with urls internal to organization. If it is not available at internal , then it will connect to external proxies to get the artifacts to local. As soon as it finds the artifacts either at local organization or external to organization it will be downloaded to local cache for future reference and it will be used in the subsequent builds.

Version Control:

Suppose the current artifact name is nexus while creating the development branch in the github the version is nexus-1.0.0-SNAPSHOT. It will be uploaded to nexus SNAPSHOT repository as nexus-1.0.0-171020241136. If we do next build then the artifact name is nexus-1.0.0-171020241138 (this buld is after 2 mins). Once we are done with all these cycles if we fix the version by removing the SNAPSHOT then the build version is nexus-1.0.0 and this will be uploaded to releases repository. Once we are done with this artifact if any one in the organization refers the GAV [12]of this artifact , this will be downloaded to their local from this nexus release repository. For the next cycle the new git hub branch will created with the version at development as nexus-1.0.1-SNAPSHOT. For eachj development build time stamped artifact will be created like nexus-1.0.1-171020241142 , nexus-1.0.0-171020241145 etc. The release artifact nexus-1.0.1 will be created at release build and it will be uploaded to nexus release repository. We can refer any of these versions using the GAV (GroupId, artifactId and version) parameters at pom.xml file.

Private Docker Registry:

Nexus can act as a private Docker registry for storing and managing container images, ensuring secure access to containerized applications. We need to create docker registry [13] at nexus instance as we have created like snapshot and release repositories, and login to this docker registry from the terminal and start pushing the docker images to same registry. We can pull the images back to our location based on our deployments.

Integration with CI/CD Pipelines:

Nexus integrates seamlessly with continuous integration [14] and continuous deployment tools like Jenkins and Bamboo, automating the retrieval and publication of build artifacts. Once the build starts by getting the github code into local , the artifact will get created. This needs to be uploaded to Nexus repository using the Nexus plugin which is already integrated to Jenkins. This should be authenticated using nexus credentials. Once the development team commit the code it will automatically trigger the job to run for build so that the code will be checked out to the base of the Jenkins (CI/CD tool). Once the build is successful , this process will create the artifact based on the type of the pom.xml file i.e it might be jar , war or zip .

This will be saved at the local cache of the machine where the build is happening. As part of the next phase of the build the artifact will be uploaded to Nexus repository. If the build created the snapshot artifact, then it will be uploaded to snapshot repository. If it is release artifact, it will be created to release repository. Even we can do the same process for npm artifacts or docker images. As soon as the docker image got created it will get pushed to docker registry available at the nexus repository. Please find the architecture of the flow what we have discussed so far on the code commit till artifact uploaded to NXRM repository at Fig 2.

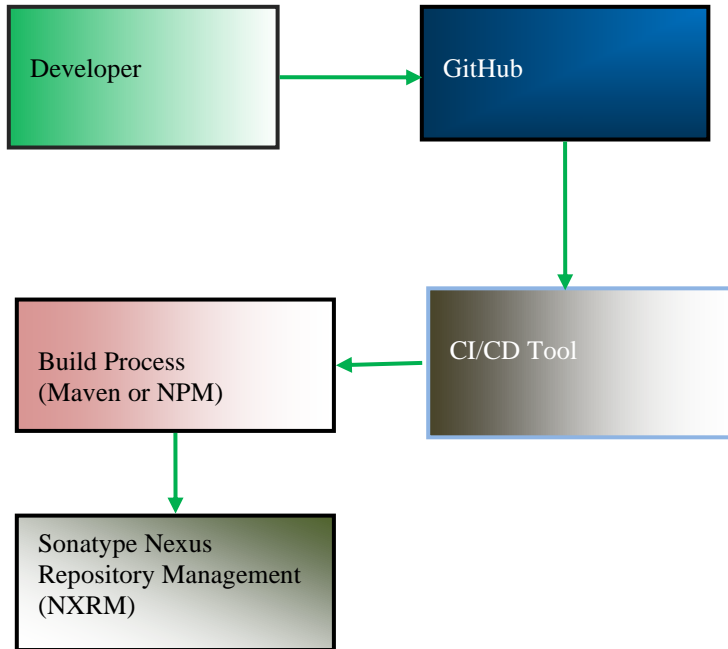


Fig 2: CI/CD Pipeline Integration

Security and Access Control:

We can have admin creds for configuring the NXRM instance. Once we are done with config operation we will be asked to update the creds by providing the password in encrypted form. We can have our own creds. We can integrate Linear Data Access Protocol (LDAP) [15] to the NXRM instance so that we can add AD Groups as well to the instance. For each project team one AD Group will be created and added to instance. User management becomes very easy by adding and removing the members to the ADGroup.

| SNo | System IP | Time Stamp | Operation | URL | Size | Status |
|-----|-----------|------------|-----------|-----|------|---------|
| 1 | 10.0.* | 100924 | PUT | url | 150 | success |
| 2 | 10.1.* | 110924 | DEL | url | 200 | success |
| 3 | 10.2.* | 120924 | GET | url | 250 | success |
| 4 | 10.4.* | 130924 | PUT | url | 300 | success |

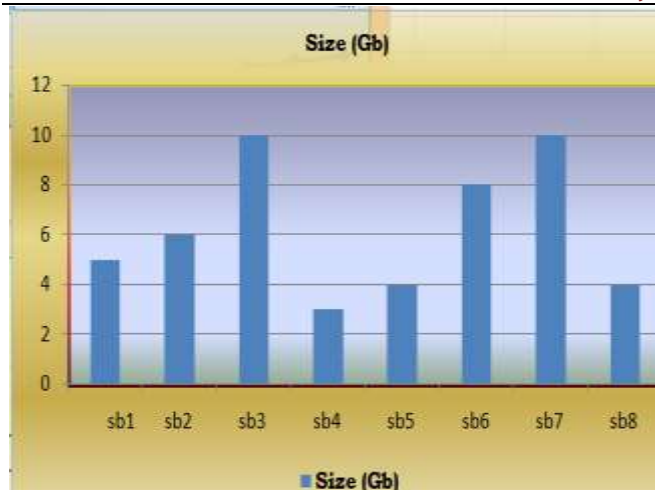
| | | | | | | |
|---|--------|--------|-----|-----|-----|---------|
| 5 | 10.3.* | 140924 | GET | url | 350 | success |
| 6 | 10.6.* | 150924 | PUT | url | 400 | success |
| 7 | 10.5.* | 160924 | DEL | url | 450 | success |
| 8 | 10.4.* | 160924 | GET | url | 500 | success |

Table 1. Nexus Repository log file entries.

For number of repositories Table 1 shows the different type of metrics like name of the repository , time stamp , sub group of the repository, type of the operation GET, PUT or DELETE [16] , volume of the upload, success or failure of the upload. In this we need to segregate the metrics using the type of the operation. Among these operations only PUT operation will directly influence the space of the disk. We need to figure out the total volume of the uploads happening on the daily , weekly and monthly analysis [17]. For example we have taken sample values only having PUT operation for each sub group from the repository. Please find the same in Table 2. The values which we have mentioned in the Table 1 might not be related to values from Table 2.

| subgroup | URL | Time Stamp | Size (Gb) | Operation | Status |
|----------|------|------------|-----------|-----------|---------|
| Sb1 | url1 | 100924 | 5 | PUT | success |
| Sb2 | url2 | 120924 | 6 | PUT | success |
| Sb3 | url3 | 130924 | 10 | PUT | success |
| Sb4 | url4 | 140924 | 3 | PUT | success |
| Sb5 | url5 | 150924 | 4 | PUT | success |
| Sb6 | url6 | 160924 | 8 | PUT | success |
| Sb7 | url7 | 170924 | 10 | PUT | success |
| Sb8 | url8 | 180924 | 4 | PUT | success |

Table 2: Business unit vs Upload operations



Graph 1: Business unit vs Upload operations

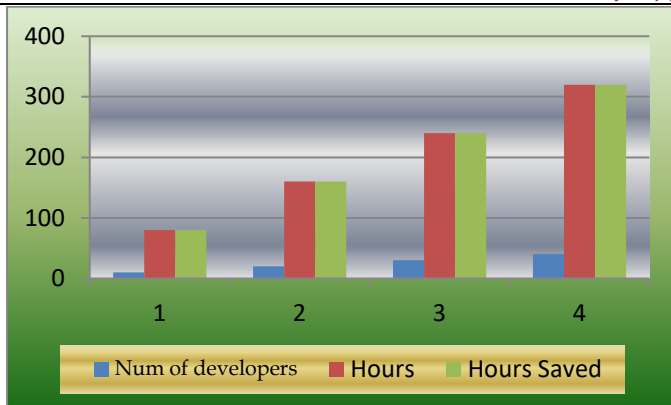
Please observe the usage on daily, weekly, monthly basis. We need to manually delete the unwanted artifacts like snapshots and very old fixed version volumes. If we didn't delete the same, then the upcoming uploads will not be successful which will effect the huge number of developers. Either admin has to delete them manually or by using admin task available at NXRM or he needs to connect with development team to pick the unwanted artifacts and remove the same. The amount of space gaining in the deletion of unwanted artifacts should be in sync with the amount of volume uploaded to NXRM repositories. Since the entire NXRM instance is shared by number of business units and there will be very high volume consumers or low volume consumers. Either we need to provide separate instance for high volume consumers or have them to provide unwanted artifacts list in high frequency so that admin can delete the same to gain the space for further activities. Failing of any of these actions will lead to abnormal failure of uploads by other business units though they are consuming very less amount of space.

Even though we have huge volume available in the beginning of the NXRM activities, one day we need to provide the solution for any of these issues. Instead of working on this process where the manual intervention is very much demanding, if we can predict the usage based on the daily, weekly and monthly usage analysis, it would be easy forecast the expected space required for further activities. If we know well in advance that the available space sufficient only for certain period of time, then it is easy for us to take the appropriate action like alerting the developers to provide unwanted artifacts or increasing the disk space.

If we follow this procedure this will avoid the abnormal incidents like NXRM is not available for operations. This will save the huge number of man hours from development team. If there are 20 members in one team and 10 to 15 teams are depending on the same NXRM then 2400 man hours we will save on daily basis. If there are any abnormal activities because of not having this prediction analysis then we need to face the wastage of man hours. Please find the table shows the man hours analysis.

| SN | Num of developers | Hours | Hours Saved |
|----|-------------------|-------|-------------|
| 1 | 10 | 80 | 80 |
| 2 | 20 | 160 | 160 |
| 3 | 30 | 240 | 240 |
| 4 | 40 | 320 | 320 |

Table 3: Number of developers vs Man hours



Graph 2: Number of developers vs Man hours

Graph 2 shows that number of hours saved by the developers

once we get any solution to avoid abnormal incidents. What ever the man hours available in each team , 100% savings we can observe.

PROPOSAL METHOD

Problem Statement

Usage of space using number of uploads will lead to abnormal incidents of the NXRM instance. Instead of waiting till the end of space and facing the abnormal incident we need to predict the space usage on daily , weekly or monthly basis will help us to patch the disk space so that we can save the number of man hours.

Proposal

Machine Learning algorithm will be used to predict the usage based on the historical analysis of the space usage. We need to consider past one year log files to analyze and figure out the dependent attributes and independent attributes (Feature Engineering). NXRM log files carries repository info in the form of url , time stamp , type of operation , status of the operation, size of the upload, time required to complete the upload operation. We need to analyze the data python data analytics tools like Pandas and numpy to get the exact repository name from the url which the dev team has used to upload the artifacts to nexus repository manager. Finally we will get the repo name , time stamp , artifact size , type of operation and time required to perform the corresponding activity.

In the second time analysis we need to filter the data using type of operation. With this we will get data for PUT operation (upload operation). Like this we need to get the data for the entire month for each repository, Ex: If there are repositories like TCT , RFR , CFC , PRD , GST, TGB, and SCY. Number of users are available in each business unit (repository name for each business unit), and each user is uploading artifacts to NXRM repository. In the analysis we need to filter the data on repo wise , followed by analysis on monthly basis for each business unit. Please find the data at Table 3 having detailed analysis on what we have described so far.

| SNo | URL | Name | Activity | Size(Gb) | Status |
|-----|---|--------|----------|----------|---------|
| 1 | http://ui | Ui | PUT | 4 | Success |
| 2 | http://back | Back | GET | --- | Success |
| 3 | http://middle | Middle | DEL | --- | Success |
| 4 | http://ui | Ui | GET | --- | Success |
| 5 | http://back | Back | GET | --- | Success |
| 6 | http://ui | Ui | GET | --- | Success |
| 7 | http://middle | middle | PUT | 6 | Success |

Table 4: Nexus log file data attributes and values

Ui , middle , back like this they have different type of repo names created inside TCT business unit. The main problem is TCT business unit is facing the space issue. We have started the analysis on getting the repo name and type of operation along with time stamp , size of the artifact for uploading to repository. Monthly wise we will get total volume has been uploaded to each repository. This is how we will get the total space consumption by one business unit (including all repos ui, middle , back end).

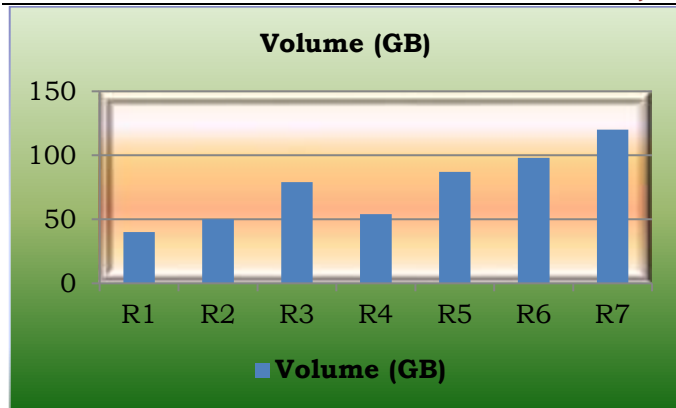
IMPLEMENTATION

We have used python numpy , pandas , Scikit learn to analyze the data . The data has been taken based on daily basis , weekly basis and monthly basis usage analysis. Using this we can get the approx space required to run the NXRM instance without facing the abnormal incidents.

Once we have historical usage for the single business unit on the monthly basis , we can use the same parameters in the machine learning algorithm to predict the dependent variable for future activities. Like this we need to do for all business unit and apply the same feature engineering process and figure out the parameters to apply inside machine learning algorithm (Linear Regression Analysis). This will provide prediction on dependent variable . In our case time is the independent variable and the size of the artifact (total size consumed by all artifacts business unit wise in a month) is the dependent variables.

| SNo | Repo Name | Volume (GB) |
|-----|-----------|-------------|
| 1 | Repo1 | 40 |
| 2 | Repo2 | 50 |
| 3 | Repo3 | 79 |
| 4 | Repo4 | 54 |
| 5 | Repo5 | 87 |
| 6 | Repo6 | 98 |
| 7 | Repo7 | 120 |

Table 5: Repo vs Volume consumption

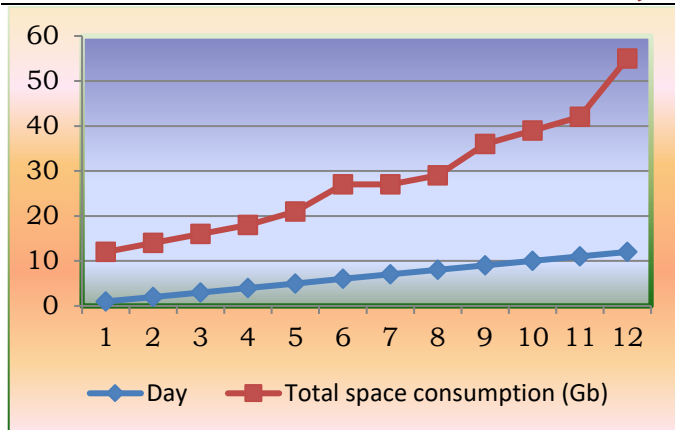


Graph 3: Repo vs Volume consumption

Graph3 shows the volume consumption for each repo (Repo1 → R1, Repo2 → R2, Repo3 → R3, Repo4 → R4, Repo5 → R5, Repo6 → R6 and Repo7 → R7). Please find the usage data for all repos in one business unit per one month. We need to take the consolidate amount space consumed by one business unit in one month. As part of the analysis we have taken daily basis usage for 12 days, so that we will predict the volume required for next x number of days. Please find the table as follows for 12 days in October month

| SNo | Day | Total space consumption (Gb) |
|-----|-----|------------------------------|
| 1 | 1 | 12 |
| 2 | 2 | 14 |
| 3 | 3 | 16 |
| 4 | 4 | 18 |
| 5 | 5 | 21 |
| 6 | 6 | 27 |
| 7 | 7 | 27 |
| 8 | 8 | 29 |
| 9 | 9 | 36 |
| 10 | 10 | 39 |
| 11 | 11 | 42 |
| 12 | 12 | 55 |

Table 6: Space consumption on daily basis



Graph 4: Day vs Volume consumption

Calculating the Slope (m) and Intercept (b):

To compute the values of m and b, we use the least squares method [18], which minimizes the sum of the squared differences between the actual and predicted values. The formulas for the slope (m) and intercept (b) are

$$m = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

$$b = \frac{(\sum y)(\sum x^2) - (\sum x)(\sum xy)}{n(\sum x^2) - (\sum x)^2}$$

n is the number of data points (days in this case), $\sum xy$ is the sum of the product of each x and y, $\sum x$ is the sum of all x values, $\sum y$ is the sum of all y values, $\sum x^2$ is the sum of the squares of all x values. In this case, based on the provided space consumption data over the months (with numerical encoding for months), we calculated the m and b 3.50 and 5.27. Final linear regression equation [19] for space consumption is $y=3.50x+5.27$. If $x=1$ the predicted y is 8.77Gb and if $x=2$ y value is 12.7 Gb. This allows you to predict future space consumption based on this historical trend.

Fig 3: shows the actual and predicted values for the daily usage of volume based on the regression equation derived in the previous step.

| Day | Actual Space Consumption (Gb) | Predicted Space Consumption (Gb) |
|-----|-------------------------------|----------------------------------|
| 0 | 1 | 8.769231 |
| 1 | 2 | 12.265734 |
| 2 | 3 | 15.762238 |
| 3 | 4 | 19.258741 |
| 4 | 5 | 22.755245 |
| 5 | 6 | 26.251748 |
| 6 | 7 | 29.748252 |
| 7 | 8 | 33.244755 |
| 8 | 9 | 36.741259 |
| 9 | 10 | 40.237762 |
| 10 | 11 | 43.734266 |
| 11 | 12 | 47.230769 |

Fig 3: Actual vs Predicted values

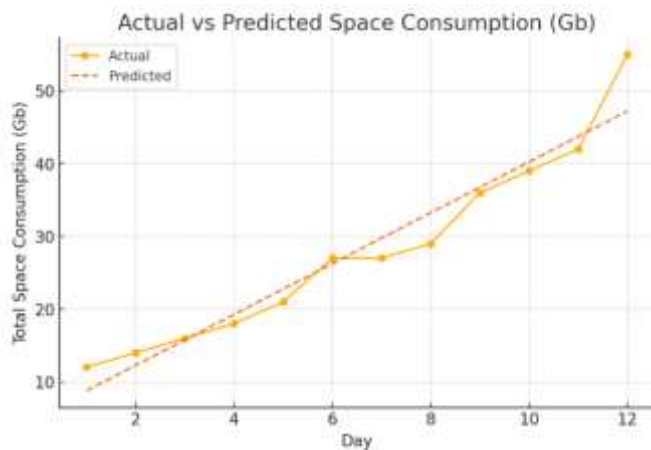
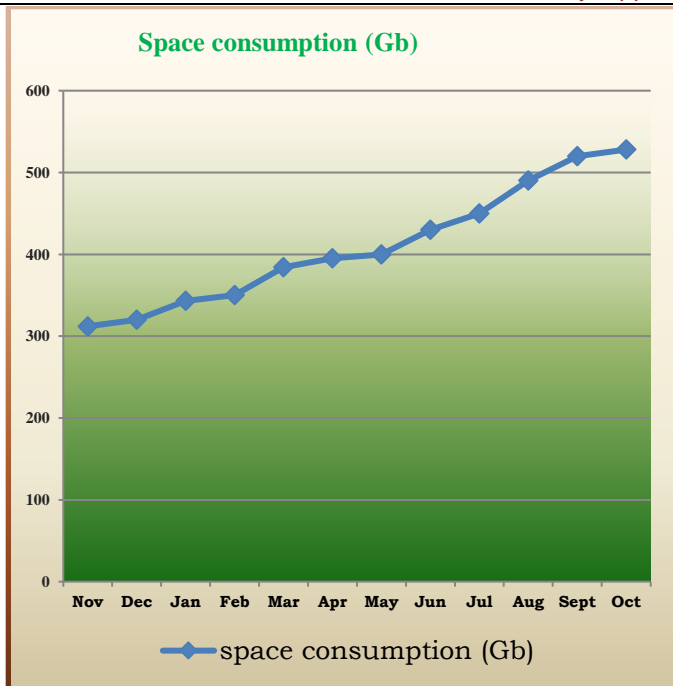
**Graph 5: Actual vs Predicted values**

Table 7 shows the monthly wise sample for volume consumption.

| SNo | Month | Total space consumption (Gb) |
|-----|-----------|------------------------------|
| 1 | November | 312 |
| 2 | December | 320 |
| 3 | January | 343 |
| 4 | February | 350 |
| 5 | March | 384 |
| 6 | April | 395 |
| 7 | May | 400 |
| 8 | June | 430 |
| 9 | July | 450 |
| 10 | August | 490 |
| 11 | September | 520 |
| 12 | October | 528 |

Table 7: Monthly volume consumption (First Sample)

This shows the data for one year. If required we can go for more samples as well. In this data time is the dependent variable and space is the independent variable. We need to find out how much consumption will be there after 4 months or 5 months so that it is easy for us to predict the expected space.



Graph 6: Time vs Space consumption (First Sample)

Graph 6 shows the graph representation of monthly usage analysis. Usage is getting increased on monthly basis even though we have removed number of unwanted artifacts. This is because dev team is not ready to give back old artifacts and they are expecting to preserve them for couple of years.

Univariate Linear Regression:

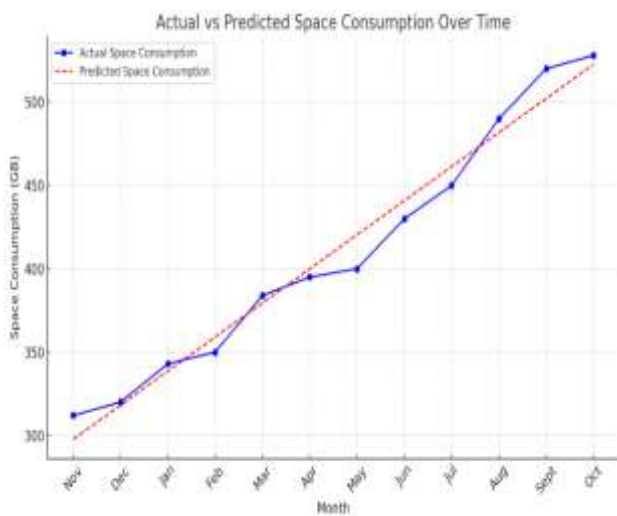
Y is dependent variable [20] and x is independent variable. $Y = mx + b$ is the linear regression equation where we need to find out the slope m and intercept b. In our example if we can find out the intercept and slope, we can predict the volume required for 4 months or 5 months.

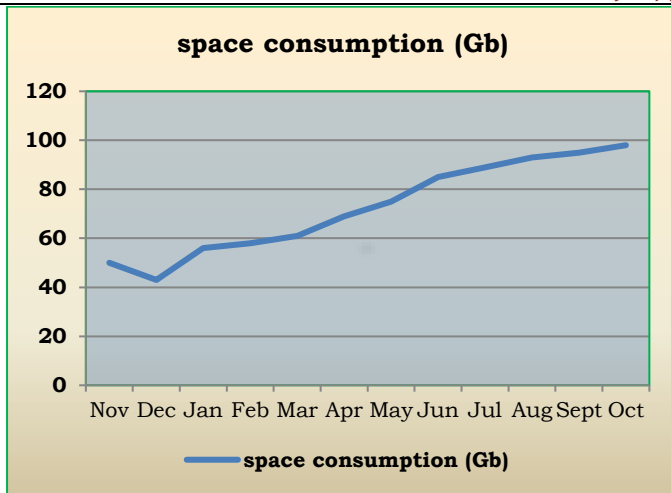
n is the number of data points [21] (months in this case), $\sum xy$ is the sum of the product of each x and y, $\sum x$ is the sum of all x values, $\sum y$ is the sum of all y values, $\sum x^2$ is the sum of the squares of all x values.

In this case, based on the provided space consumption data over the months (with numerical encoding for months), we calculated the m and b 20.45 and 277.26. It means every month memory consumption is increased by 20.45 GB and 277.26 is the predicted space consumption at the baseline (if $x=0$ which could represent a point before the given months).

Final linear regression equation for space consumption is $y = 20.45x + 277.26$. If $x=1$ (for November month) the predicted y is $20.45 \cdot 1 + 277.26$ i.e, 297.71 Gb and if $x=12$ (October) $20.45 \cdot 12 + 277.26 = 522.66$ Gb. This allows you to predict future space consumption based on this historical trend.

| SNo | Month | Actual value (Gb) | Expected Value (GB) |
|-----|-----------|-------------------|---------------------|
| 1 | November | 312 | 297.71 |
| 2 | December | 320 | 318.16 |
| 3 | January | 343 | 338.61 |
| 4 | February | 350 | 359.06 |
| 5 | March | 384 | 379.51 |
| 6 | April | 395 | 399.96 |
| 7 | May | 400 | 420.41 |
| 8 | June | 430 | 440.86 |
| 9 | July | 450 | 461.31 |
| 10 | August | 490 | 481.76 |
| 11 | September | 520 | 502.21 |
| 12 | October | 528 | 522.66 |

Table 8: Actual vs Predicted values for monthly analysis (First Sample)**Graph 7:** Actual vs Prediction values (First Sample)



Graph 8: Monthly usage volume (Second sample)

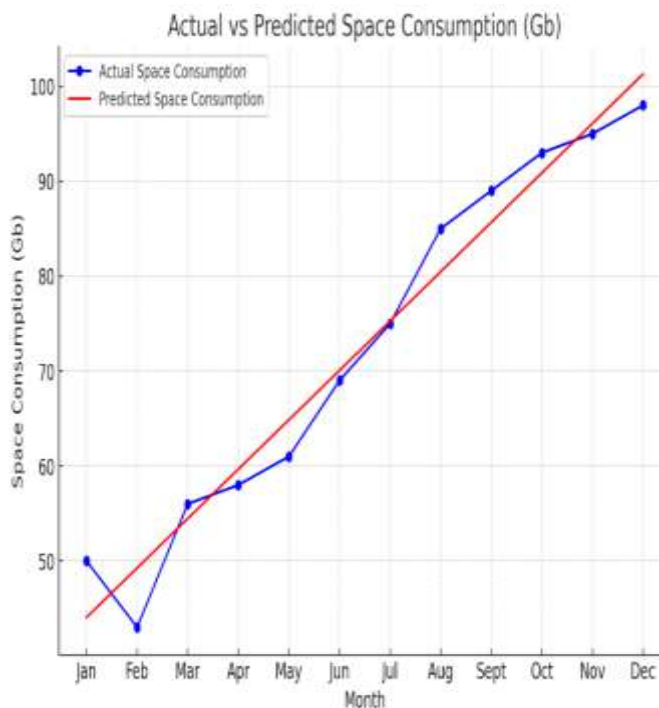
| SNo | Month | Total space consumption (Gb) |
|-----|-----------|------------------------------|
| 1 | November | 50 |
| 2 | December | 43 |
| 3 | January | 56 |
| 4 | February | 58 |
| 5 | March | 61 |
| 6 | April | 69 |
| 7 | May | 75 |
| 8 | June | 85 |
| 9 | July | 89 |
| 10 | August | 93 |
| 11 | September | 95 |
| 12 | October | 98 |

Table 9: Monthly usage volume (Second sample)

| SNo | Month | Actual value (Gb) | Expected value(Gb) |
|-----|----------|-------------------|--------------------|
| 1 | November | 50 | 44.05 |
| 2 | December | 43 | 49.25 |
| 3 | January | 56 | 54.46 |
| 4 | February | 58 | 59.66 |

| | | | |
|----|-----------|----|--------|
| 5 | March | 61 | 64.86 |
| 6 | April | 69 | 70.07 |
| 7 | May | 75 | 75.27 |
| 8 | June | 85 | 80.47 |
| 9 | July | 89 | 85.67 |
| 10 | August | 93 | 90.88 |
| 11 | September | 95 | 96.08 |
| 12 | October | 98 | 101.28 |

Table 10: Actual vs Predicted values for monthly analysis (Second Sample)



Graph 9: Actual vs Prediction values (Second Sample)

As per the given data regression equation is $y=4.06x+44.52$

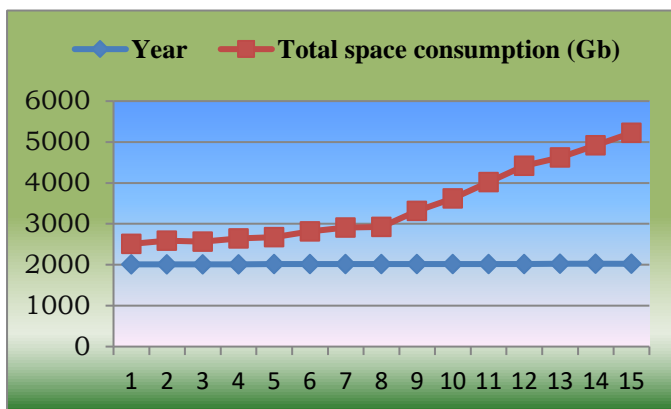
Prediction: Using this equation, you can predict the space consumption for any month within the provided data range. For example, to predict space consumption for March ($x=5$).

$$y=4.06*5+44.52=20.30+44.52=64.82 \text{ Gb}$$

The positive slope indicates that space consumption tends to increase over time, suggesting a growing usage trend.

We have worked on CFC sample having the following dataset.

| SNo | Year | Total space consumption (Gb) |
|-----|------|------------------------------|
| 1 | 2010 | 500 |
| 2 | 2011 | 576 |
| 3 | 2012 | 556 |
| 4 | 2013 | 634 |
| 5 | 2014 | 659 |
| 6 | 2015 | 798 |
| 7 | 2016 | 892 |
| 8 | 2017 | 908 |
| 9 | 2018 | 1300 |
| 10 | 2019 | 1600 |
| 11 | 2020 | 2000 |
| 12 | 2021 | 2400 |
| 13 | 2022 | 2600 |
| 14 | 2023 | 2900 |
| 15 | 2024 | 3200 |

Table 11: Yearly volume consumption**Graph 10:** Yearly volume consumption

Regression Analysis:

Regression equation obtained from the data is

$y=200.58x-403139.32$, y is the predicted total space consumption (in Gb), x is the year.

| SNo | Year | Actual value (Gb) | Expected value (Gb) |
|-----|------|-------------------|---------------------|
| 1 | 2010 | 500 | 215.54 |
| 2 | 2011 | 576 | 231.37 |
| 3 | 2012 | 556 | 431.96 |
| 4 | 2013 | 634 | 632.54 |
| 5 | 2014 | 659 | 833.12 |
| 6 | 2015 | 798 | 1033.70 |
| 7 | 2016 | 892 | 1234.28 |
| 8 | 2017 | 908 | 1434.87 |
| 9 | 2018 | 1300 | 1635.45 |
| 10 | 2019 | 1600 | 1836.03 |
| 11 | 2020 | 2000 | 2036.61 |
| 12 | 2021 | 2400 | 2237.20 |
| 13 | 2022 | 2600 | 2437.78 |
| 14 | 2023 | 2900 | 2638.36 |
| 15 | 2024 | 3200 | 2838.94 |

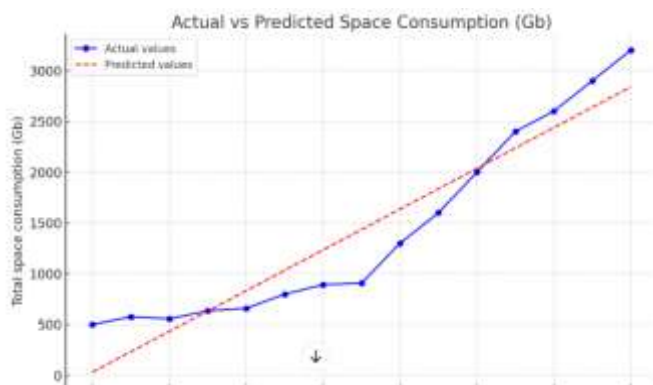
Table 12: Actual vs Prediction values**Graph 11:** Actual vs Prediction values

Table 11 and Graph 10 shows that yearly volume consumption and Table 12 and Graph 11 shows the actual and prediction value after working on the regression equation.

CONCLUSION

NXRM instance has been facing issues with usage of the space. To avoid the abnormal incidents we have prepared the linear regression analysis so that we can predict the space required for the next 3 months 4 months or whatever the duration you need it. Once we have the prediction analysis we can take the necessary precautions like connecting with dev team to get the list of unwanted artifacts so that we can remove the list to gain the space. If it is not satisfies with the predicted values, then extra space request need to be created so that there will not be any abnormal incidents on

NXRM. If we doesn't have this type of analysis, we will not be known until the system goes down. Having this type of analysis saves lot of man hours. We have number of options like connecting with dev team to get the unwanted artifacts to remove and gaining the space or running the cron tasks from the admin tasks. But all these things workout until there is no addition of business. On day to day, business will get increased or the developers are not ready to give the artifacts proportional to uploading the artifacts to NXRM. As long as there is proportionality between upload artifacts and delete artifacts definitely we need to add extra space to disk, else one day we need to face the system hang issue, where all the developers start connecting admin for the help or quick resolution which will not be possible. To avoid all these scenarios the solution which we have provided will help us a lot. Future work includes finding out the solution for two independent variables. In this paper we have taken time is the independent variable and volume is the dependent variable. So this is called univariate linear regression analysis. If we have one more independent variable along with time i.e, number of users in this case as well we need to findout the solution to avoid the abnormal incidents.

REFERENCES

- [1] Apache Maven Cookbook By Raghuram Barathan.
- [2] Docker In Action By Jeff Nickoloff, Stephen Kuenzli, Second Edition, Manning Publications.
- [3] Sonatype Nexus Repository Management Documentation.
- [4] Improving Docker Registry Design Based On Production Workload Analysis, Ali Anwar, Mohamed Mohamed, Vasily Tarasov, Michael Littely, Feb 2018.
- [5] Jenkins 2: Up And Running- Evolve Your Deployment Pipeline For Next Generation Automation, Brent Laster.
- [6] The Maven Repository Dataset Of Metrics, Changes, And Dependencies, Steven Raemaekers, Arie Van Deursen, Joost Visser.
- [7] Linux Administration A Beginners Guide, Wale Soyinka, 1 July 2017.
- [8] The Implementation Of Postgres, Michael Stonebraker, Lawrence A. Rowe And Michael Hirohama Eecs Department, University Of California, Berkeley.
- [9] International Research Journal Of Engineering And Technology (Irjet), Volume: 06 Issue: 04 | Apr 2019.
- [10] An Empirical Study Of Maven Archetype, Xinlei Ma Yan Liu School Of Software Engineering, Tongji University, Shanghai, China.
- [11] Maven Mission Perspectives And Approaches To Inclusion, S.M. Curry.
- [12] A Comprehensive Study Of Bloated Dependencies In The Maven Ecosystem, Nicolas Harrand, Martin Monperrus & Benoit Baudry, 25 March 2021, Volume 26, Article Number 45, (2021).
- [13] An Introduction To Docker And Analysis Of Its Performance, Babak Bashari Rad, Harrison John Bhatti, Mohammad Ahmadi.
- [14] Ci/Cd Pipeline Using Jenkins Unleashed: Solutions While Setting Up Ci/Cd Processes, By Pranoday Pramod Dingare.
- [15] Understanding And Deploying Ldap Directory Services, By Timothy A. Howes Ph.D.
- [16] B. Purnachandra Rao, Dr. N. Nagamalleswara Rao, Hdfs Write Operation Using Fully Connected Digraph Datanode Network Topology, International Journal Of Applied Engineering Research Issn 0973- 4562 Volume 12, Number 16 (2017) Pp. 6076-6090, © Research India Publications., [Http://Www.Ripublication.Com](http://www.Ripublication.Com)
- [17] Hadoop Distributed File System Write Operations, Renukadevi Chuppala, Dr. B. Purnachandra Rao, International Journal Of Intelligent Systems And Applications In Engineering, Issn:2147-6799.
- [18] Method Of Least Square, Alexandria Engineering Journal, 2011, Science Direct.
- [19] Linear Regression Analysis Study, January 2018journal Of The Practice Of Cardiovascular Sciences, Kushbu Kumari, Suniti Yadav.
- [20] Multi Variate Linear Regression Analysis For Sonatype Nexus Repository Space Management, Kishore Kumar Jinka, Dr. B. Purnachandrarao, International Journal Of Intelligent Systems And Applications In Engineering,

Issn:2147-6799.

- [21] Analytical Study Of Two Feature Extraction Methods In Comparison With Deep Learning Methods For Classification Of Small Metal Objects, Somaieh Amraee, Maryam Chinipardaz & Mohammadali Charoosaei , Article Number: 13 (2022).