

Anomaly detection using Unsupervised Machine Learning for Grid computing site operation

Tomoe Kishimoto
ICEPP, University of Tokyo

Nov. 4 2019



ICEPP regional analysis center

- ▶ International Center for Elementary Particle Physics (ICEPP) at U.Tokyo
 - ATLAS experiment at LHC, MEG experiment at PSI, R&D for ILC
- ▶ ICEPP operates a computing center:
 - Support ATLAS VO in WLCG as Tier2 and provide ATLAS–Japan dedicated resources (local use)
 - ~ 10k CPU cores and ~15 PB disk storage
 - Tokyo Tier2 is the only WLCG site in ATLAS–Japan
- ▶ Grid middleware (services):
 - Computing Element: ARC–CE
 - Storage Element: DPM
 - ✓ One of the biggest DPM site in WLCG

Regional analysis center



~300 m²



Grid service logs

- ▶ “Text logs” produced by the Grid services provide useful information for understanding the status of services
 - Important for a stable and reliable site operation
 - e.g DPM logs at Tokyo Tier2:

Time stamp	Severity	Message
May 29 14:05:50 lcg-se01 python2[65569]:	[0]	dmlite Config ConfigFactory : — ConfigFactory started. Starting configuration phase. DMLite v1.12.1
May 29 14:05:50 lcg-se01 python2[65569]:	[4]	dmlite BuiltInAuthnFactory configure : Key: LogLevel Value: 0
May 29 14:05:50 lcg-se01 python2[65569]:	[4]	dmlite BuiltInCatalogFactory configure : Key: LogLevel Value: 0
May 29 14:05:50 lcg-se01 python2[65569]:	[4]	dmlite ConfigFactory LogCfgParm : Key: LogLevel Value: 0
May 29 14:05:50 lcg-se01 python2[65569]:	[0]	dmlite config configure : Setting global log level to :0
May 29 14:05:50 lcg-se01 python2[65569]:	[0]	dmlite config configure : Processing config directory:/etc/dmlite.conf.d/*.conf
May 29 14:05:51 lcg-se01 xrootd[65478]:	!!!	dmlite dome processreq : DN '/C=JP/O=KEK/OU=CRC/OU=ICEPP/CN=host/lcg-se01.icepp.jp' has NOT been authorized.

- Can understand that there is a problem with authentication
- Time-consuming tasks for site administrators to monitor and analyze the service logs everyday...
- A support framework to detect anomaly logs has been developed and examined

Machine Learning approach

- ▶ Log analysis using script does not scale for complicated logs
 - > 1M errors are observed for DPM everyday. Which error is real anomaly...?
 - Difficult to set criteria for all possible error logs
 - **Examine Machine Learning approach based on log “similarity”**
- ▶ “Unsupervised” Machine Learning is used because:
 - Typical classifications require pre-defined labels, but it is difficult to collect a large amount of anomaly logs that cover all possible anomalies
- ▶ Idea of anomaly detection:
 - Text logs → Preprocessing → **Vectorization** → **Clustering**
 - **Aim to detect anomalies in DPM logs in this study**

Preprocessing → Vectorization → Clustering

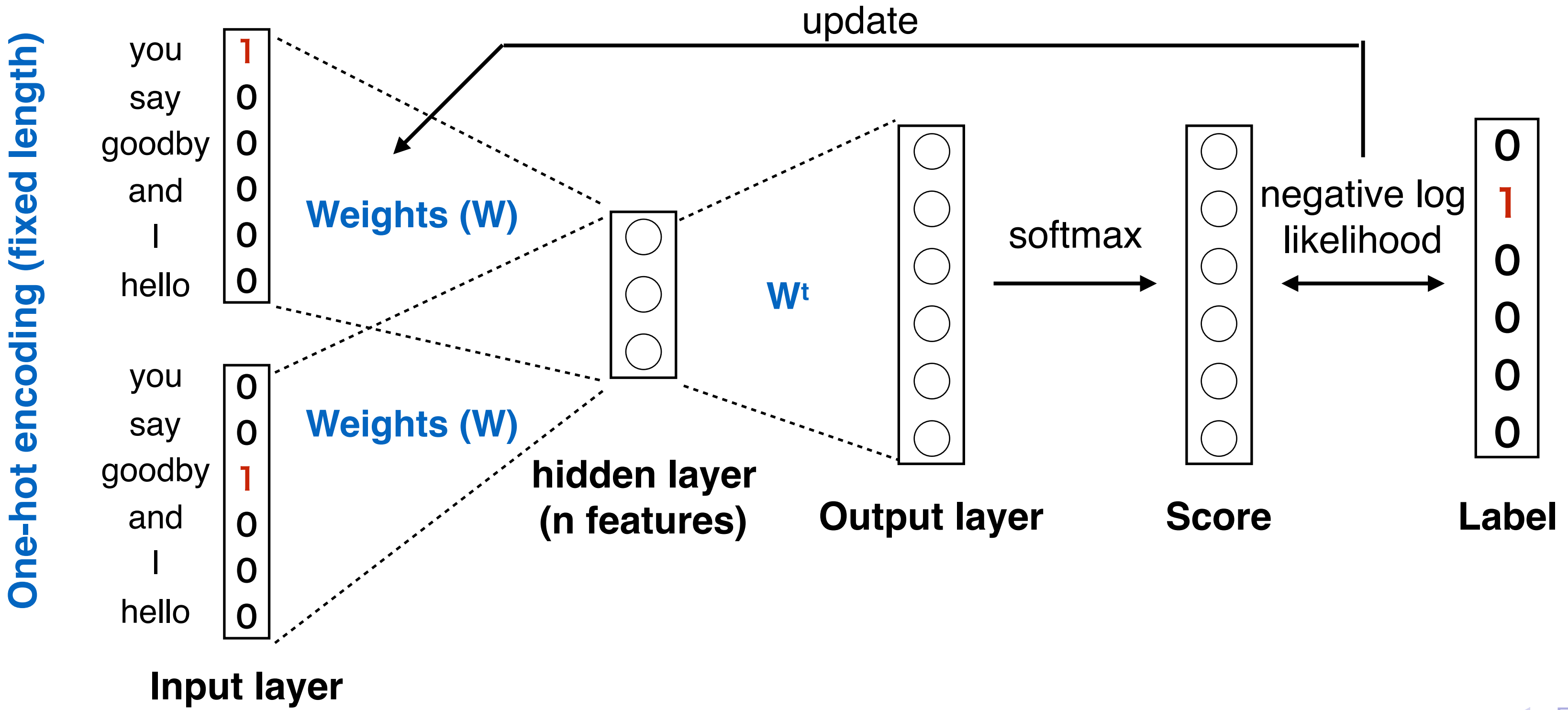
- ▶ Service logs are aggregated using “logwatch” everyday
 - Need to make dedicated scripts for each service
 - For DPM, logwatch outputs “Error counts : Error message” format

Error counts	Error message
336780 time(s)	globus-gridftp-server : dmlite DmStatus : Info: [#00.000002] DPM_LFN not found
333902 time(s)	xrootd : dmlite DmStatus : Info: [#00.000002] replica DPM_LFN not found
333900 time(s)	globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to http://lcg-se01.icepp.jp:1094/domehead/command/dome_getstatinfo . Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'File not found on rfn: DPM_LFN err: NUMBER what: '[#00.000002] replica DPM_LFN not found'
...	...

- Meaningless words are replaced with common words
 - e.g) 0xabcdef (memory address) → NUMBER
- The outputs are saved in sqlite3 file for subsequent ML processes

Preprocessing → **Vectorization** → Clustering

- ▶ Texts are converted to vectors using word embedding technique
 - doc2vec algorithm ([link](#)) is used in this study
 - Use a neural network that predict a target word from context
 - ✓ You ○○ goodbye and I say hello → what is ○○ ?

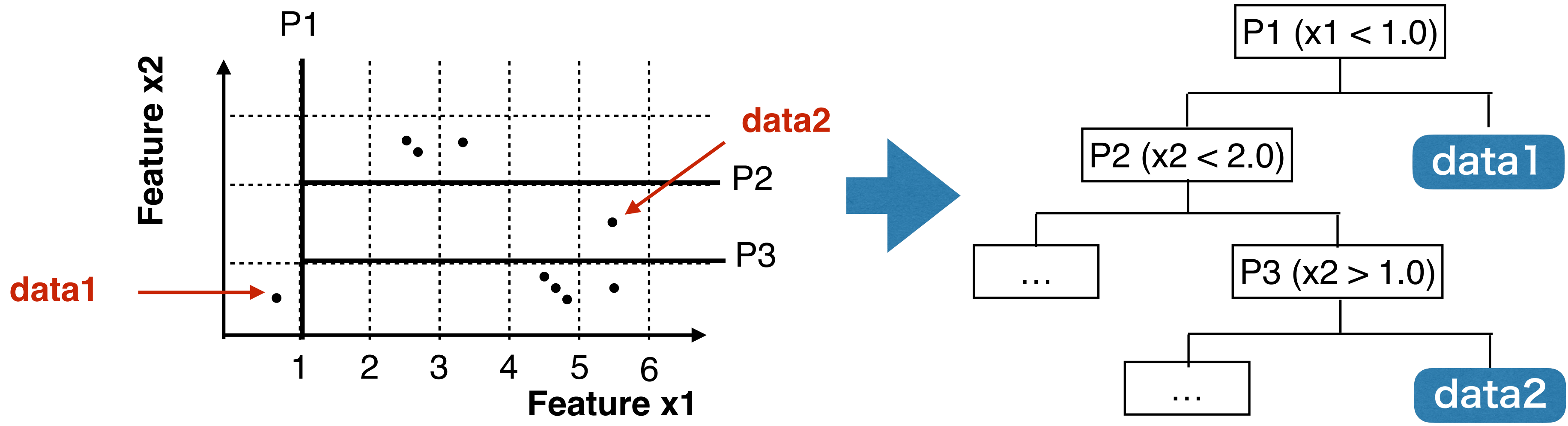


- Weights (W) represents word vectors
- For paragraph, additional matrix is defined for paragraph-id, which represents paragraph vectors

One-hot encoding (fixed length)

Preprocessing → Vectorization → Clustering

- ▶ Clustering n-dimension (feature) vectors in the hidden layer
 - IsolationForest algorithm (link) is used
 - Decision trees are defined with random feature and value
 - Anomaly data tend to result in small number of partitions



1 partition for data1
3 partitions for data2

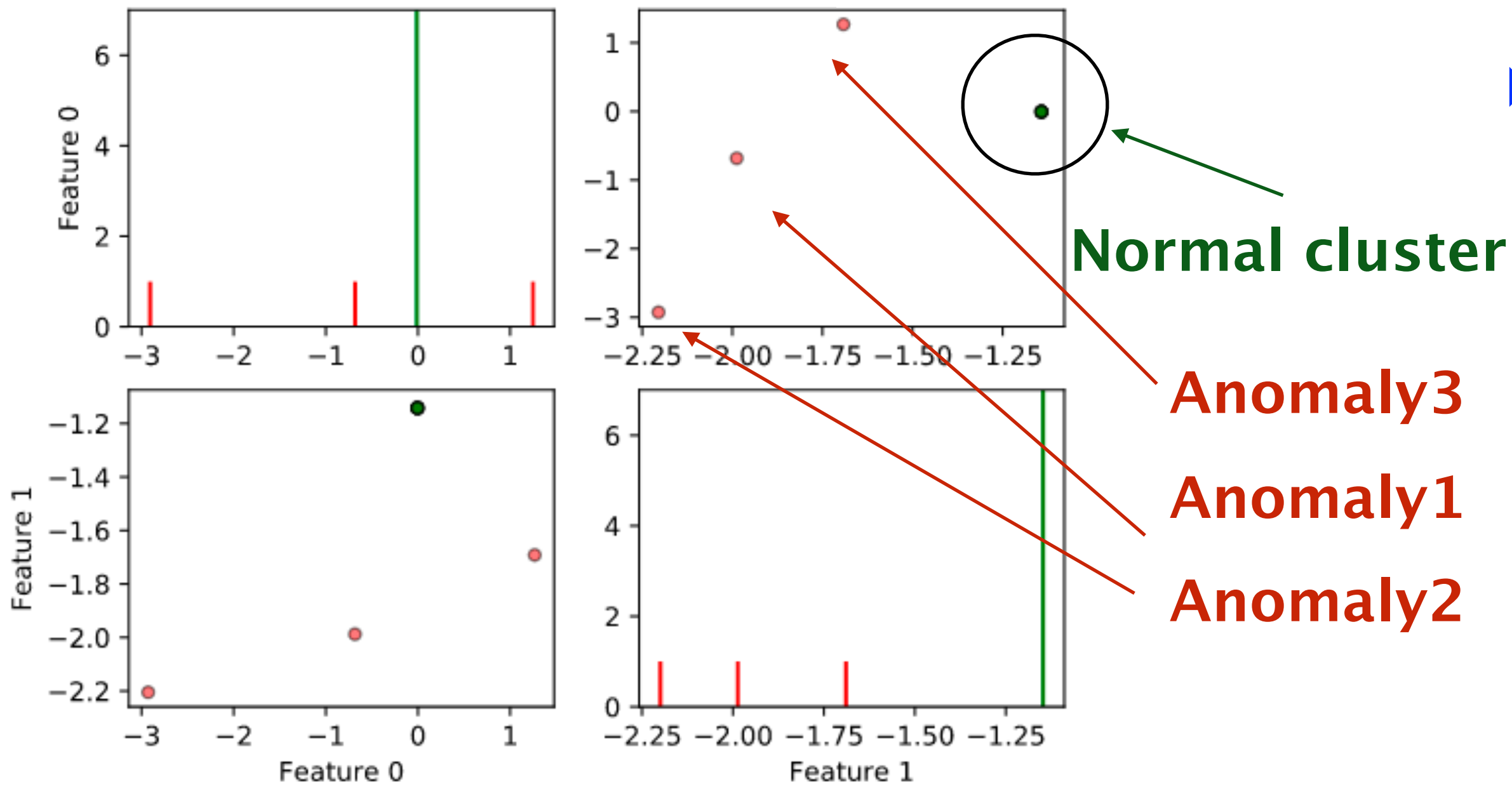
→ **data1 is likely to be anomaly**

- Not need to define the number of clusters as input parameter

Simple test using dummy data

▶ Simple dummy data (sshd log) are defined to check the idea

Label	Text logs	# of samples
Normal	Users logging in through sshd: root: 2001:2f8:102d:589:2:1:2:1 : 288 times	7
Anomaly1	Users logging in through sshd: root: 2001:2f8:102d:589:2:1:2:1 : 300 times	1
Anomaly2	Users logging in through sshd: root: 2001:2f8:102d:589:2:1:2:1 : 1000 times	1
Anomaly3	Users logging in through sshd: root: 2001:2f8:102d:589:2:1:2:1 : 288 times 2001:2f8:102d:589:2:1:2:2 : 10 times	1



- ▶ Converted to two-dimension vectors
 - Numbers are converted to other words to show similarity
 - ✓ 25 → digit2to2 bit2to2 digit2to1 bit2to2to5
 - ✓ 26 → digit2to2 bit2to2 digit2to1 bit2to2to5
 - Anomaly logs are detected as expected using clustering algorithm

DPM logs

- ▶ Use DPM logs from May 29 2019 to Aug 31 2019 at Tokyo Tier2
 - Logs are aggregated everyday using logwatch → 1 day = 1 sample

Date	# of samples	
May 29 - Aug 7	71	Tuning data
Aug 8 - Aug 31	23	Validation data

- ▶ Hyper parameters of doc2vec and IsolationForest are tuned using tuning data
 - There were 7 days operational issues in this tuning data period
 - ✓ e.g.) Configuration of quota was wrong, then data transfers failed

DPM log

```
dmlite DmException : DmException(..):[#00.000028] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_put'. Status 507. DavixError: 'HTTP NUMBER : Insufficient Storage '. Response (NUMBER bytes): 'Unable to complete put for DPM_LFN - quotatoken 'ATLASGROUPDISK' has insufficient free space. minfreespace_bytes: NUMBER
```

- These 7 days are labeled as anomaly, then hyper parameters are tuned to maximize performances

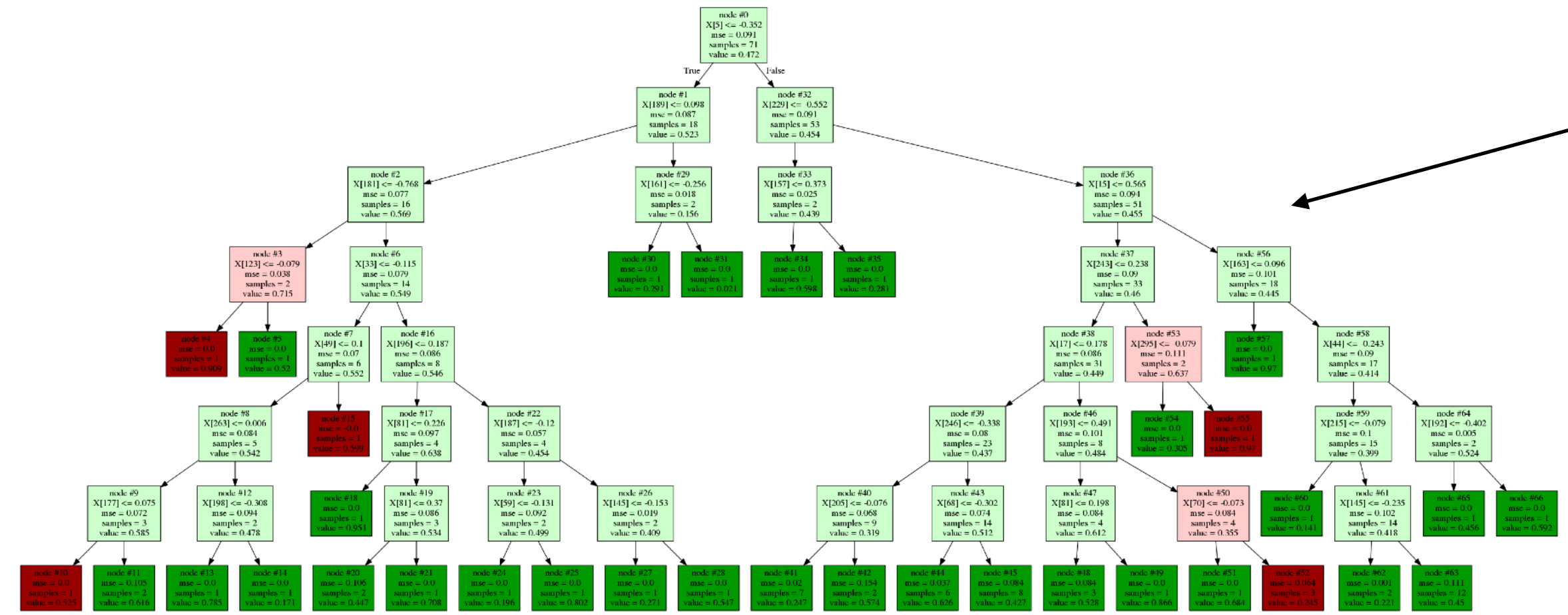
Hyper parameters

- ▶ doc2vec:

- Vector size (dimensionality of feature vectors): 2,3,4,5,10,30,50,100,200,**300**,400,500
- Vector window (distance between the current and predicted word): 3,4,5,6,**7**,8,9,10
- Epochs: 100, **500**, 1000, 2000, 3000

- ▶ IsolationForest

- N estimators (# of base estimators in the ensemble): 5,**10**,30
- Contamination (amount of anomaly data): 0.10, **0.11**, 0.12, 0.13, 0.14



Decision tree of IsolationForest

Green: normal leaf
Red: anomaly leaf

Anomaly leaf tends to have short path length

Results using tuning/validation data

▶ Tuning data:

Accuracy of anomaly detection	Accuracy of normal detection
$6/7 = 0.86$	$62/64 = 0.97$

- Reasonable accuracies are achieved for pre-labeled anomalies
- Will have an anomaly alert ~once a week in this example

▶ Validation data:

- 4 days (Aug 24, 27, 28, 29) are detected as anomaly
- DPM was upgraded on Aug 27, DPM showed new type of logs after the upgrade
 - **Aug 27, 28, 29 are detected as anomaly as expected**
- Can not understand clear reason for Aug 24 (false alert)

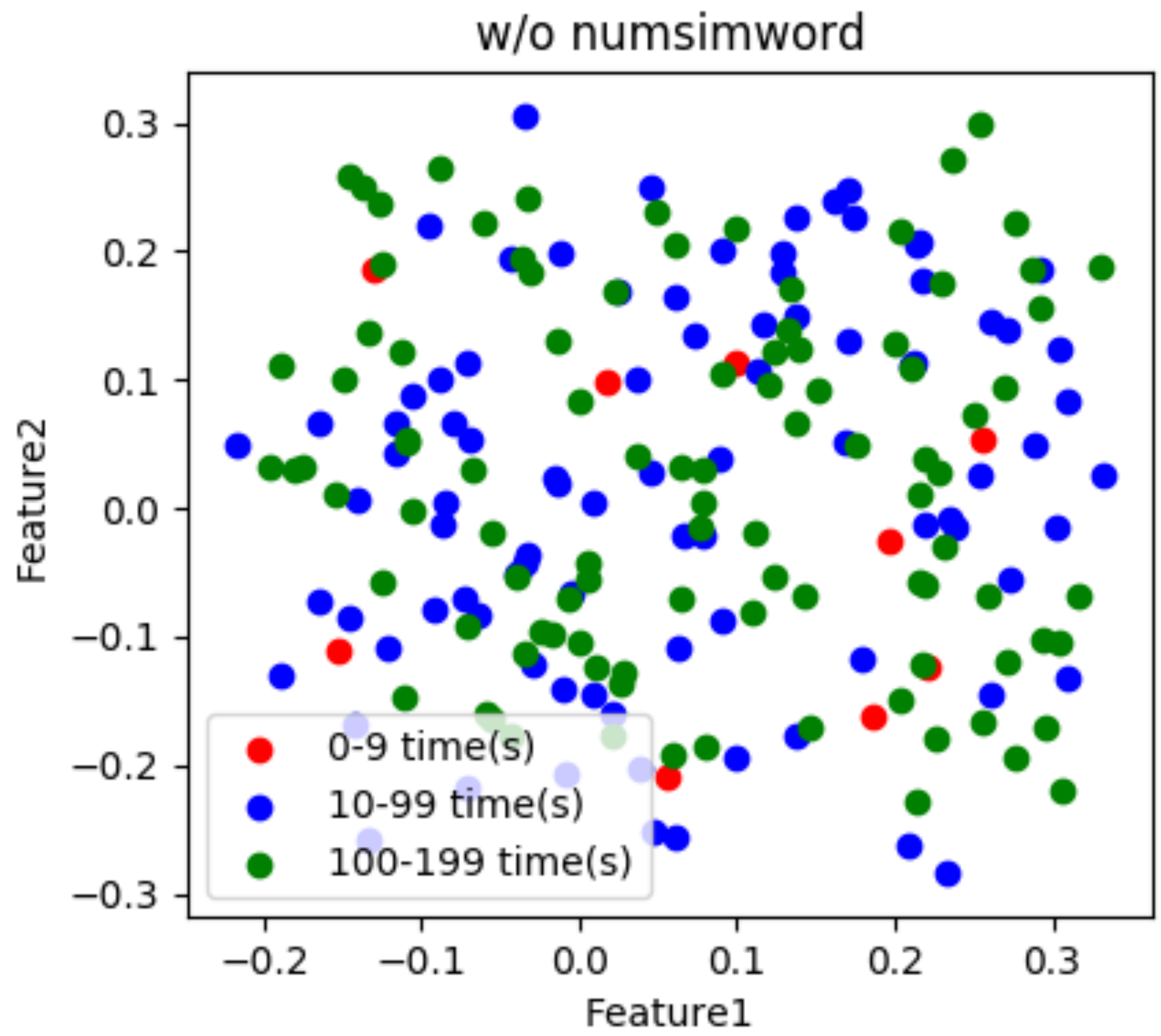
Summary and consideration

- ▶ Anomaly detection using unsupervised Machine Learning technique for Grid site operation is discussed
 - doc2vec and IsolationForest algorithms are used and this workflow is validated using simple data
 - This workflow is applied to the real DPM logs at Tokyo Tier2:
 - ✓ Reasonable accuracy of anomaly detection is observed, $6/7 = 0.86$ in tuning data
 - ✓ Even if unsupervised Machine Learning is used, data labeling is still required to tune hyper parameters
 - ✓ Difficult to understand reason of false alerts
- ▶ Plan: Study a possibility to include information other than text logs, such as CPU usage, memory usage, etc, to detect anomalies

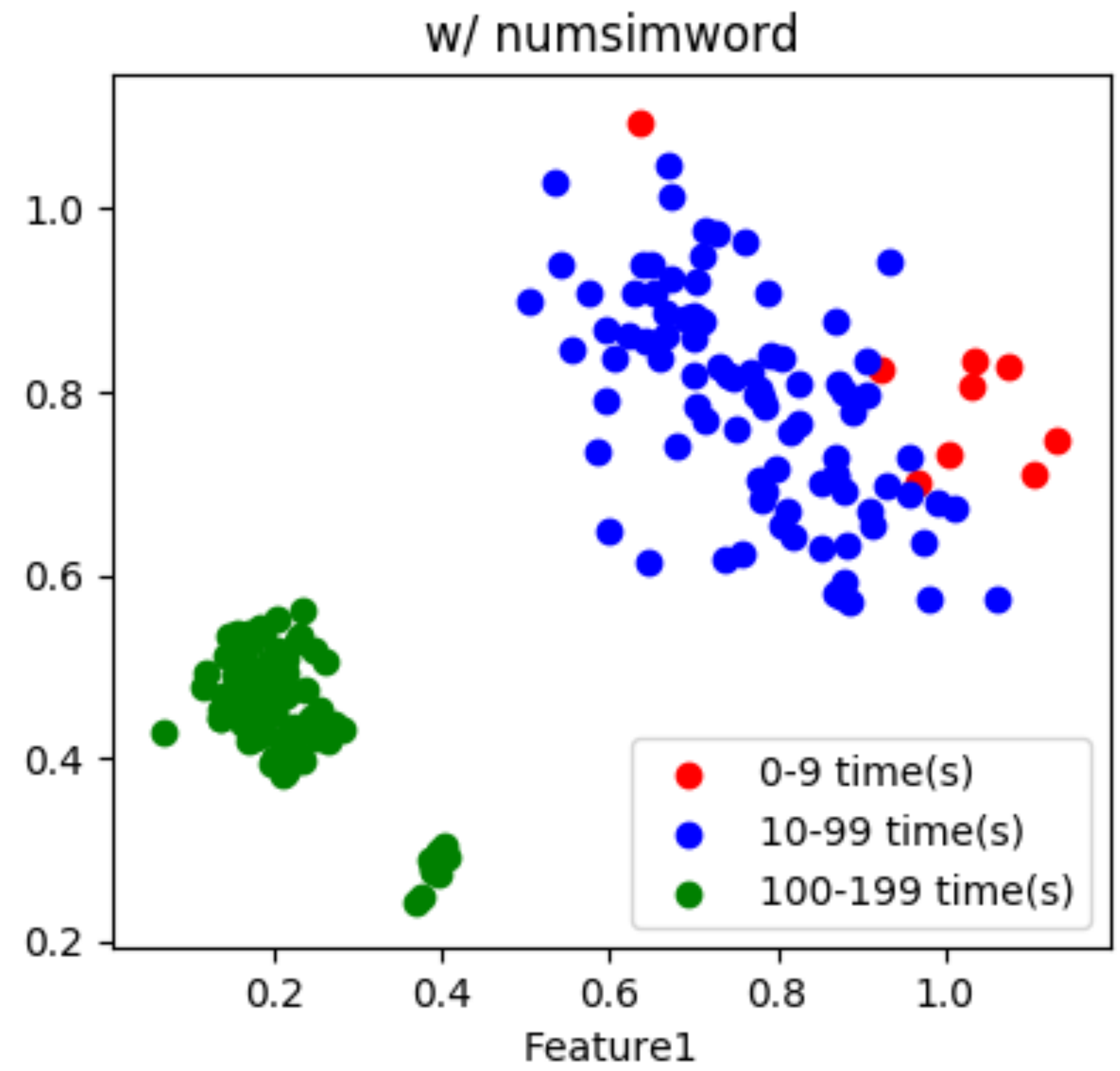
Backup

Number conversion

w/o number conversion



w/ number conversion



Example of DPM logs

```
384756 time(s) : xrootd : dmlite DomeMetadataCache::purgeExpired purgeExpired_fileid : Cached empty record (fileid: NUMBER)
72794 time(s) : globus-gridftp-server : dmlite DmStatus : Info: [#00.000002] DPM_LFN not found
72258 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_getstatinfo'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'File not found on rfn: DPM_LFN err: NUMBER what: '[#00.000002] replica DPM_LFN not found'
29337 time(s) : xrootd : dmlite DmStatus : Info: [#00.000002] replica DPM_LFN not found
27281 time(s) : xrootd : dmlite DmStatus : Info: [#00.000002] file FILE_NAME not found (cached)
24205 time(s) : xrootd : dmlite DmStatus : Info: [#00.000002] file FILE_NAME not found
21888 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000022] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_mkdir'. Status 400. DavixError: 'HTTP NUMBER : Server Error '. Response (NUMBER bytes): 'Cannot create dir DPM_LFN - 22-[#00.000022] Can't create folder RANDOM_DIR
16738 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000013] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_get'. Status 500. DavixError: 'HTTP NUMBER : Unexpected server error: NUMBER '. Response (NUMBER bytes): 'Only pending replicas are available.
12918 time(s) : httpd : dmlite DmStatus : Info: [#00.000002] DPM_LFN not found
9234 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_unlink'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'Cannot stat DPM_LFN : 2-[#00.000002] file FILE_NAME not found (cached)
8954 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000028] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_put'. Status 507. DavixError: 'HTTP NUMBER : Insufficient Storage '. Response (NUMBER bytes): 'Unable to complete put for DPM_LFN - quotatoken 'ATLASGROUPDISK' has insufficient free space. minfreespace_bytes: NUMBER
6851 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_unlink'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'Cannot stat DPM_LFN : 2-[#00.000002] file FILE_NAME not found
5970 time(s) : xrootd : dmlite setMessage : DmException(..):[#03.001062] Duplicate entry FILE_KEY for key 'file_full_id' at #012[bt]: (NUMBER) /usr/lib64/libdome-4.so : dmlite::Statement::throwException()+0x43 HEXNUMBNER (NUMBER) /usr/lib64/libdome-4.so : dmlite::Statement::execute()+0x558 HEXNUMBNER (NUMBER) /usr/lib64/libdome-4.so : DomeMySQL::create(dmlite::ExtendedStat&)+0x3f2 HEXNUMBNER (NUMBER) /usr/lib64/libdome-4.so : DomeMySQL::mkdir(dmlite::ExtendedStat const&, std::string, unsigned int, int, int)+0x1ec HEXNUMBNER
5963 time(s) : xrootd : dmlite DmStatus : Info: [#00.000017] File RANDOM_DIR parent: NUMBER already exists - mysql duplicate key. err: 1062-[#03.001062] Duplicate entry FILE_KEY for key 'file_full_id'
5959 time(s) : xrootd : dmlite DmStatus : Info: [#00.000022] Can't create folder RANDOM_DIR
5958 time(s) : xrootd : dmlite dome dome_mkdir : Cannot create dir DPM_LFN - 22-[#00.000022] Can't create folder RANDOM_DIR
4972 time(s) : xrootd : dmlite DmStatus : Info: [#00.000022] '' is not a valid checksum type.
485 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_unlink'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'Cannot get parent of DPM_LFN : 2-[#00.000002] Entry RANDOM_DIR not found under DPM_LFN
222 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000022] DavixError (NUMBER) Failure HTTP NUMBER : File not found after NUMBER attempts at #012[bt]: (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeTunnelHandler::checkErr(Davix::DavixError*)+0xea HEXNUMBNER (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeTunnelHandler::DomeTunnelHandler(dmlite::DavixCtxPool&, std::string const&, int, unsigned int)+0x135 HEXNUMBNER (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeIODriver::createIOHandler(std::string const&, int, dmlite::Extensible const&, unsigned int)+0x31b HEXNUMBNER (NUMBER) /lib64/libdmlite.so.0 : dmlite_fopen+0x249 HEXNUMBNER
178 time(s) : xrootd : dmlite dome fillSecurityContext : Cannot add unknown group ''
171 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_mkdir'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'Can't find parent path of DPM_LFN
57 time(s) : xrootd : dmlite DmStatus : Info: [#00.000002] DPM_LFN not found
ed by DPM_LFN
23 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000013] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_put'. Status 403. DavixError: 'HTTP NUMBER : Permission refused '. No response to show. at #012[bt]: (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeAdapterPoolManager::whereToWrite(std::string const&)+0x123f HEXNUMBNER (NUMBER) /lib64/libdmlite.so.0 : dmlite_put+0x9c HEXNUMBNER (NUMBER) /usr/lib64/libglobus_gridftp_server_dmlite.so : dmlite_gfs_check_node+0x101 HEXNUMBNER (NUMBER) /usr/lib64/libglobus_gridftp_server_dmlite.so : +0x51dd HEXNUMBNER
19 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000022] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_chksum'. Status 500. DavixError: 'HTTP NUMBER : Unexpected server error: NUMBER '. Response (NUMBER bytes): 'Found a previous finished checksum calculation that likely failed. namekey: DPM_LFN checksum calculations in queue right now: NUMBER
19 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000013] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_mkdir'. Status 403. DavixError: 'HTTP NUMBER : Permission refused '. No response to show. at #012[bt]: (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeAdapterHeadCatalog::makeDir(std::string const&, unsigned int)+0x323 HEXNUMBNER (NUMBER) /lib64/libdmlite.so.0 : dmlite_mkdir+0x8b HEXNUMBNER (NUMBER) /usr/lib64/libglobus_gridftp_server_dmlite.so : +0x4176 HEXNUMBNER (NUMBER) /usr/lib64/libglobus_gridftp_server_dmlite.so : +0x6ac4 HEXNUMBNER
10 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_put'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'Cannot stat the parent of lfn: DPM_LFN
4 time(s) : xrootd : dmlite dome calculateChecksum : Found a previous finished checksum calculation that likely failed. namekey: DPM_LFN checksum calculations in queue right now: NUMBER
4 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_mkdir'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'Can't find parent path of '/dataXX/ops/20XX-XX-XX'
4 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000002] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_mkdir'. Status 404. DavixError: 'HTTP NUMBER : File not found '. Response (NUMBER bytes): 'Can't find parent path of '/dataXX/ops'
3 time(s) : httpd : dmlite setMessage : DmException(..):[#00.000013] Missing token on pfn: /.well-known/oauth-authorization-server at #012[bt]: (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeIODriver::createIOHandler(std::string const&, int, dmlite::Extensible const&, unsigned int)+0xae4 HEXNUMBNER (NUMBER) /lib64/libdmlite.so.0 : dmlite_fopen+0x249 HEXNUMBNER (NUMBER) /usr/lib64/httpd/modules/mod_lcgdm_disk.so : +0x70d8 HEXNUMBNER (NUMBER) /usr/lib64/httpd/modules/mod_lcgdm_dav.so : +0x455f HEXNUMBNER
3 time(s) : httpd : dmlite setMessage : DmException(..):[#00.000013] Missing token on pfn: /favicon.ico at #012[bt]: (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeIODriver::createIOHandler(std::string const&, int, dmlite::Extensible const&, unsigned int)+0xae4 HEXNUMBNER (NUMBER) /lib64/libdmlite.so.0 : dmlite_fopen+0x249 HEXNUMBNER (NUMBER) /usr/lib64/httpd/modules/mod_lcgdm_disk.so : +0x70d8 HEXNUMBNER (NUMBER) /usr/lib64/httpd/modules/mod_lcgdm_dav.so : +0x455f HEXNUMBNER
1 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000013] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_unlink'. Status 0. DavixError: '(Neon): Could not read status line: connection was closed by server'. No response to show. at #012[bt]: (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeAdapterHeadCatalog::unlink(std::string const&)+0x25d HEXNUMBNER (NUMBER) /lib64/libdmlite.so.0 : dmlite_unlink+0x83 HEXNUMBNER (NUMBER) /usr/lib64/libglobus_gridftp_server_dmlite.so : +0x41ab HEXNUMBNER (NUMBER) /usr/lib64/libglobus_gridftp_server_dmlite.so : +0x6ac4 HEXNUMBNER
1 time(s) : globus-gridftp-server : dmlite DmException : DmException(..):[#00.000022] Error when issuing request to 'http://lcg-se01.icepp.jp:1094/domehead/command/dome_getidmap'. Status 0. DavixError: '(Neon): Could not read status line: connection was closed by server'. No response to show. at #012[bt]: (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeAdapterAuthn::uncachedGetIdMap(std::string const&, std::vector<std::string, std::allocator<std::string>> const&, dmlite::UserInfo*, std::vector<dmlite::GroupInfo, std::allocator<dmlite::GroupInfo>> *)+0xf26 HEXNUMBNER (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeAdapterAuthn::getIdMap(std::string const&, std::vector<std::string, std::allocator<std::string>> const&, dmlite::UserInfo*, std::vector<dmlite::GroupInfo, std::allocator<dmlite::GroupInfo>> *)+0x314 HEXNUMBNER (NUMBER) /usr/lib64/dmlite/plugin_domeadapter.so : dmlite::DomeAdapterAuthn::createSecurityContext(dmlite::SecurityCredentials const&)+0xa9 HEXNUMBNER (NUMBER) /lib64/libdmlite.so.0 : dmlite::StackInstance::setSecurityCredentials(dmlite::SecurityCredentials const&)+0x2f6 HEXNUMBNER
1 time(s) : xrootd : dmlite DmStatus : Info: [#00.000020] '' is not a directory, and is referenced by DPM_LFN
```