

ApplianceTelemetryCorrelationAnalysis

October 19, 2024

0.1 Provided as-is (w/o support)

Kubernetes clusters collect various application and infrastructure statistics. While this information is useful, it's very difficult to identify which metrics are useful for monitoring and troubleshooting. The Goal here is to collect this information, and use a statistical model to identify which metrics should be included in reports/dashboard such that: * Unnecessary overhead and sensory overload can be reduced. * Time can be saved by prioritising monitoring the correct metrics.

This process needs to assume **zero knowledge** of the workings of the cluster, workload being run and any other information. This way, **generic** clusters can be monitored without explicitly programming dashboards based on internal knowledge. This is also a good method to discover/verify application knowledge/bottlenecks with statistical data analysis.

0.2 Step1: Data Loading

We will load cpu, memory, task_queue information along with stats from structured and unstructured scans from csv files stored on disk using the `dataframeLoader` helper.

*# The dataframeLoader helper function implements the loadApplianceTimeSeriesData method.
This method loads the csv files, and pivots them to generate distinct "metrics" timeseries.
see <https://github.com/amitgupta7/docker-jupy-ntbk-s3-reporting/blob/main/dataframeLoader.py>*

```
[1]: import sys
      sys.path.append('../')

      import dataframeLoader as dfl
      import pandas as pd
      from importlib import reload
      reload(dfl)

      # Provide csv data location and appliance and timerange information.
      root = '../..dataDir'
      fromDt = '2024-08-15'
      toDt = '2024-10-15'

      # Provide list of prometheus metrics to load.
      # metricsArr = ['cpu_used', 'download_workers_count', 'memory_used',
      ↪ 'task_queue_length', 'infra_access_latency', 'pod_cpu_usage',
      ↪ 'pod_memory_usage']
      metricsArr = ['cpu_used']
```

```

        , 'task_queue_length'
        , 'memory_used'
    ]

    daterange=[fromDt, toDt]
    df = dfl.loadApplianceTimeSeriesData(root, metricsArr, daterange)

```

```

loading Unstrctured Data from file: SCANPROC-*.csv
loading Strctured Data from file: STRUCTURED-*.csv
processing securiti_appliance_cpu_used-max*.csv
processing securiti_appliance_cpu_used-avg*.csv
processing securiti_appliance_task_queue_length-max*.csv
processing securiti_appliance_task_queue_length-avg*.csv
processing securiti_appliance_memory_used-max*.csv
processing securiti_appliance_memory_used-avg*.csv
loading Unstrctured Data from file: UNSTRUCTURED-*.csv

```

0.3 Step2: Data Pivoting

We now aggregate the data by `appliance_id` (unique identifier for our cluster) and `ts` timestamp, to get different metrics values as separate columns. Notice there are: * 21 metrics * Tracked every hour

```

[2]: dfp = df.pivot_table(index=['appliance_id','ts'], columns=['metrics'],
    ↪values='value', aggfunc='sum').reset_index()
dfp.head()

```

```

[2]: metrics                appliance_id                ts \
0      0036f473-ad7f-4439-8d37-f65fdeb50b2d  2024-10-13  14:00:00
1      0036f473-ad7f-4439-8d37-f65fdeb50b2d  2024-10-13  15:00:00
2      0036f473-ad7f-4439-8d37-f65fdeb50b2d  2024-10-13  16:00:00
3      0036f473-ad7f-4439-8d37-f65fdeb50b2d  2024-10-13  17:00:00
4      0036f473-ad7f-4439-8d37-f65fdeb50b2d  2024-10-13  18:00:00

```

```

metrics  IdleTimeInHrs  avgFileSizeInMB  cpu_used_avg  cpu_used_max \
0                NaN                NaN        3.021810        21.46
1                NaN                NaN        1.569917         3.14
2                NaN                NaN        1.748750         2.98
3                NaN                NaN        1.740000         1.93
4                NaN                NaN        1.740000         1.93

```

```

metrics  dataScannedinGB  fileDownloadTimeInHrs  linkerq_avg  linkerq_max \
0                NaN                NaN        NaN        NaN
1                NaN                NaN        NaN        NaN

```

2		NaN		NaN		NaN		NaN
3		NaN		NaN		NaN		NaN
4		NaN		NaN		NaN		NaN

metrics	...	memory_used_max	numFilesScanned	numberOfChunksScanned	\
0	...	74.53	NaN		NaN
1	...	66.62	NaN		NaN
2	...	66.54	NaN		NaN
3	...	66.33	NaN		NaN
4	...	66.33	NaN		NaN

metrics	numberOfColsScanned	scanTime	taskq_avg	taskq_max	tmp_taskq_avg	\
0	NaN	NaN	NaN	NaN		NaN
1	NaN	NaN	NaN	NaN		NaN
2	NaN	NaN	NaN	NaN		NaN
3	NaN	NaN	NaN	NaN		NaN
4	NaN	NaN	NaN	NaN		NaN

metrics	tmp_taskq_max	uniqPodCount
0	NaN	NaN
1	NaN	NaN
2	NaN	NaN
3	NaN	NaN
4	NaN	NaN

[5 rows x 21 columns]

0.4 Step 3: Data transformation and correlation

We need to achieve two main goals: 1. Isolate data for individual appliance. 2. Remove ghost correlation between unrelated metrics. * We will calculate percentage change between adjacent timeseries values. 3. Calculate absolute correlation between metrics for each single appliance. * Transpose every metrics correlation. 4. Generate correlation for every `appliance_id` and `metric` identifier using steps 1, 2 and 3

```
[3]: # appliance = '01c75278-9c0d-41be-b693-c970b18dbedc'
# for metric in metrics_category_order:
dfc_arr = []
for pod in dfp.appliance_id.unique():
    dfa = dfp[(dfp.appliance_id == pod)]
    dfa = dfa.drop(['appliance_id', 'ts'], axis=1)
    dfa = dfa.pct_change(periods=1, fill_method=None)
    dfca = dfa.corr().abs()
    # print(type(dfca))
    for col in dfca.columns:
        # print(col)
        dfc = dfca[col].to_frame().T
```

```

dfc.insert(0, 'metric', col )
dfc.insert(0, 'appliance_id', pod )
dfc_arr.append(dfc)
dfc = pd.concat(dfc_arr, ignore_index=True)
dfc.set_index('appliance_id', inplace=True)
dfc.head()

```

```

[3]: metrics                                metric  IdleTimeInHrs  \
appliance_id
0036f473-ad7f-4439-8d37-f65fdeb50b2d      IdleTimeInHrs      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d  avgFileSizeInMB      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      cpu_used_avg      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      cpu_used_max      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d  dataScannedinGB      NaN

metrics                                avgFileSizeInMB  cpu_used_avg  \
appliance_id
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      1.000000
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      0.536825
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN

metrics                                cpu_used_max  dataScannedinGB  \
appliance_id
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      0.536825      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      1.000000      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN

metrics                                fileDownloadTimeInHrs  linkerq_avg  \
appliance_id
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN

metrics                                linkerq_max  memory_used_avg  \
appliance_id
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      0.323098
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      0.907248
0036f473-ad7f-4439-8d37-f65fdeb50b2d      NaN      NaN

```

metrics appliance_id	memory_used_max	numFilesScanned \
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	0.622750	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	0.256122	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN

metrics appliance_id	numberOfChunksScanned \
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN

metrics appliance_id	numberOfColsScanned	scanTime \
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN

metrics appliance_id	taskq_avg	taskq_max	tmp_taskq_avg \
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN	NaN

metrics appliance_id	tmp_taskq_max	uniqPodCount
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN
0036f473-ad7f-4439-8d37-f65fdeb50b2d	NaN	NaN

0.5 Step 4: Isolate related metrics using correlation

We now iterate over each `metric`, to see if there is any significant statistical correlation to be found across `appliance_ids`. This is done with two steps:

1. Removing outliers:
 - Remove any metrics with `mean correlation value` below the cut-off. The cut-off can be varied for depending on use cases:

- 0.9 for Exec Dashboards
- 0.7 for Customer Ops
- 0.5 for L1 - support
- 0.3 for L2 - suport

Please note that we are filtering metrics with **mean** correlation below the **low cut-off**. This ensures that atleast half of the values are correlated to reduce outliers.

2. Plot box chart to visually represent metrics with any correlation (for cutoff as 0.3).
3. Decide between **max** or **avg** values if both are present. We chose to display **avg** values metrics in this case.

0.6 Final List of metrics

The below table shows the list of **metrics** that are useful with respective correlation **cutoff**. The cut-off values can be interpreted as follows: * below 0.3 negligible correlation * 0.3 to 0.5 Low positive (negative) correlation * 0.5 to 0.7 Moderate positive (negative) correlation * 0.7 to 0.9 High positive (negative) correlation * 0.9 to 1.0 Very High positive (negative) correlation

0.9	0.7	0.5	0.3
linkerq_avg	linkerq_avg	linkerq_avg	linkerq_avg
numberOfChunksScanned	numberOfChunksScanned	numberOfChunksScanned	numberOfChunksScanned
numberOfColsScanned	numberOfColsScanned	numberOfColsScanned	numberOfColsScanned
numFilesScanned	numFilesScanned	numFilesScanned	numFilesScanned
	scanTime	scanTime	scanTime
	tmp_taskq_avg	tmp_taskq_avg	tmp_taskq_avg
	fileDownloadTimeInHrs	fileDownloadTimeInHrs	fileDownloadTimeInHrs
		avgFileSizeInMB	avgFileSizeInMB
		dataScannedinGB	dataScannedinGB
		memory_used_avg	memory_used_avg
		IdleTimeInHrs	IdleTimeInHrs
		cpu_used_avg	cpu_used_avg
			uniqPodCount

```
[4]: corr_vals = [0.9, 0.7, 0.5, 0.3]
line = set()
for cutoff in corr_vals:
    arr = []
    for metr in dfc.metric.unique():
        dfcm = dfc[(dfc.metric == metr)]
        dfcm = dfcm.drop('metric', axis=1)
        dfcm = dfcm.drop(metr, axis=1)
        dfcm = dfcm.dropna(axis = 0, how = 'all')
        dfcm = dfcm.loc[:, dfcm.median() > cutoff]
        [arr.append(x) for x in dfcm.columns]
        # dfcm = dfcm.dropna(axis = 1, thresh=getMiniumValidValues(dfcm, pct=10,
↪ ceiling=10))
```

```

# display(dfcm)
# break
if(cutoff == 0.3):
    if len(dfcm.columns) > 0:
        title=f'''Absolute correlation vs percent-change of {metr}
        (For median correlation greater than {cutoff})
        '''
        dfcm.plot(kind='box'
                    ,vert=False
                    ,title=title
                    ,colormap='tab20'
                    )
    for met in set(arr):
        if("max" not in met):
            line.add(met)
print(cutoff, line)

```

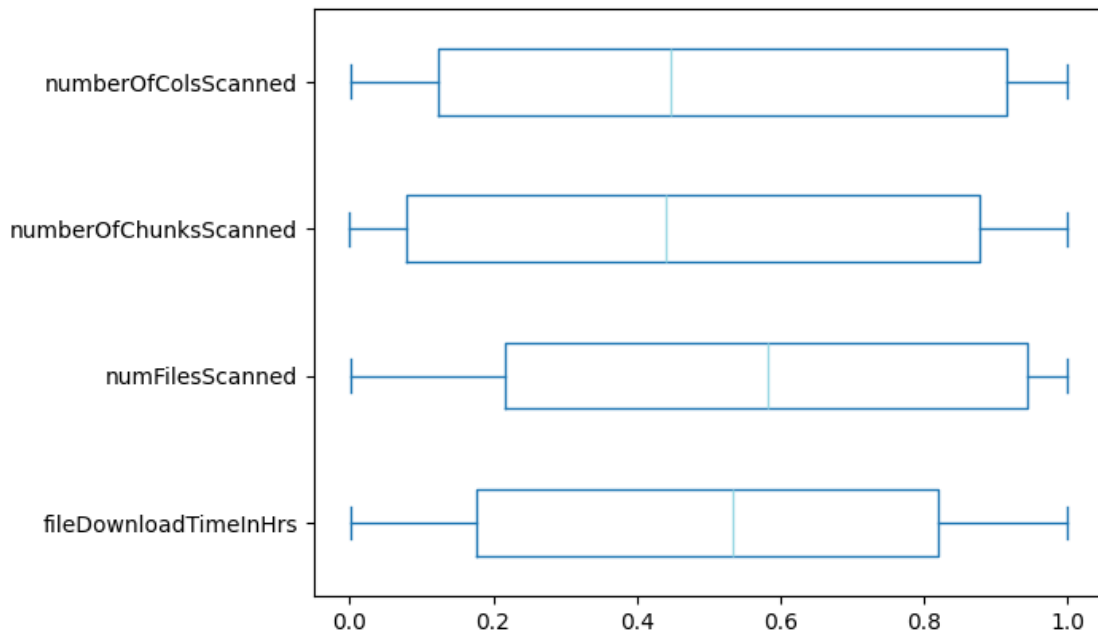
```

0.9 {'linkerq_avg', 'numberOfChunksScanned', 'numberOfColsScanned',
'numFilesScanned'}
0.7 {'fileDownloadTimeInHrs', 'numberOfColsScanned', 'scanTime', 'linkerq_avg',
'tmp_taskq_avg', 'numberOfChunksScanned', 'numFilesScanned'}
0.5 {'fileDownloadTimeInHrs', 'memory_used_avg', 'numberOfColsScanned',
'scanTime', 'dataScannedinGB', 'IdleTimeInHrs', 'linkerq_avg',
'avgFileSizeInMB', 'tmp_taskq_avg', 'numberOfChunksScanned', 'cpu_used_avg',
'taskq_avg', 'numFilesScanned'}

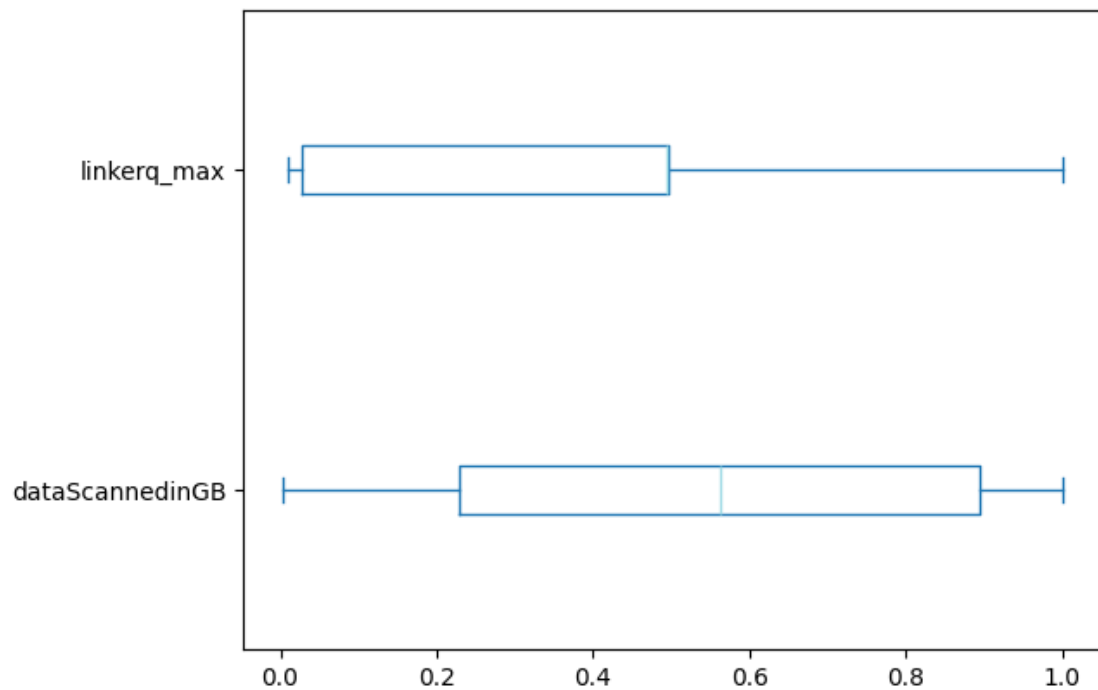
0.3 {'fileDownloadTimeInHrs', 'memory_used_avg', 'numberOfColsScanned',
'scanTime', 'dataScannedinGB', 'IdleTimeInHrs', 'uniqPodCount', 'linkerq_avg',
'avgFileSizeInMB', 'tmp_taskq_avg', 'numberOfChunksScanned', 'cpu_used_avg',
'taskq_avg', 'numFilesScanned'}

```

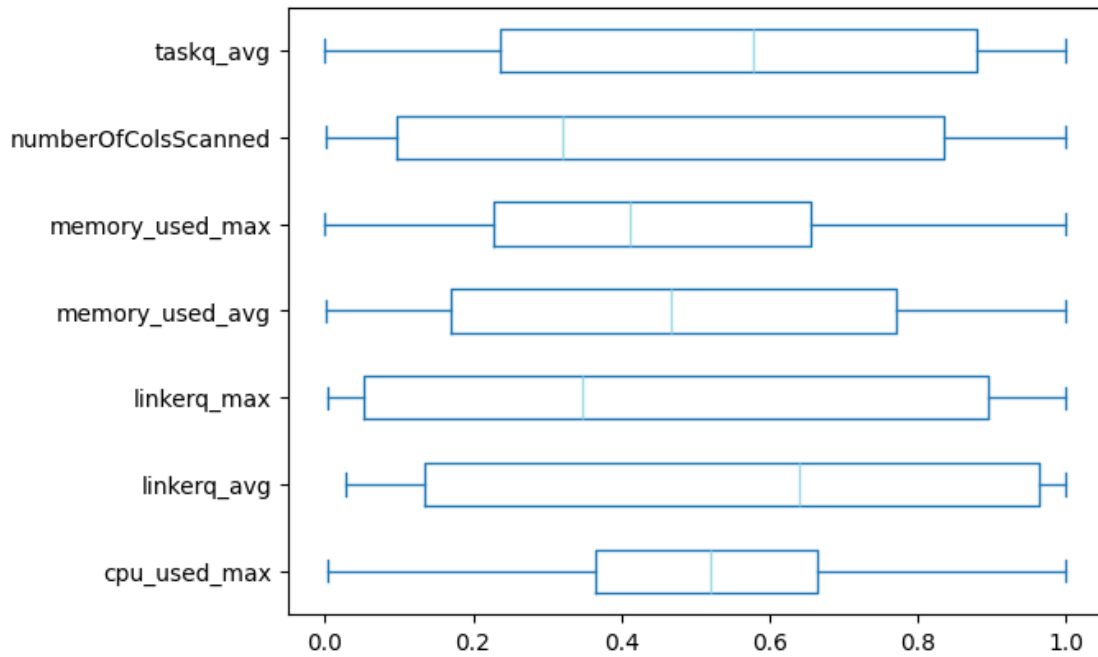
Absolute correlation vs percent-change of IdleTimeInHrs
(For median correlation greater than 0.3)



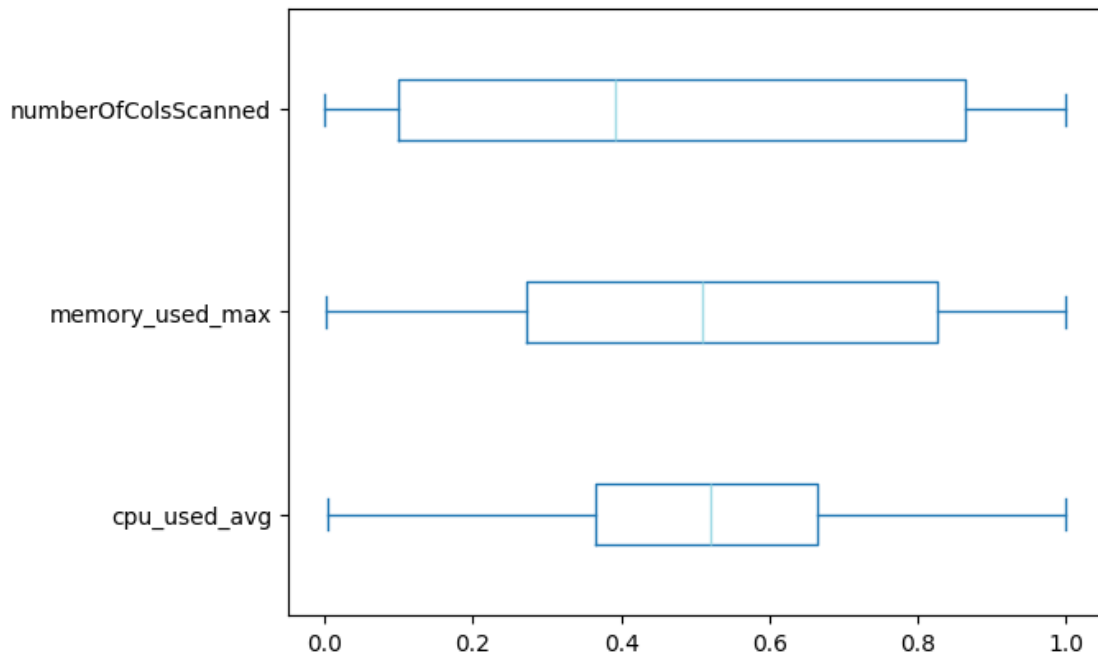
Absolute correlation vs percent-change of avgFileSizeInMB
(For median correlation greater than 0.3)



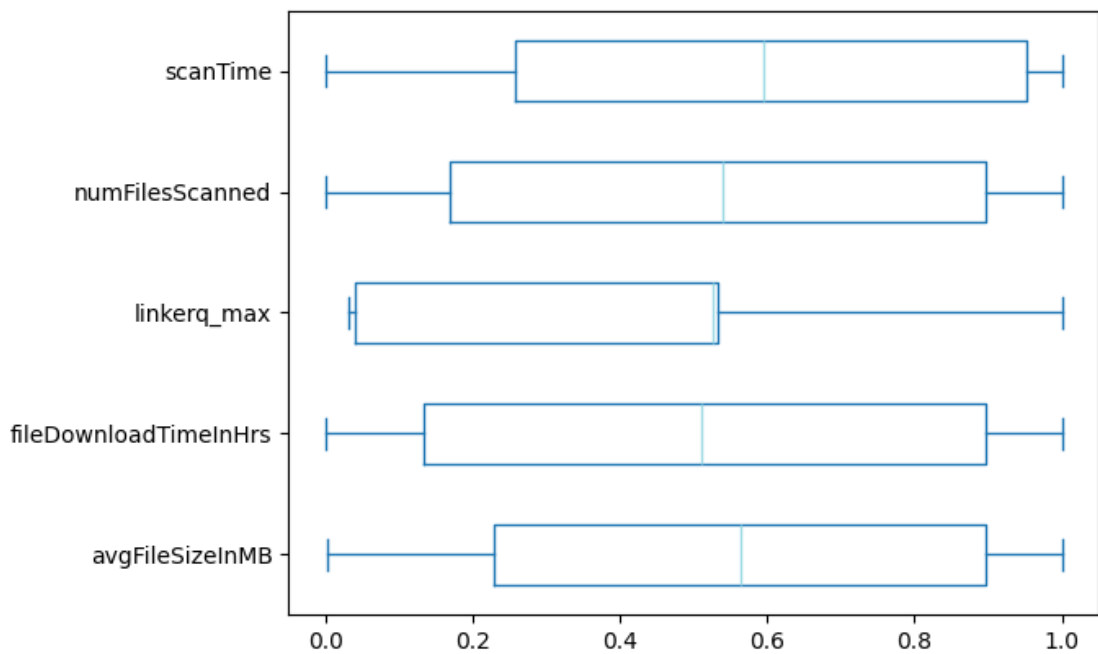
Absolute correlation vs percent-change of `cpu_used_avg`
(For median correlation greater than 0.3)



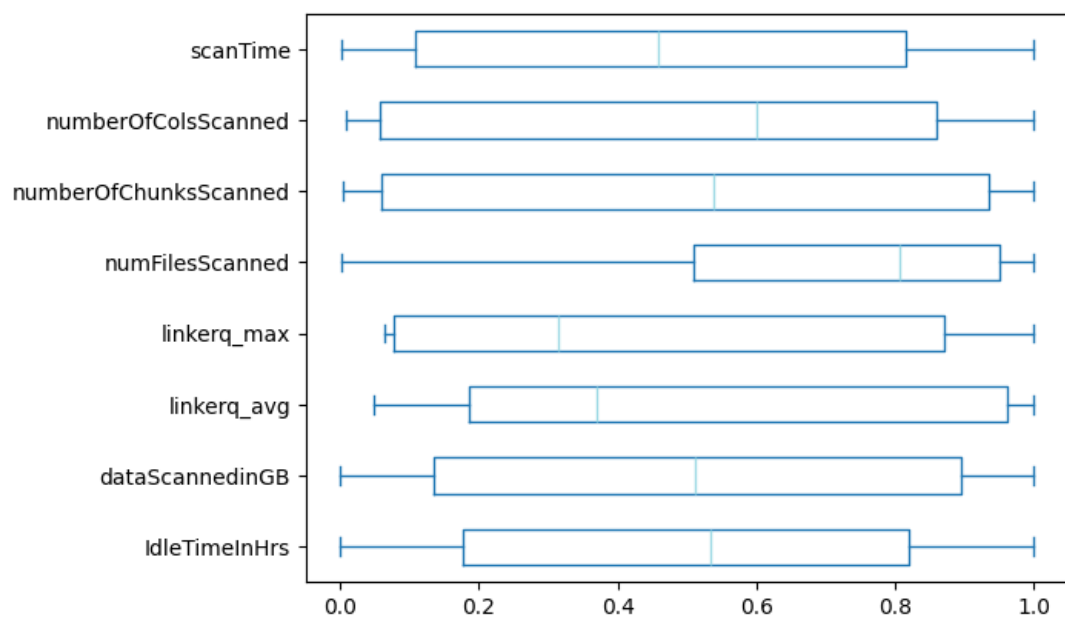
Absolute correlation vs percent-change of cpu_used_max
(For median correlation greater than 0.3)



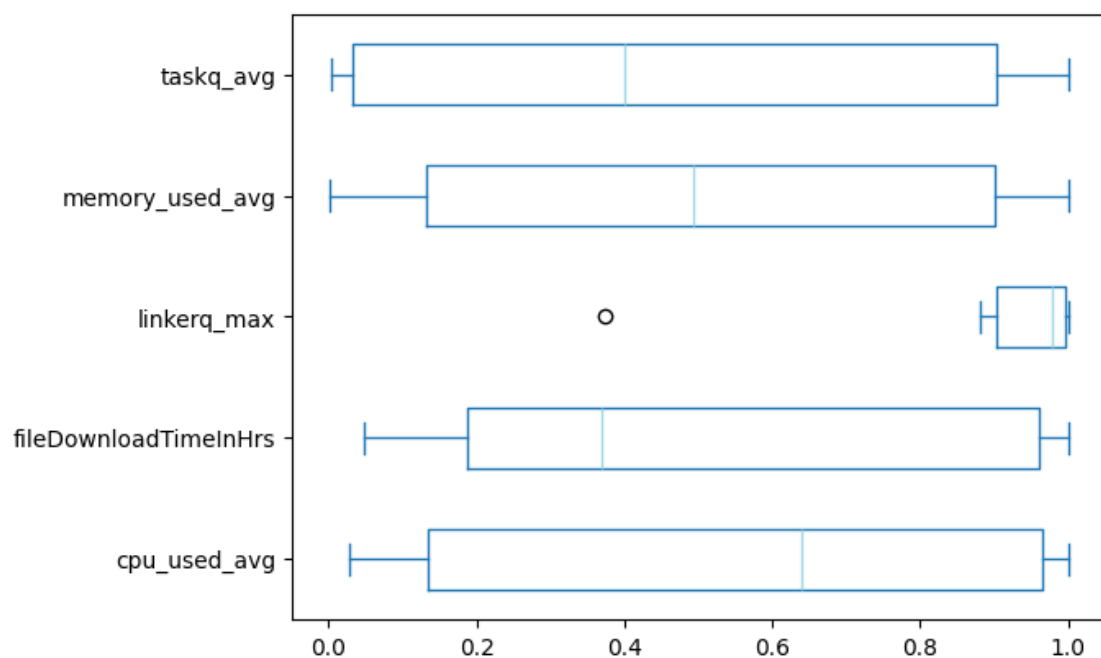
Absolute correlation vs percent-change of dataScannedinGB
(For median correlation greater than 0.3)



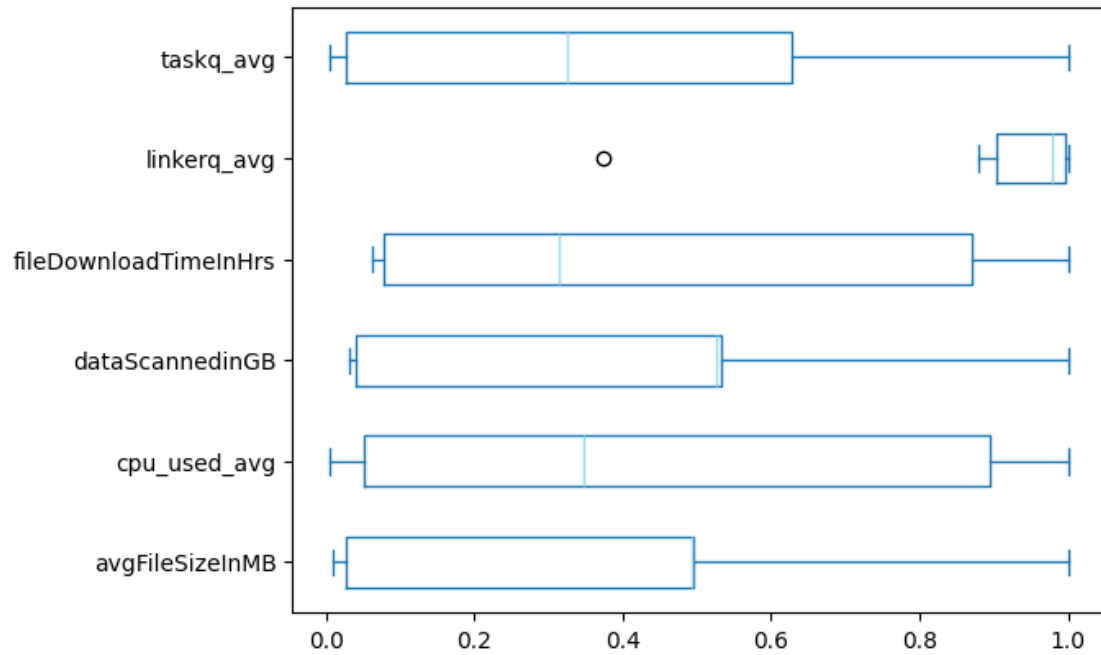
Absolute correlation vs percent-change of fileDownloadTimeInHrs
(For median correlation greater than 0.3)



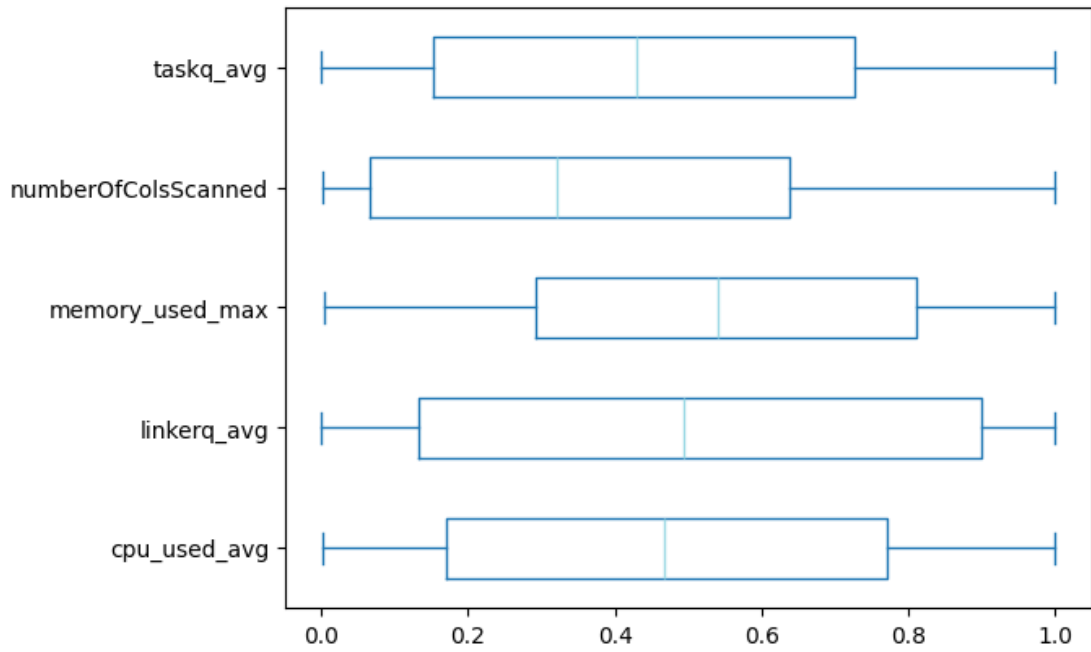
Absolute correlation vs percent-change of linkerq_avg
(For median correlation greater than 0.3)



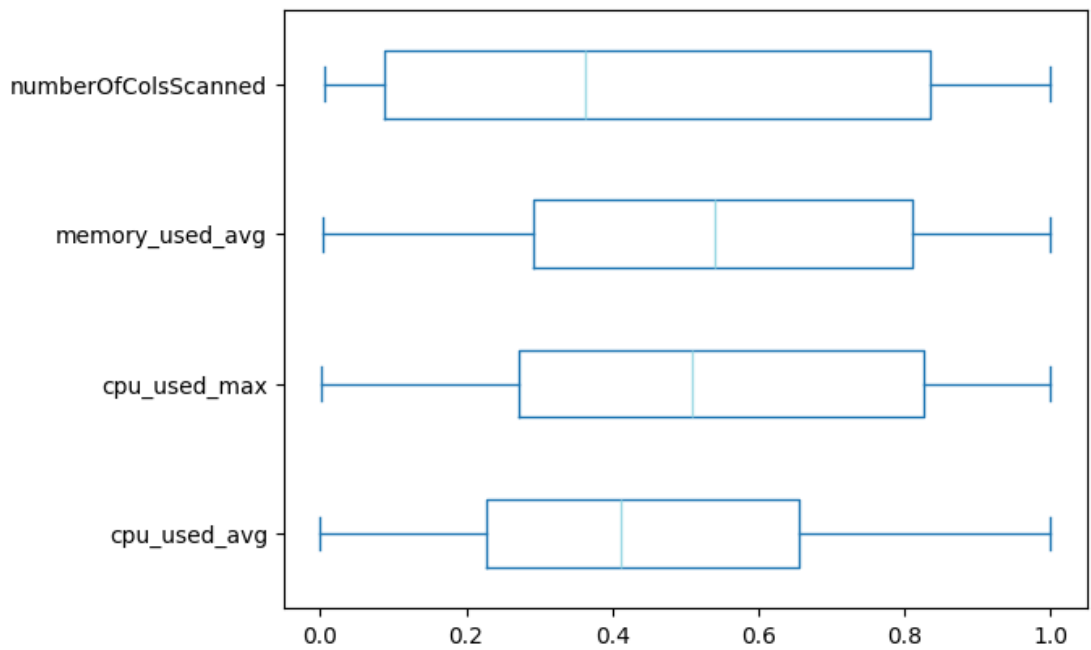
Absolute correlation vs percent-change of linkerq_max
(For median correlation greater than 0.3)



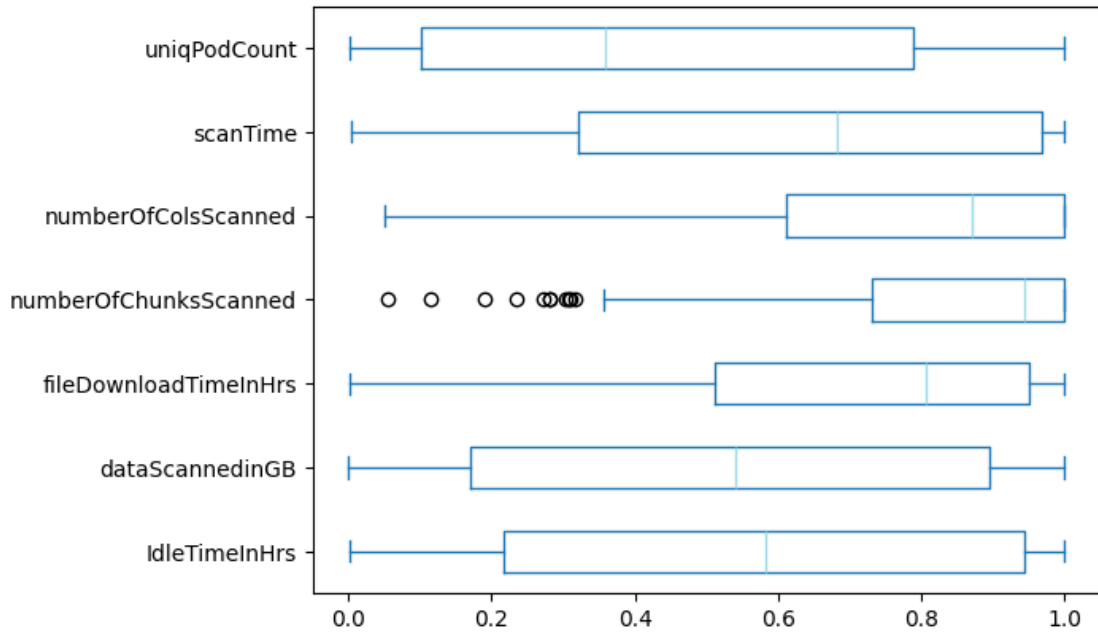
Absolute correlation vs percent-change of memory_used_avg
(For median correlation greater than 0.3)



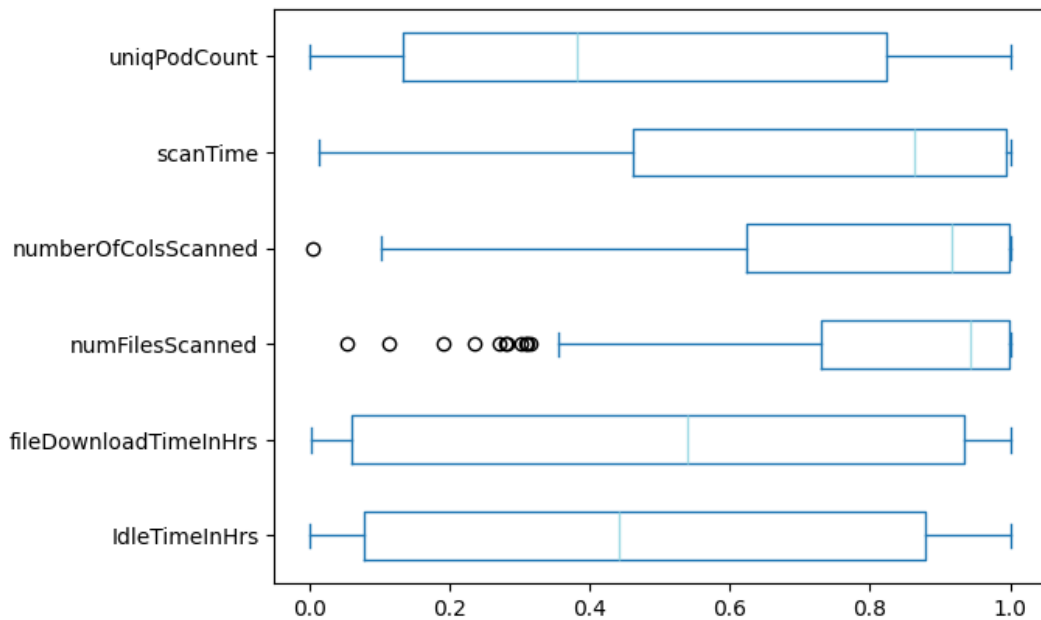
Absolute correlation vs percent-change of memory_used_max
(For median correlation greater than 0.3)



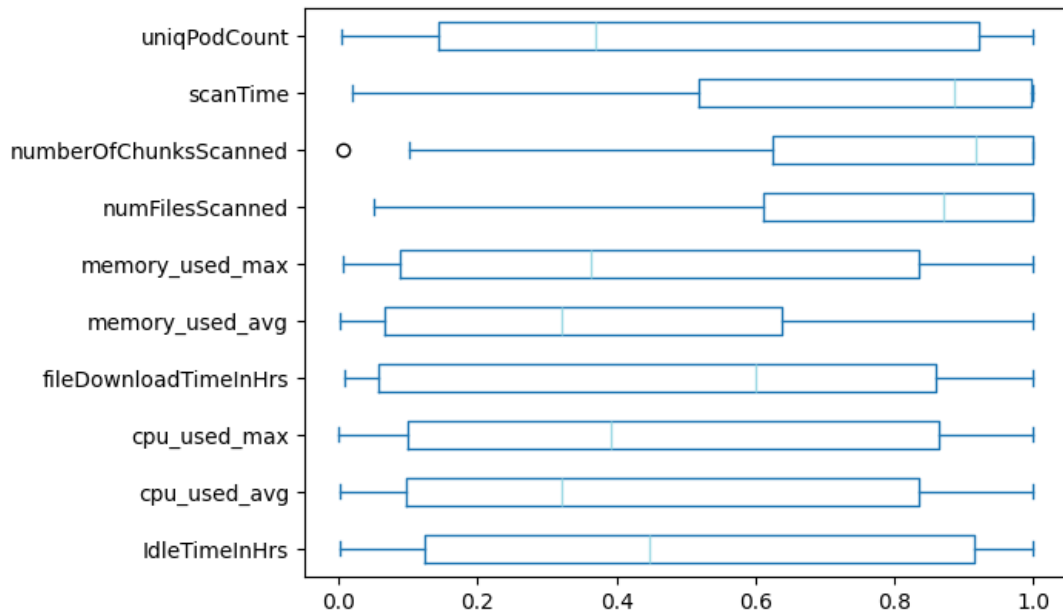
Absolute correlation vs percent-change of numFilesScanned
(For median correlation greater than 0.3)



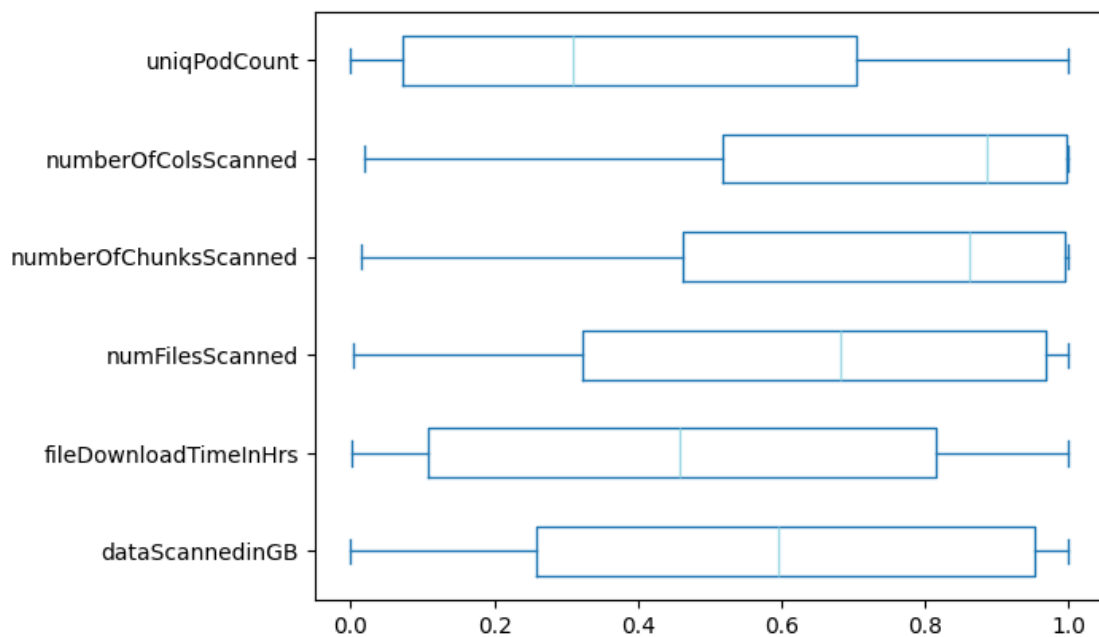
Absolute correlation vs percent-change of numberOfChunksScanned
(For median correlation greater than 0.3)



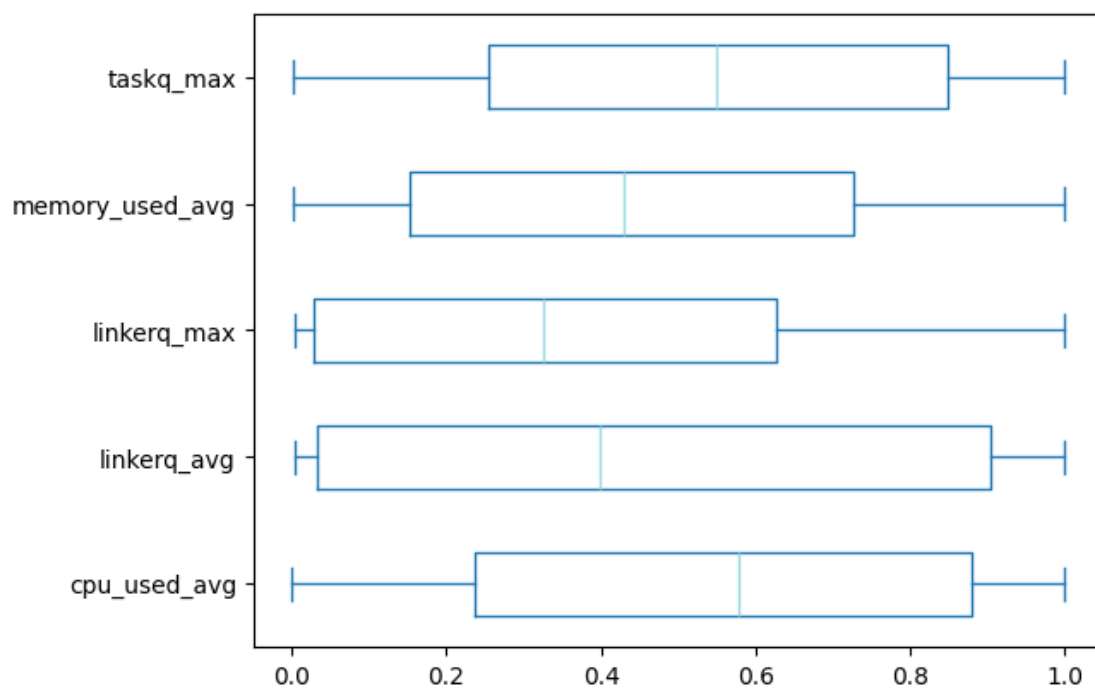
Absolute correlation vs percent-change of numberOfColsScanned
(For median correlation greater than 0.3)



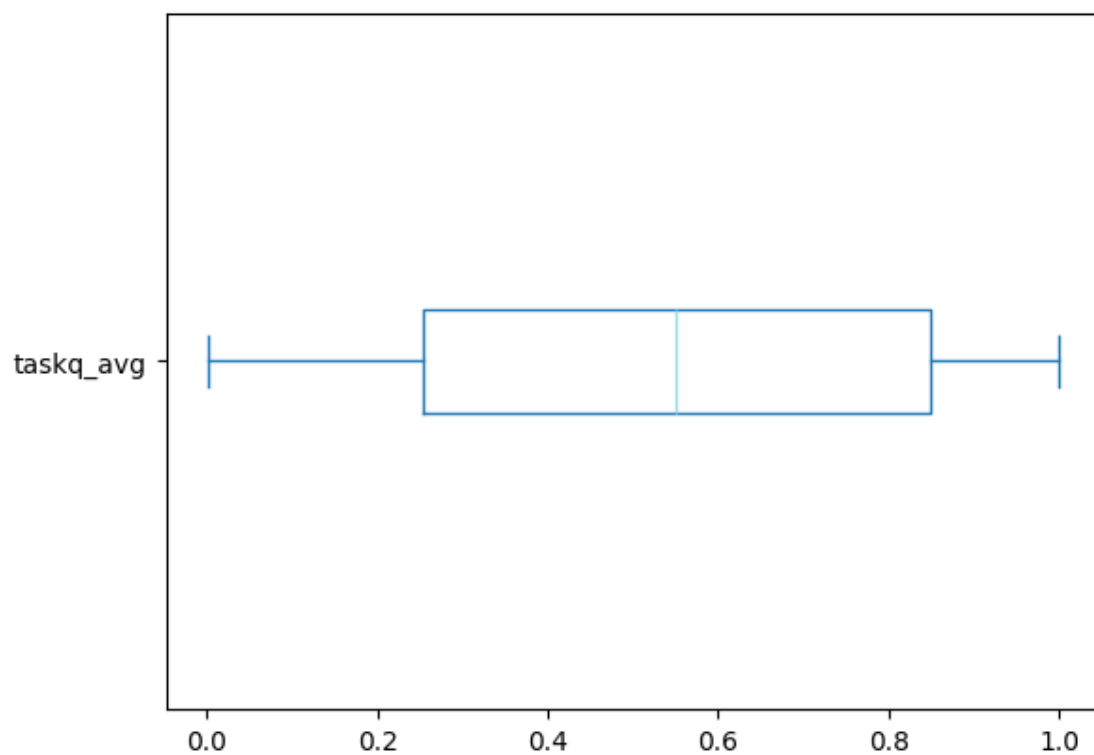
Absolute correlation vs percent-change of scanTime
(For median correlation greater than 0.3)



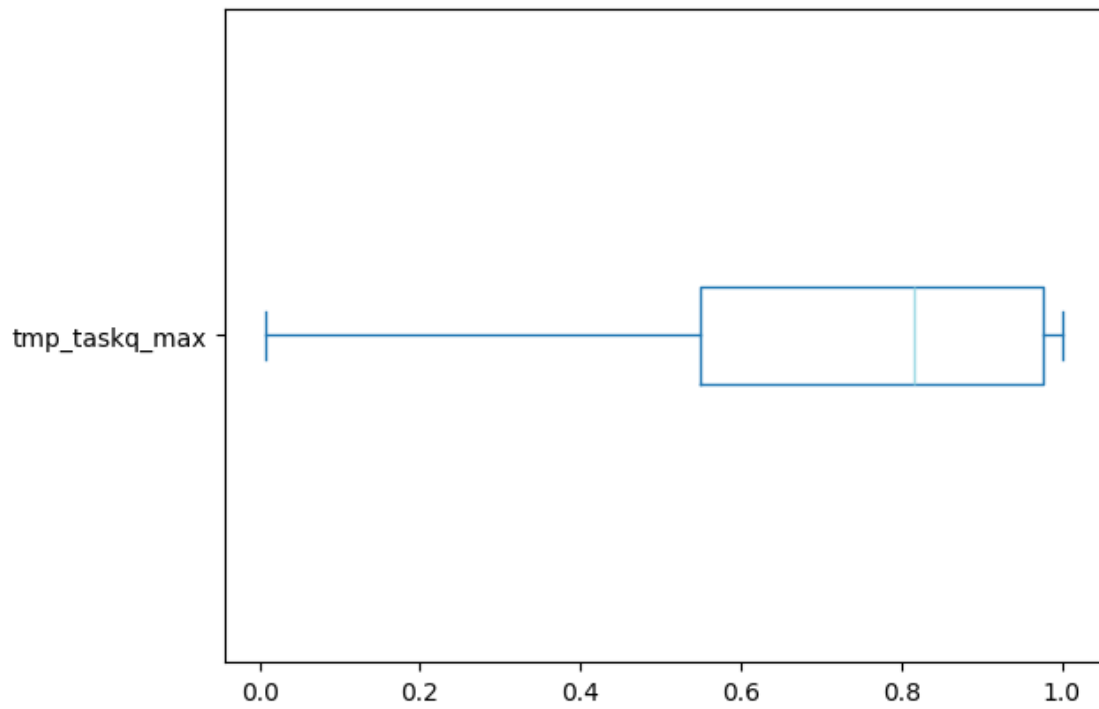
Absolute correlation vs percent-change of taskq_avg
(For median correlation greater than 0.3)



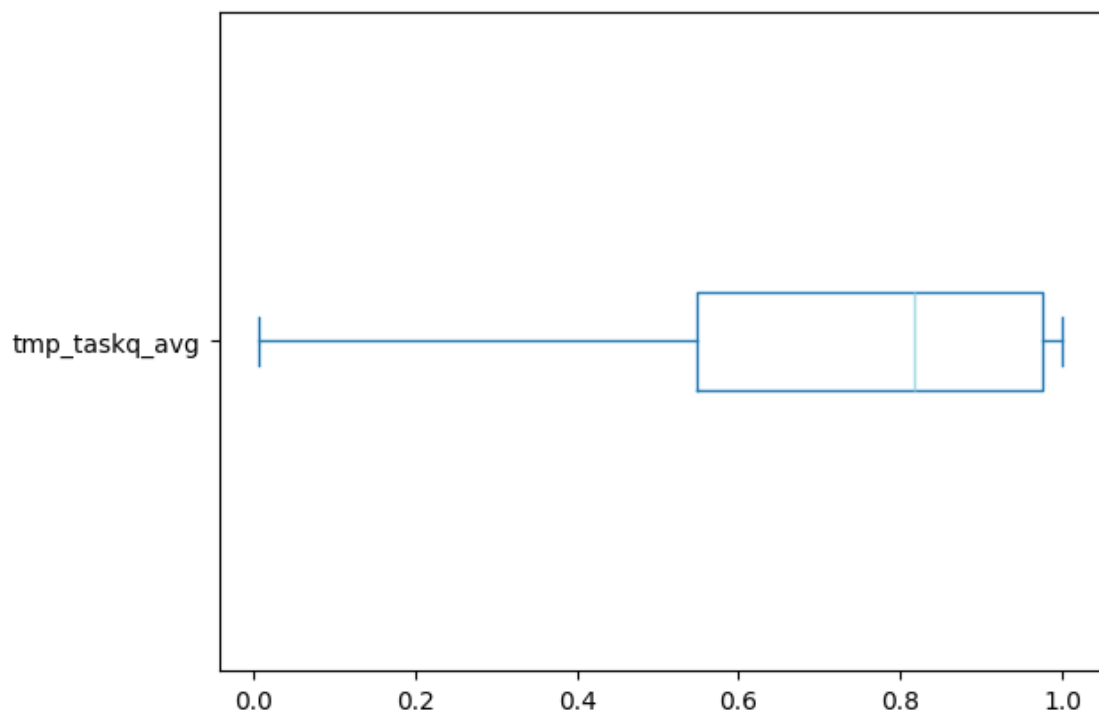
Absolute correlation vs percent-change of taskq_max
(For median correlation greater than 0.3)



Absolute correlation vs percent-change of tmp_taskq_avg
(For median correlation greater than 0.3)



Absolute correlation vs percent-change of tmp_taskq_max
(For median correlation greater than 0.3)



Absolute correlation vs percent-change of uniqPodCount
(For median correlation greater than 0.3)

