# Causal inference of multivariate server performance logs

Kewei Zhang

**2024/02/20**

# Introduction

- Data monitoring and visualization
    - Datadog

- Causal Inference
    - Data Imputation
    - Stationary Transformation
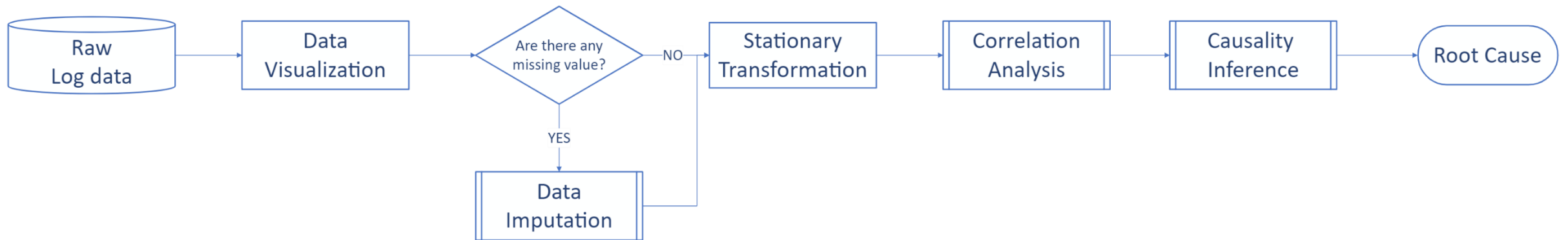    - Correlation Analysis
    - Causal and Inference

**INTERNAL USE**
Access limited to internal use only

KINAXIS®

# Data Monitoring and Visualization

Datadog

- [Logs Dashboard | Datadog (datadoghq.com)](datadoghq.com)

**INTERNAL USE**
Access limited to internal use only

KINAXIS®

# Causal Inference

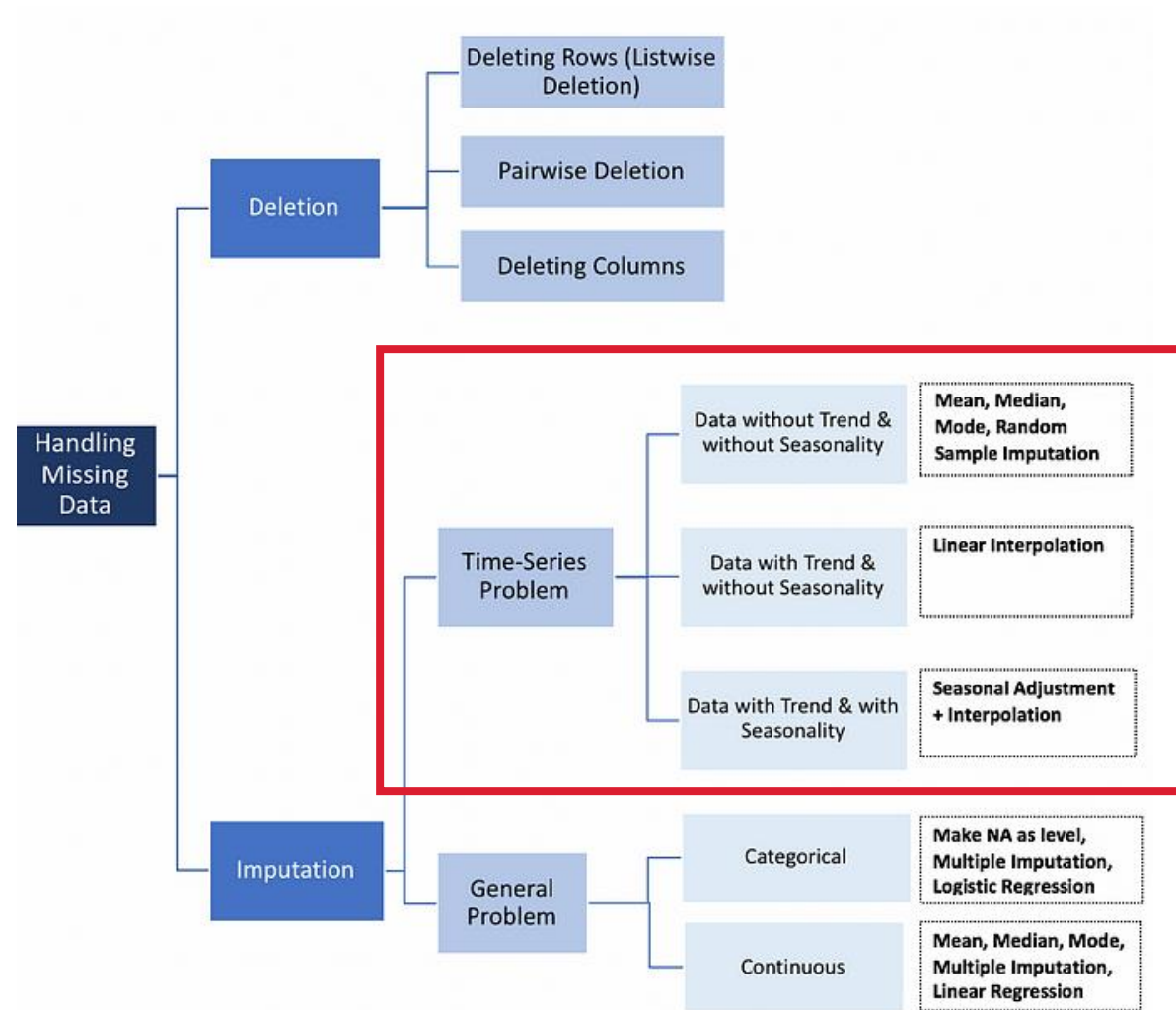Establish whether and how changes in one variable cause changes in another



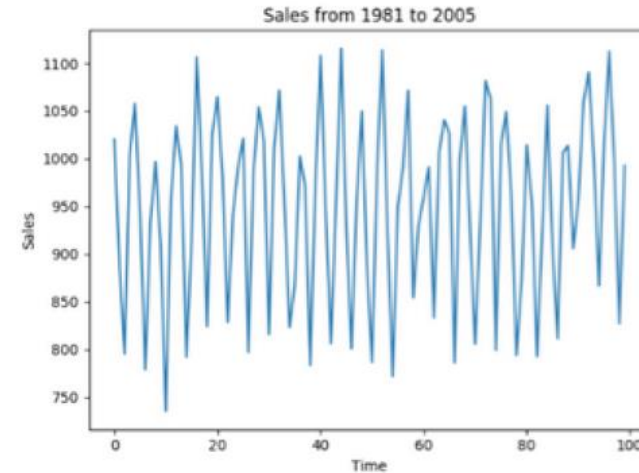Example: Processor Time – Performance Counter (Instance: RapidResponse)

**INTERNAL USE**
Access limited to internal use only

**KINAXIS®**

# Data Imputation

Imputation for missing value to generate complete datasets

**INTERNAL USE**
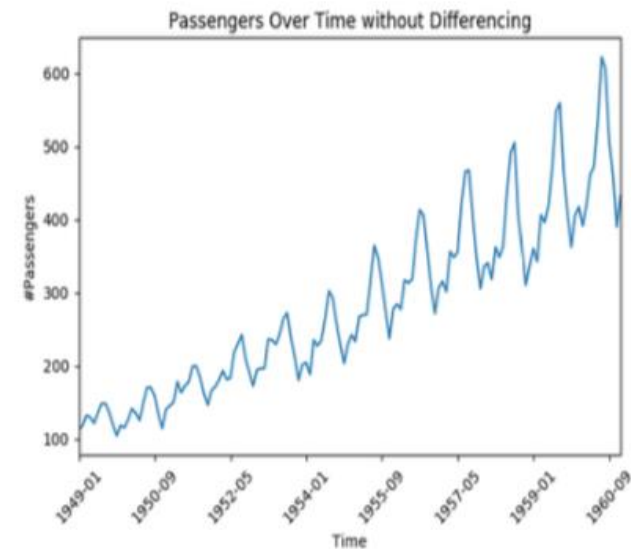Access limited to internal use only

# Stationary Transformation

Why stationary is important:
most time series models assume
that each point is independent
of one another.

**Stationary data:** mean and
variance do not vary across time



Stationary data



Non-stationary data

# Correlation Analysis

Evaluate strength and direction of the linear relationship between two variables

Question: We can only retain a month's worth of data, insufficient for any computation

Solution: Employ a variety of mathematical computations

Method:

1. Pearson correlation coefficient

2. Spearman correlation coefficient
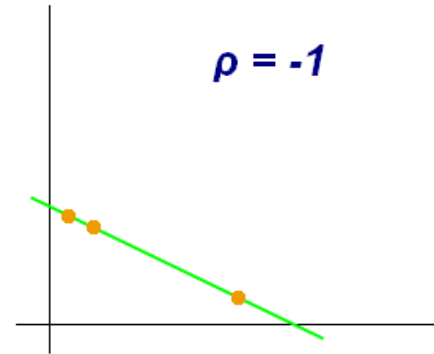
3. Kendall correlation coefficient

INTERNAL USE
Access limited to internal use only

# Correlation Analysis

**Pearson Correlation Coefficient**: [-1,1]
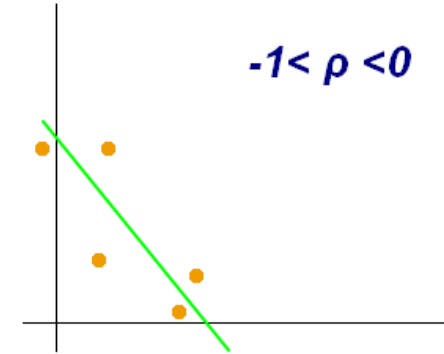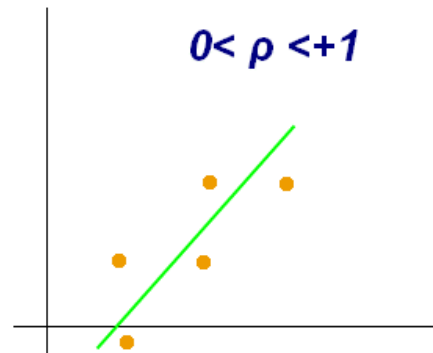
Coefficient = 1 or -1:  perfect linear relationship

Coefficient = 0 :  no linear relationship, could have a non-linear or more complex relationship

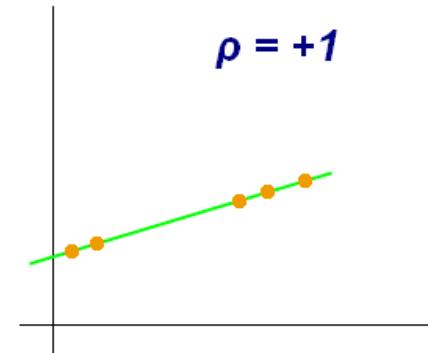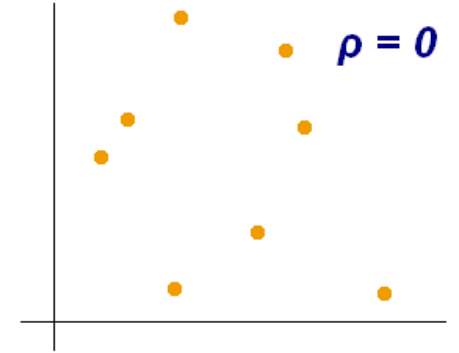$\rho = -1$

perfect negative correlation

$-1 < \rho < 0$

negative correlation

$0 < \rho < +1$

positive correlation

$\rho = +1$

perfect positive correlation

$\rho = 0$

no liner correlation

# Pearson Correlation Coefficient



Feature Importance:% Processor Time_RapidResponse

| Metrics | Pearson Linear Correlation Coefficient |
|---|---|
| Queries being executed_ | 0.7515 |
| Users logged in_ | 0.5399 |
| DctBlocks: create rate (blocks/sec)_ | 0.3946 |
| MC: Total Quota (MB)_ | 0.3368 |
| MC: Total InUse (MB)_ | 0.3308 |
| Queries open_ | 0.3083 |
| MC: Query Quota (MB)_ | 0.3058 |
| MC: Query InUse (MB)_ | 0.3058 |
| MC: Calculated Quota (MB)_ | 0.2342 |
| MC: Other Quota (MB)_ | 0.2225 |
| MC: Other InUse (MB)_ | 0.2225 |
| MC: Query Reclaim Count_ | 0.2132 |
| MC: Query Reclaim Actual (MB)_ | 0.2008 |
| Lock wait time (ms/sec)_ | 0.1901 |
| MC: Query Reclaim Target (MB)_ | 0.1803 |
| MC: Calculated InUse (MB)_ | 0.178 |
| MC: jemalloc Quota (MB)_ | 0.1644 |
| MC: jemalloc InUse (MB)_ | 0.1644 |
| MC: BigBlock Quota (MB)_ | 0.1637 |
| MC: BigBlock InUse (MB)_ | 0.1637 |
| Total query execution time_ | 0.1399 |
| Total data query executions_ | 0.0905 |
| MC: RecordBlock InUse (MB)_ | 0.0771 |
| MC: RecordBlock Quota (MB)_ | 0.0348 |
| MC: Input Quota (MB)_ | 0.0335 |
| MC: Input InUse (MB)_ | 0.0334 |
| MC: Hash Indexes InUse (MB)_ | 0.0276 |
| MC: Hash Indexes Quota (MB)_ | 0.0274 |
| % Processor Time_ServiceHost | 0.0005 |
| DctBlocks: read rate (blocks/sec)_ | -0.02 |
| DctBlocks: create memory usage rate (bytes/sec)_ | -0.0347 |

**INTERNAL USE**

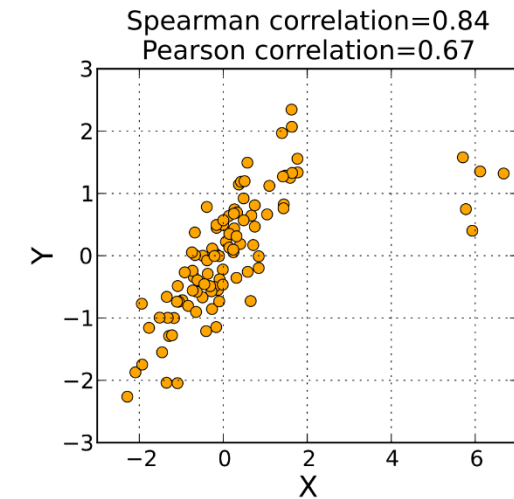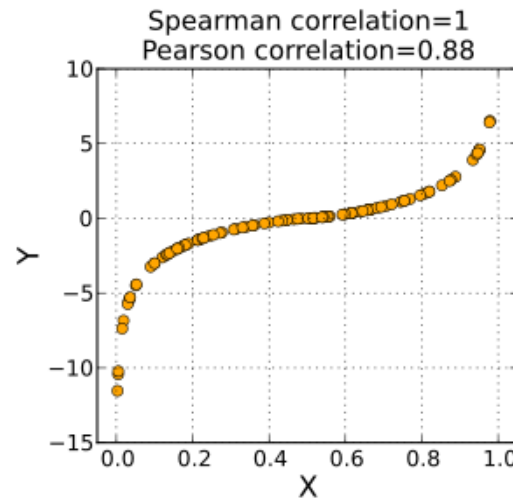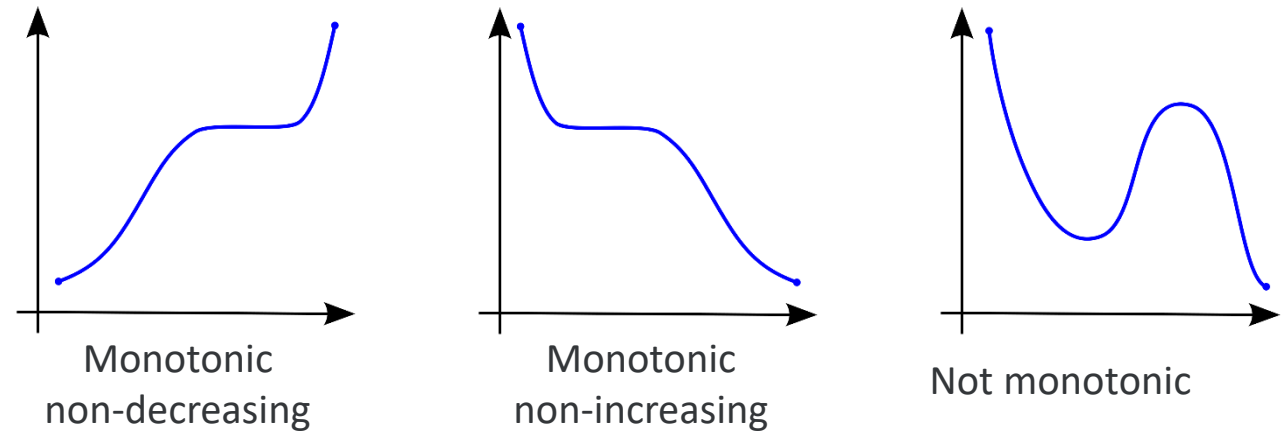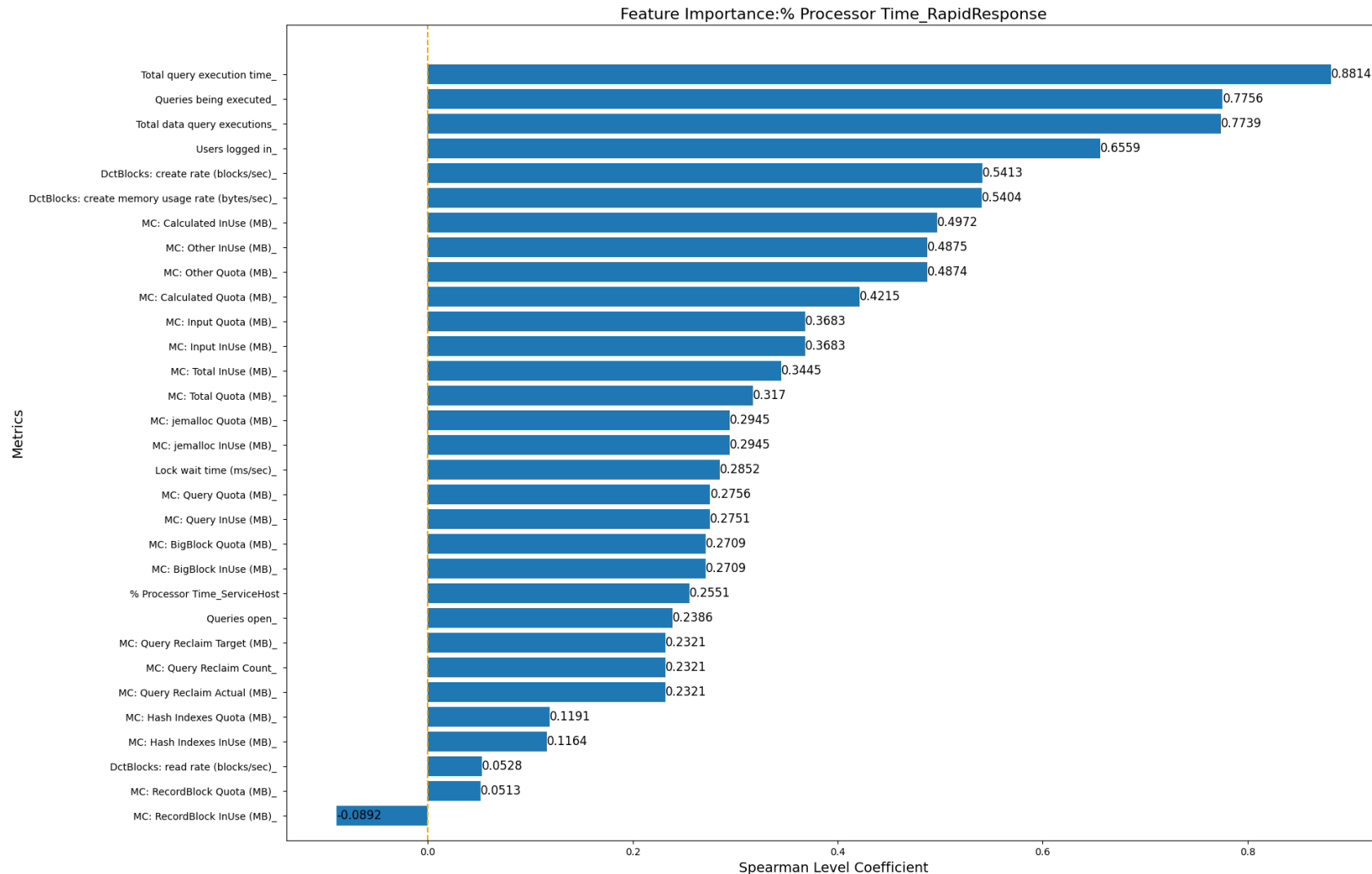Access limited to internal use only

KINAXIS®

# Correlation Analysis

**Spearman Correlation Coefficient**

coefficient = 1/-1 : perfect positive/negative monotonic relationship

Coefficient = 0 : no monotonic relationship



Monotonic non-decreasing

Monotonic non-increasing

Not monotonic

Spearman correlation=1
Pearson correlation=0.88

Spearman correlation=0.84
Pearson correlation=0.67

# Spearman Correlation Coefficient



Feature Importance:% Processor Time_RapidResponse

Copyright © 2023 Kinaxis. All Rights Reserved.

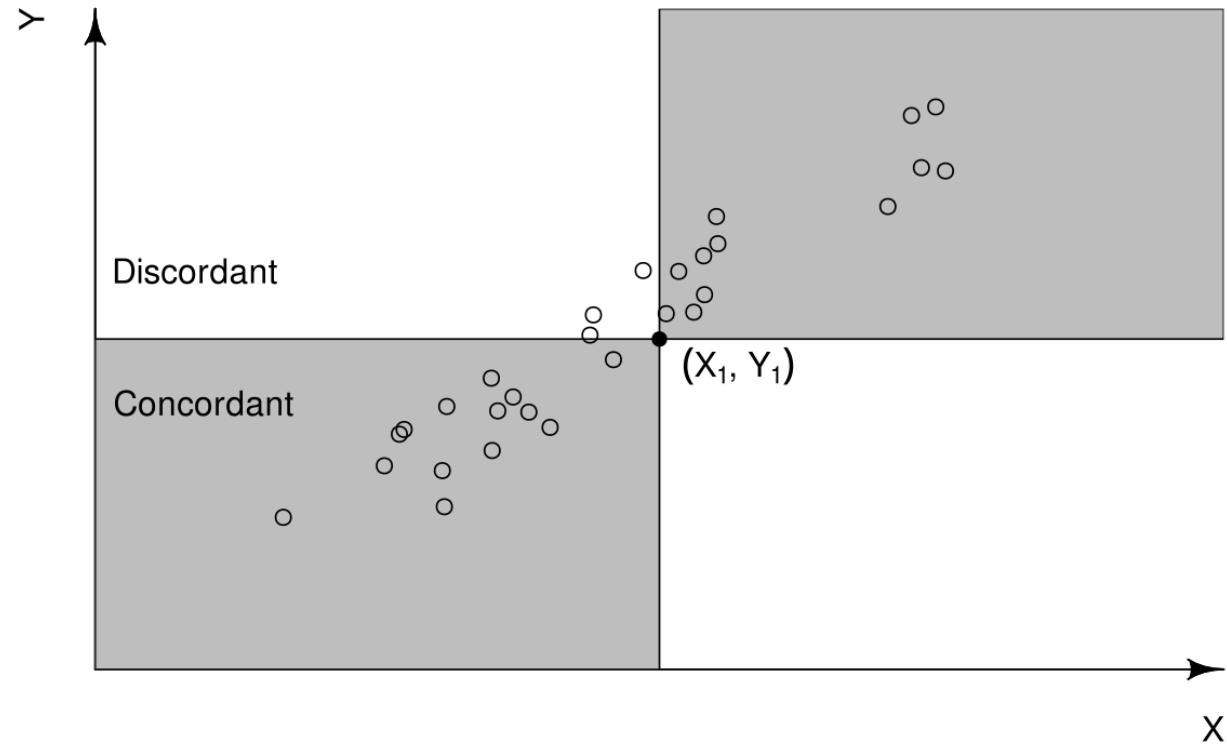**INTERNAL USE**
Access limited to internal use only

# Kendall Correlation Coefficient

coefficient = 1/-1: ranks of corresponding values within each data sample always same
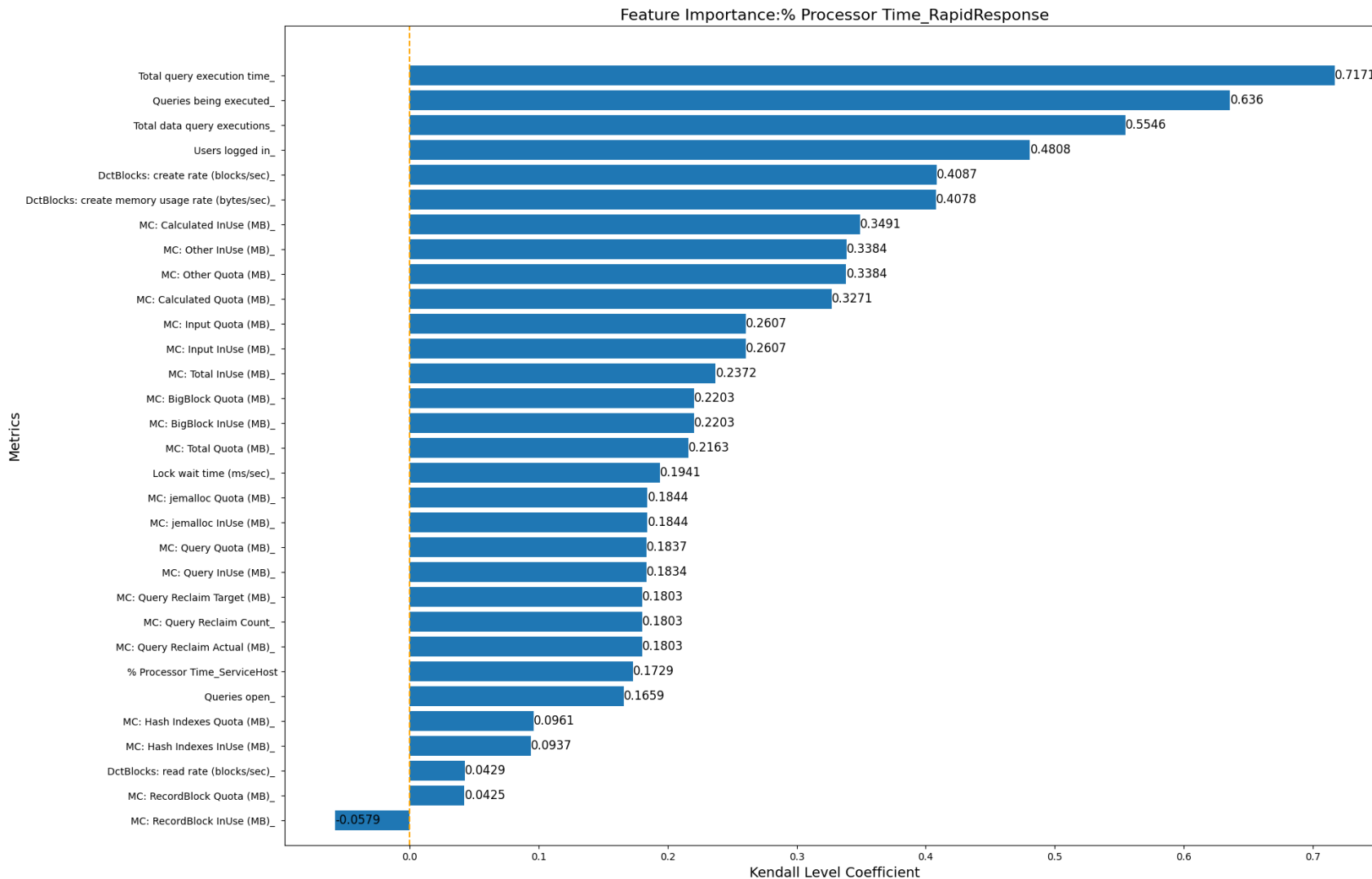
Coefficient = 0 : no association between the ranks of the values

Ex: (xi,yi), (xj,yj), i < j:

Either both xi > xj and yi > yj or xi < xj, yi < yj

# Kendall Correlation Coefficient



Feature Importance:% Processor Time_RapidResponse

# Correlation Coefficient

Top5 Results:

| Pearson Correlation: | Spearman Correlation: | Kendall Correlation: |
|---|---|---|
| Queries being executed | Total query execution time | Total query execution time |
| Users logged in | Queries being executed | Queries being executed |
| DctBlocks: create rate (blocks/sec) | Total data query execution | Total data query execution |
| MC: Total Quota (MB) | Users logged in | Users logged in |
| MC: Total Inuse (MB) | DctBlocks: create rate (blocks/sec) | DctBlocks: create rate (blocks/sec) |

**INTERNAL USE**
Access limited to internal use only

**KINAXIS®**

# Correlation VS Causation

**INTERNAL USE**

Access limited to internal use only

# Transfer Entropy

Features of time series: max, min, mean, median, variance,...

Series A : 1, 2, 1, 2, 1, 2, 1, ...

Series B : 1, 1, 2, 1, 2, 2, 1, ...

Same mean, same variance, same median, <span style="color:red">different entropy</span>

Larger the entropy, more chaotic the system
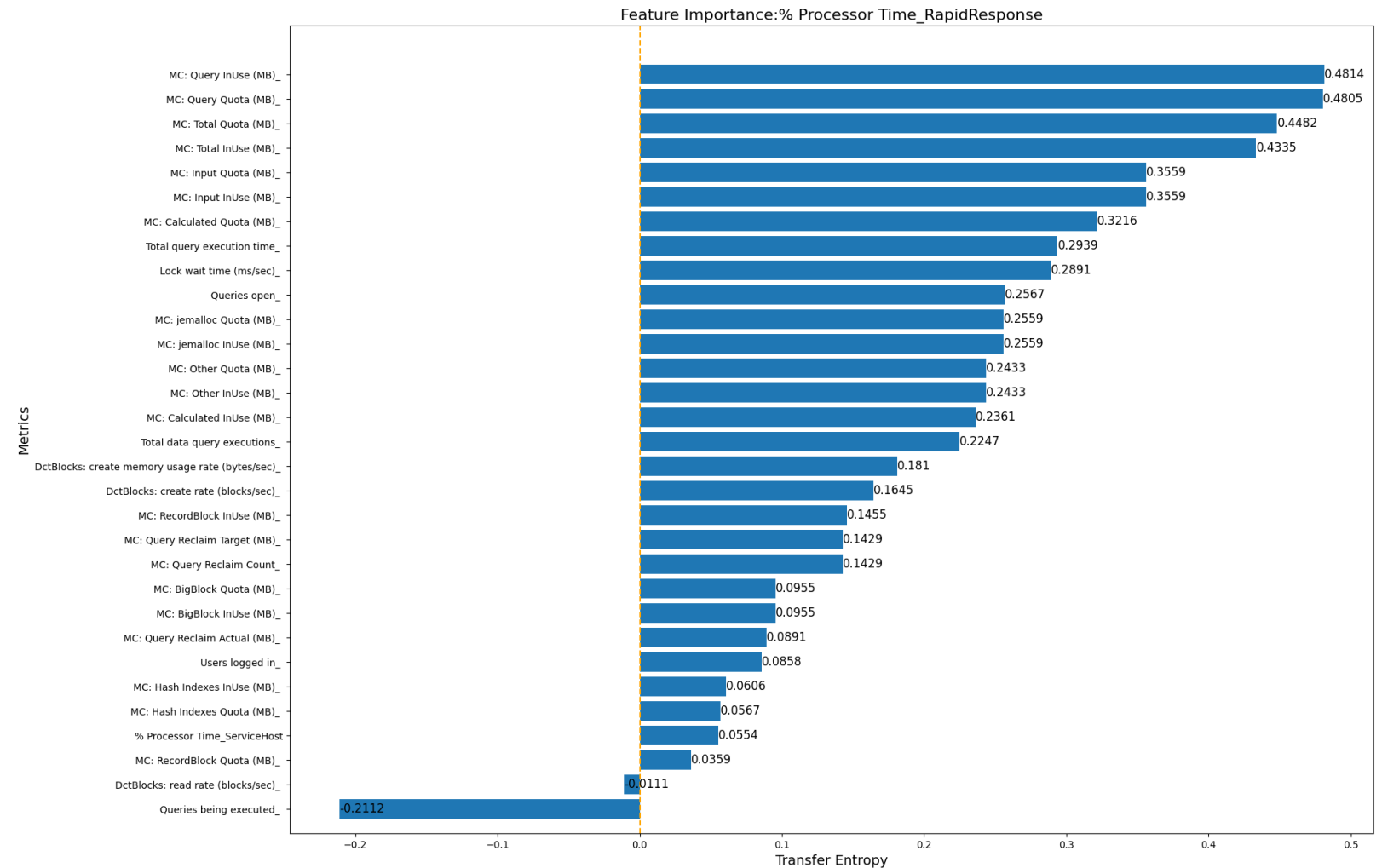
Transfer entropy: transfer of information

INTERNAL USE
Access limited to internal use only

kinaxis®

# Causal Inference

**Transfer Entropy**

Symbol: direction of transfer

+:Y -> X

-: X -> Y



Feature Importance:% Processor Time_RapidResponse

Copyright © 2023 Kinaxis. All Rights Reserved.

# Causal Inference

## Causal Inference

**Neural Network: LSTM(Long-short term memory)**

Step 1: Forecasting

Step 2: Evaluating the prediction model

Step 3: Calculating feature Importance

**INTERNAL USE**
Access limited to internal use only

**KINAXIS®**

# Causal Inference

## Step1: Forecasting

Selected 3 subsections from the
entire dataset as examples



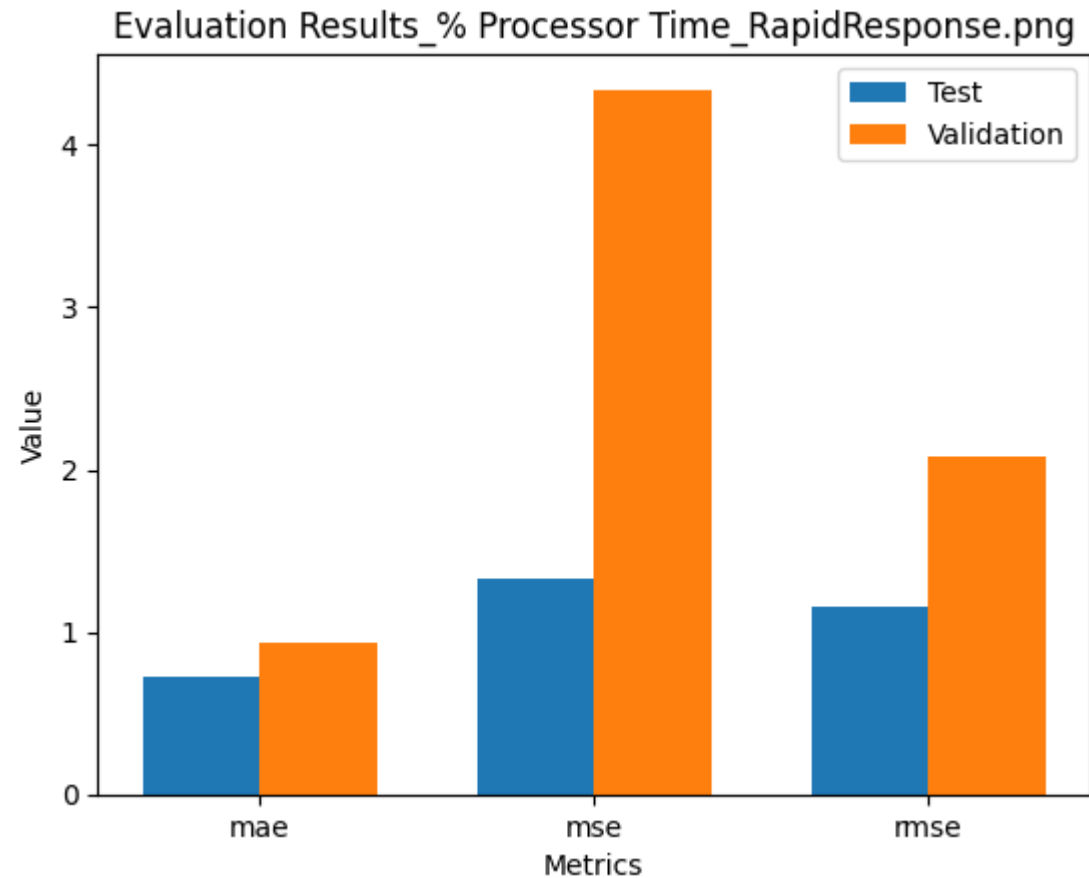Copyright © 2023 Kinaxis. All Rights Reserved.

# Causal Inference

Step 2: Evaluating the

prediction model

MAE: mean absolute error

MSE: mean squared error
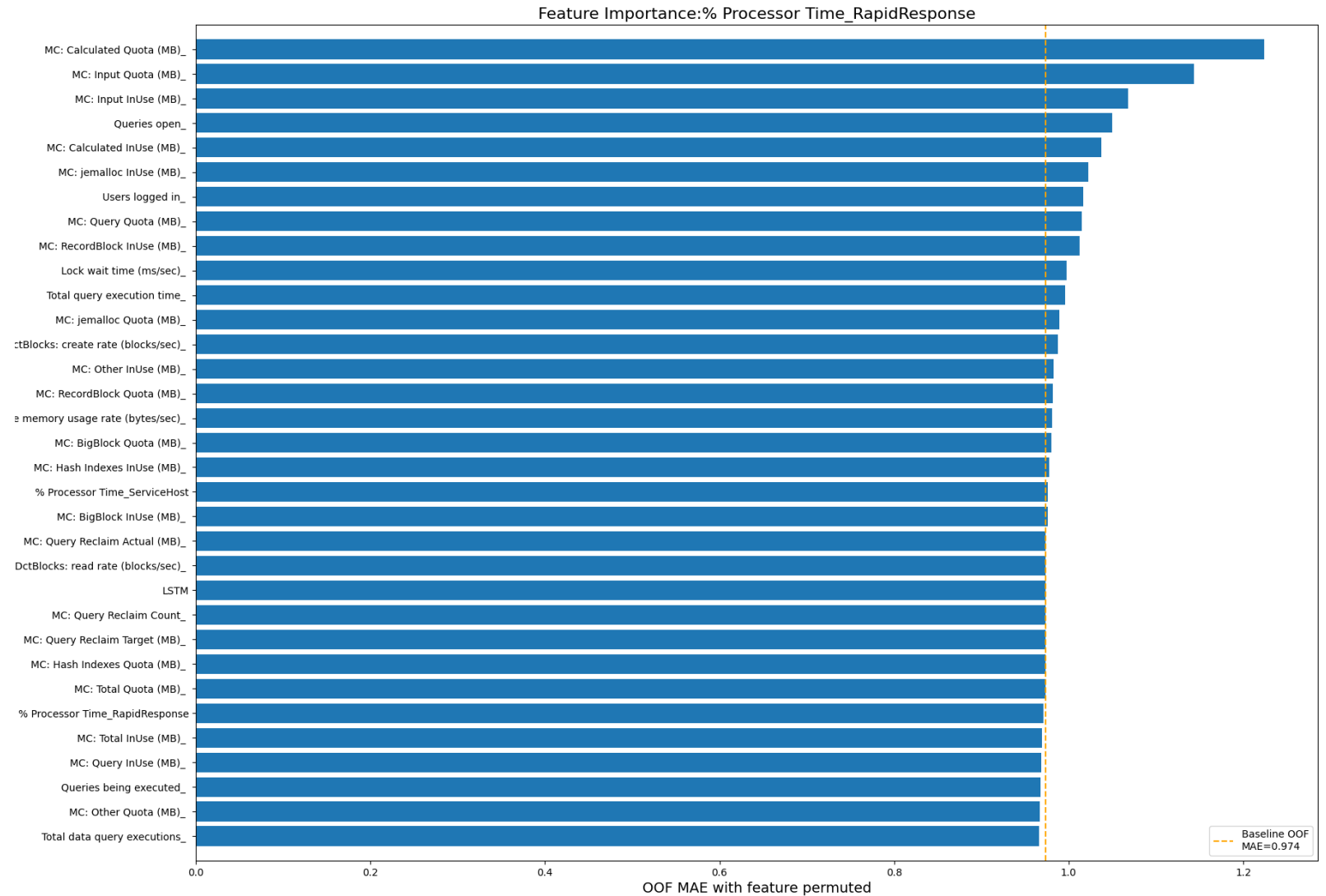
RMSE: root mean squared error

# Causal Inference

## Causal Inference

Step 3: Calculating feature Importance

OOF: Out-of-Fold, dividing a dataset into subsets, known as "folds." Each fold is then used as a validation set once while the remaining folds are utilized for training.



Feature Importance:% Processor Time_RapidResponse

**INTERNAL USE**
Access limited to internal use only

# Conclusion

- Most Correlated metrics with Processor Time:
  - Total query execution time
  - Queries being executed
  - Total data query execution
  - users logged in
  - DctBlocks: create rate (blocks/sec)

- Metrics with the largest Transfer Entropy:
  - MC metrics
  - Total query execution time
  - Lock wait time
  - Queries open
  - Total data query execution

- Metrics with the biggest influence:
  - MC metrics
  - Queries open
  - Users logged in
  - Lock wait time
  - Total query execution time

# Conclusion

Possible Reason of error:

- Insufficient data

- Inappropriate model selection

  - Not suitable

  - Overfitting, too complex

- Feature engineering issues

  - Incorrect stationary transformation

**INTERNAL USE**
Access limited to internal use only

**KINAXIS®**