



LUND
UNIVERSITY

Data Driven Methods for Train Delay Prediction

TIONG KAH YONG, LUND UNIVERSITY



Content

Introduction

Research design

Answering research questions

Conclusion and future research

Introduction



- Train delay is defined as the *deviation of actual train events from scheduled train events.*

- Train delay is defined as the *deviation of actual train events from scheduled train events*.
- *High-capacity utilisation* and *heterogeneous traffic* make the railway network susceptible to delays.



Train delay prediction model

- *Predicting* the expected train traffic conditions at a future time



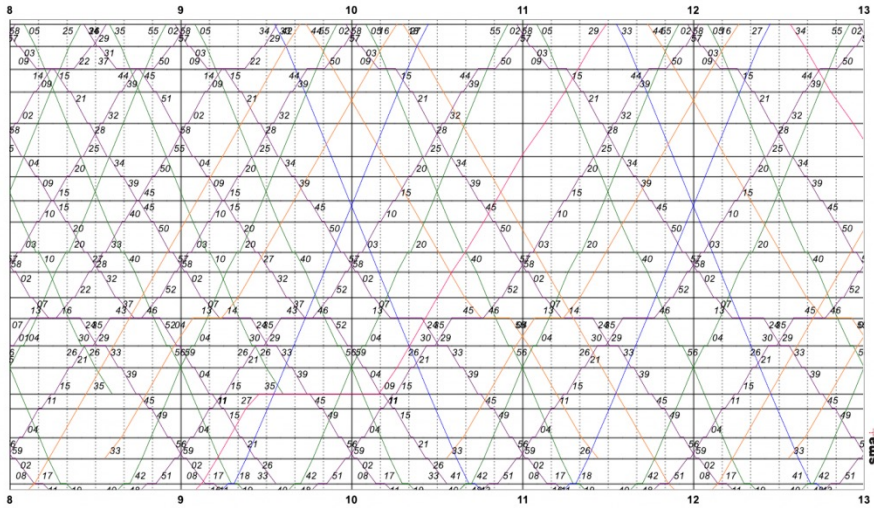
Train delay prediction model

- *Predicting* the expected train traffic conditions at a future time.
- *Input* for solving many problems related to train traffic management.

Train delay prediction model



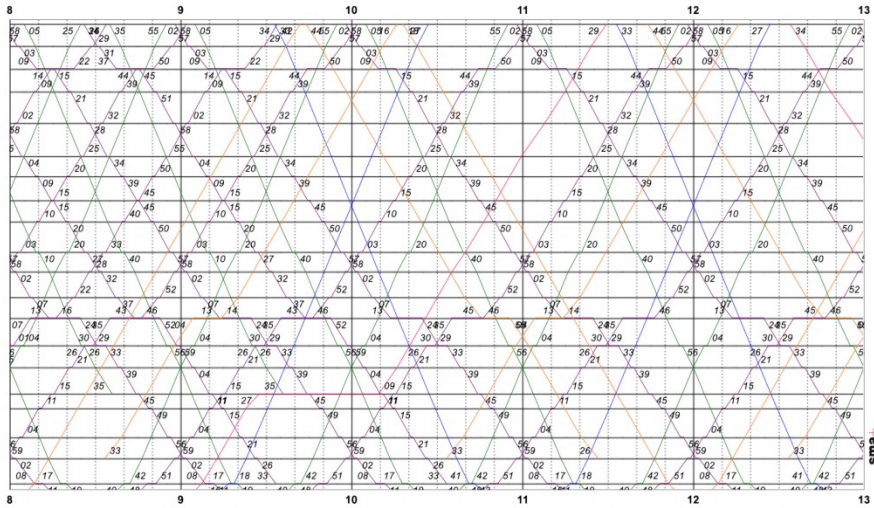
Train delay prediction model



Timetable planning



Train delay prediction model



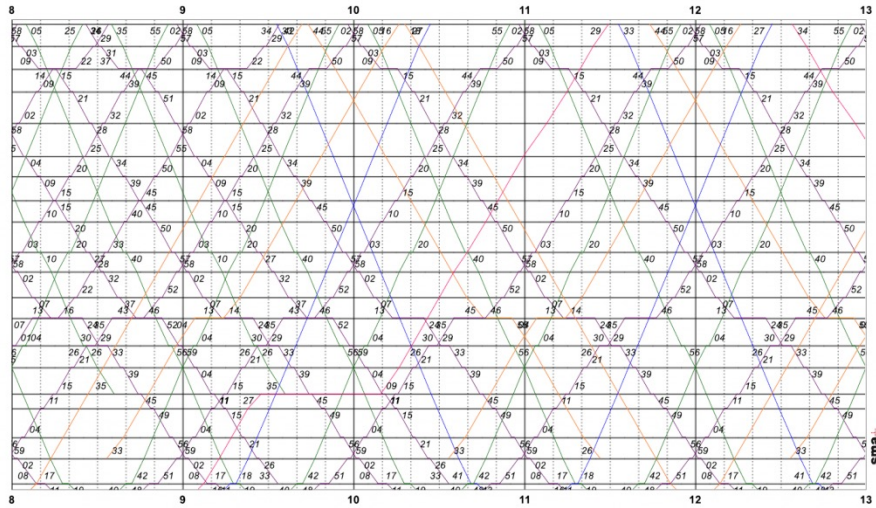
Timetable planning



Real-time train management



Train delay prediction model



Timetable planning



Real-time train management



Reliable passenger information system

Research design



Research gap



Research gap

- *Insufficient understanding* of existing train delay prediction models.

Research gap

- *Insufficient understanding* of existing train delay prediction models.
- There is a lack of innovation in developing train delay prediction models with *practical applications*.

Research aim

To increase understanding of data-driven train
delay prediction models

Research questions, papers and their connections

Research question 1

What factors need to be taken into account when building a train delay prediction model?

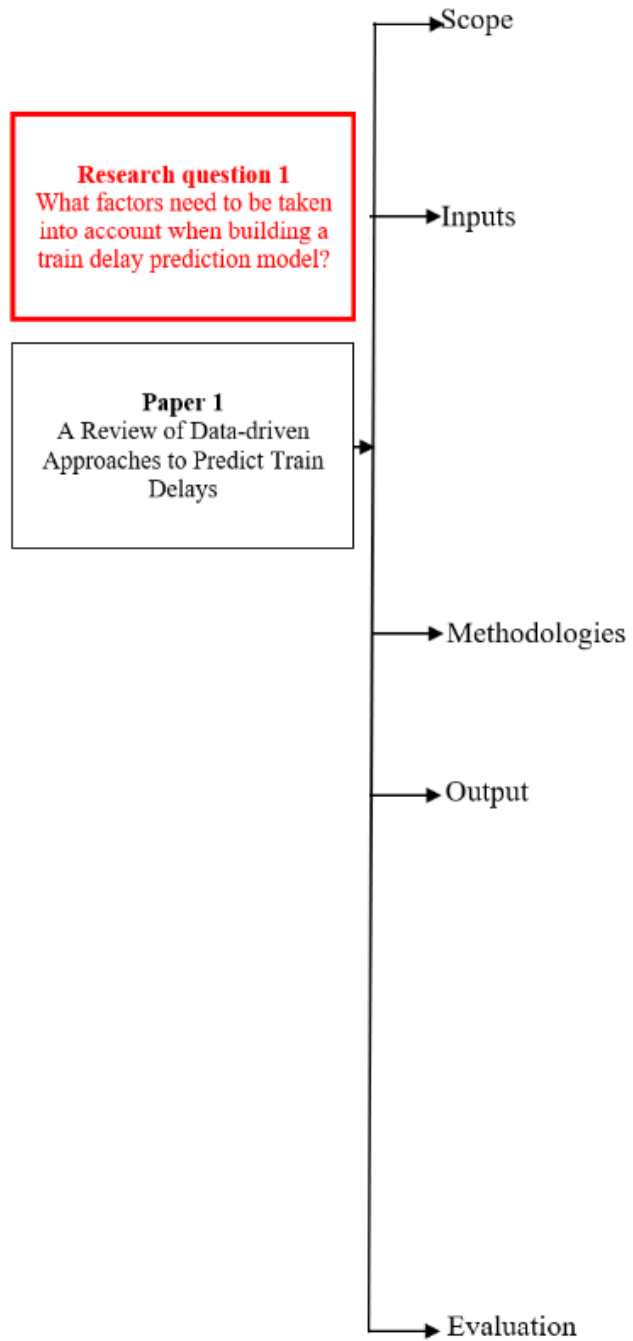
To increase *our understanding of the various aspects* that must be considered.

Research question 1

What factors need to be taken into account when building a train delay prediction model?

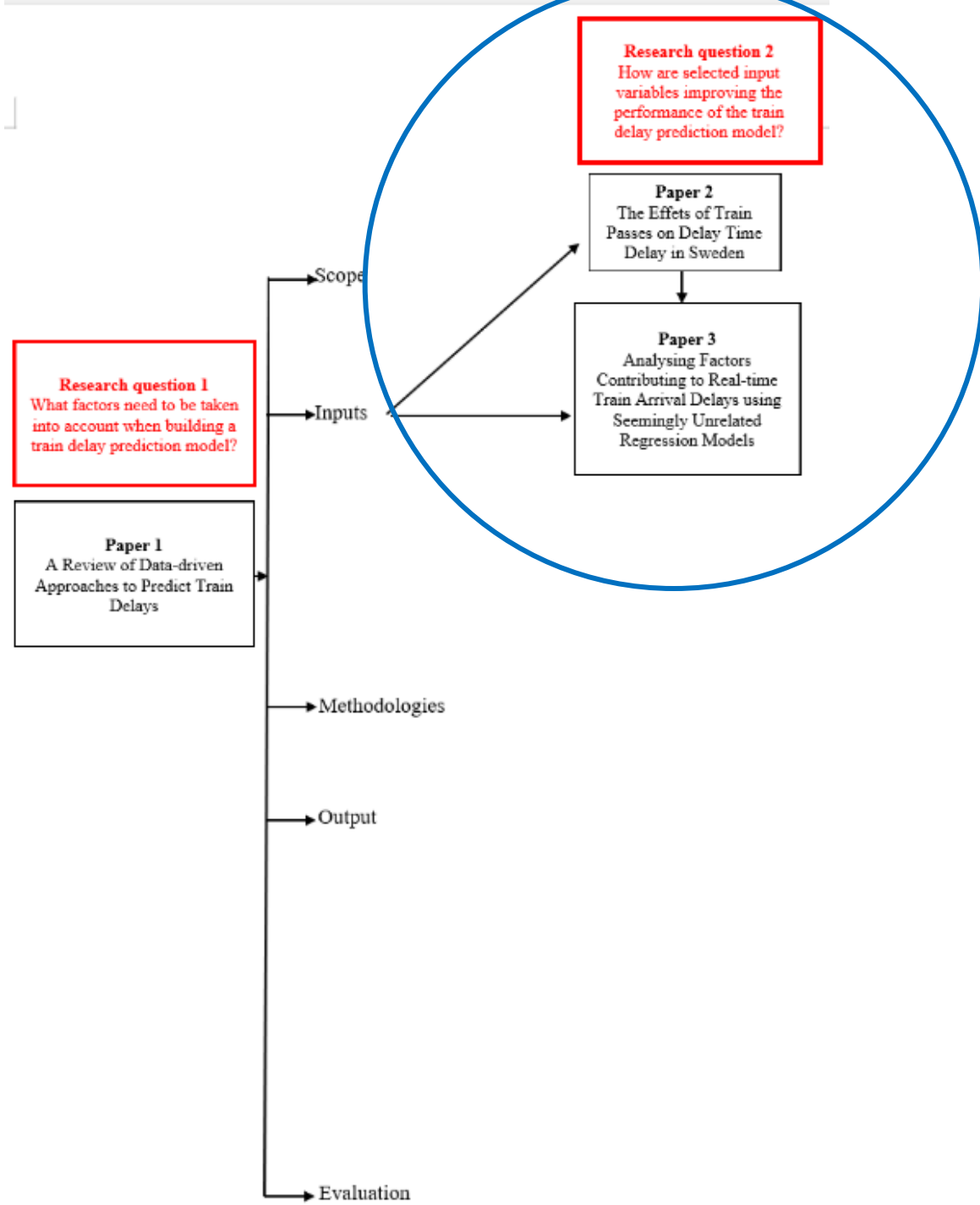
Paper 1

A Review of Data-driven Approaches to Predict Train Delays



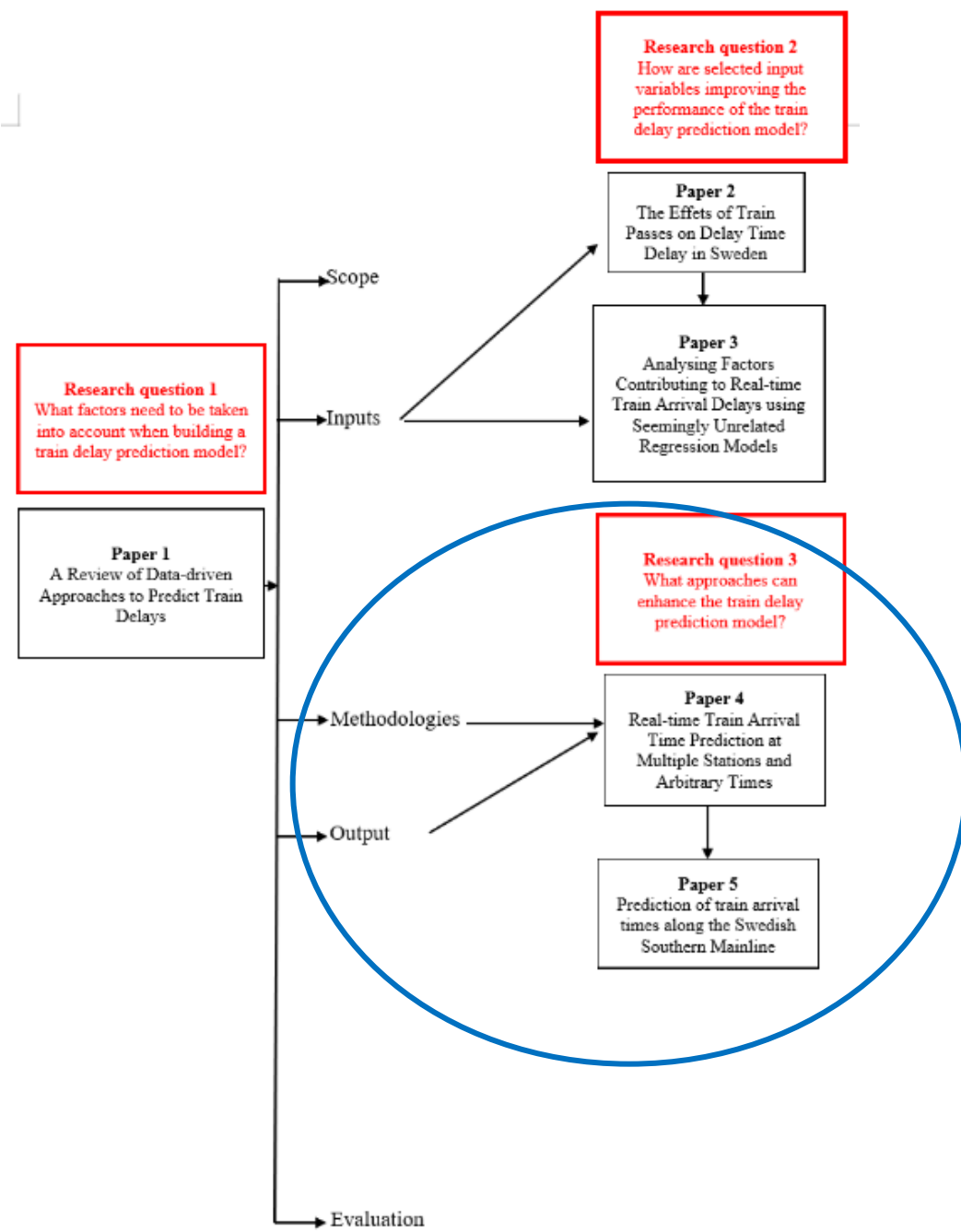
Research question 2
How are selected input variables improving the performance of the train delay prediction model?

To *identify useful input variables* to enhance model performance



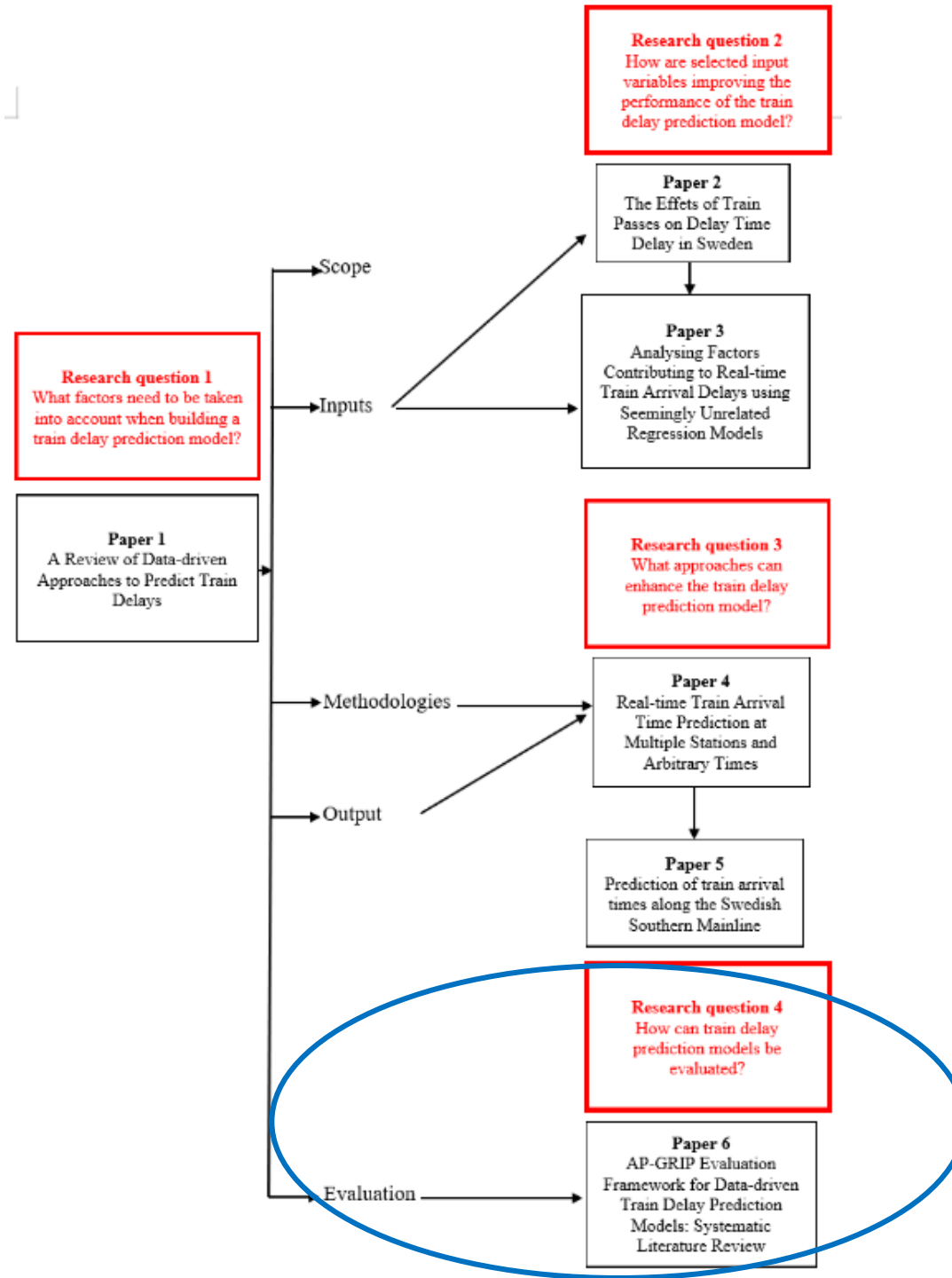
Research question 3
What approaches can
enhance the train delay
prediction model?

The *formulation of innovative technical solutions* to address
the current modelling challenges



Research question 4
How can train delay
prediction models be
evaluated?

To **thorough assessment** of train delay prediction
models

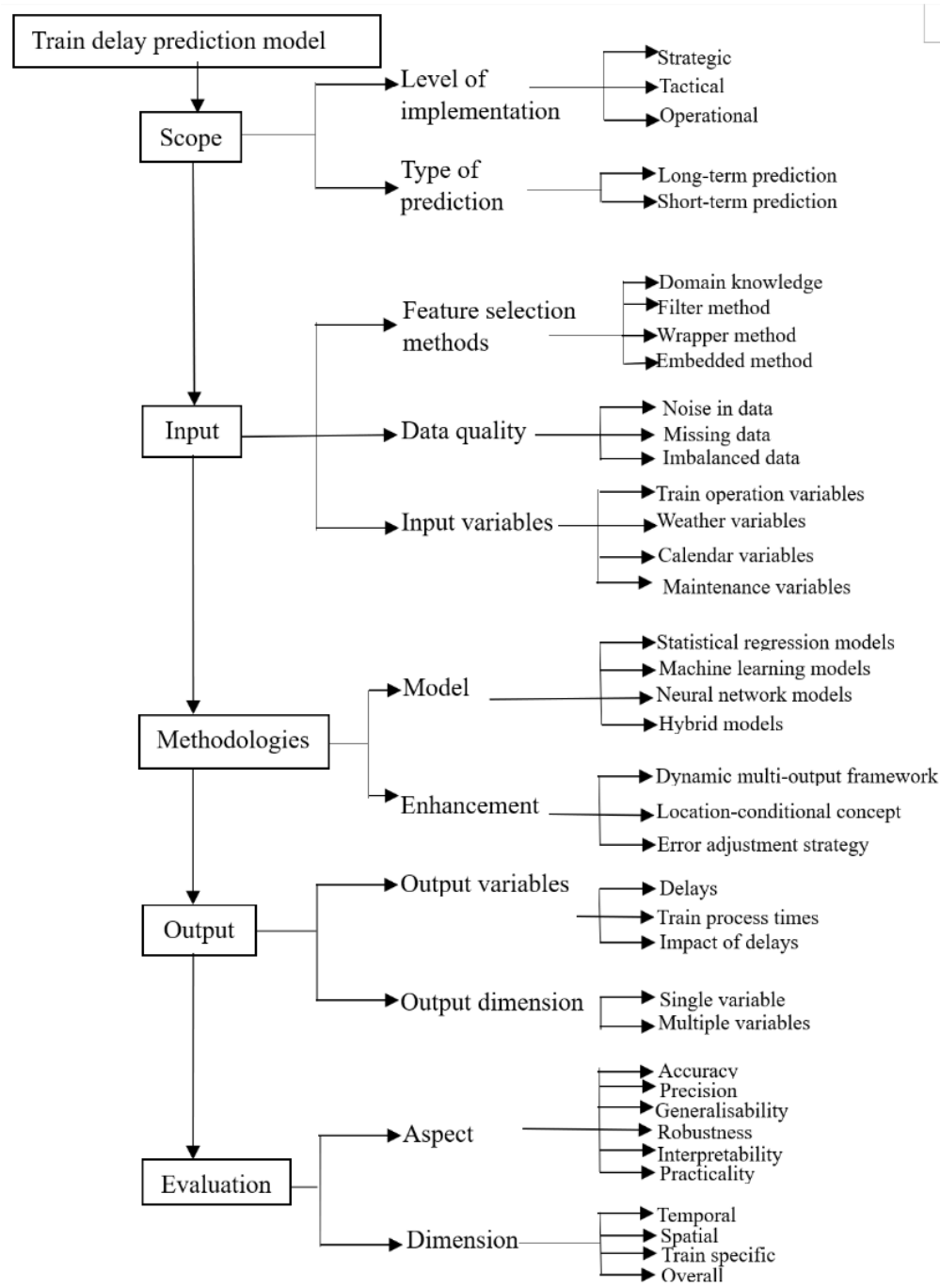


Answering research questions

Research question 1

What factors need to be taken into account when building a train delay prediction model?





Scope

Long-term prediction models

Short-term prediction models

Scope

Long-term prediction models

- Study how different factors affect train delays
- Use historical data
- Predict delays several days or months in advance
- For both strategic and tactical train traffic planning

Short-term prediction models



Scope

Long-term prediction models

- Study how different factors affect train delays
- Use historical data
- Predict delays several days or months in advance
- For both strategic and tactical train traffic planning

Short-term prediction models

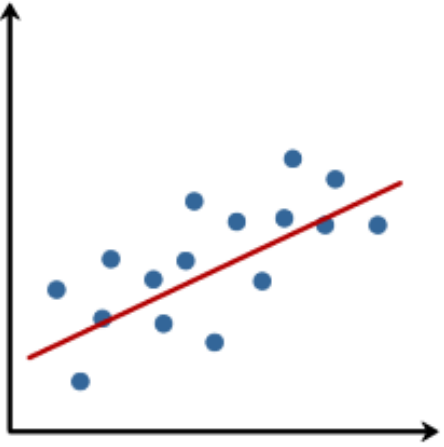
- Focus on making accurate predictions
- Use real-time and historical data
- Predict near-future train delays
- For operational level traffic management



Method



Method

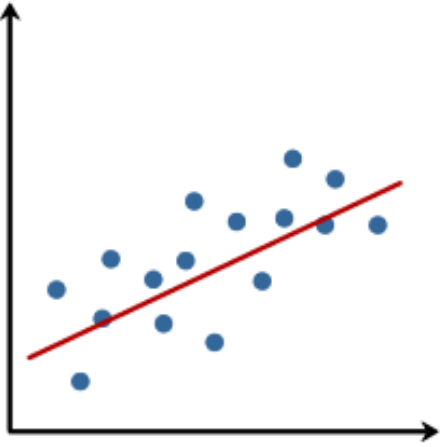


Statistical regression

- It has limitations for modelling complex and non-linear relationships

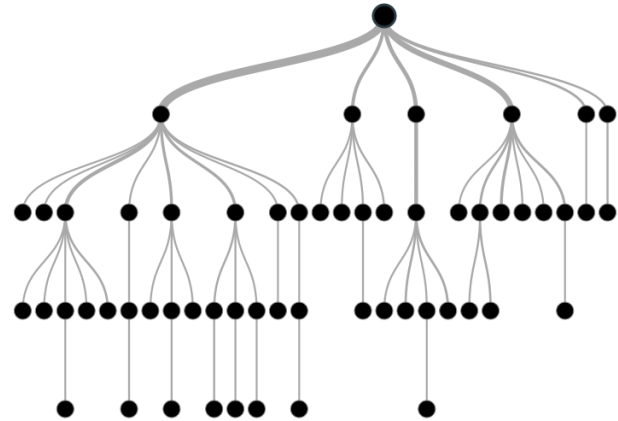


Method



Statistical regression

- It has limitations for modelling complex and non-linear relationships

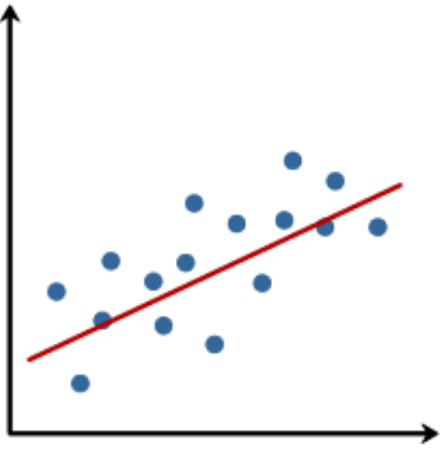


Conventional machine learning

- Less interpretable
- Requires human-engineered spatiotemporal features to capture the spatial and temporal flow patterns of data

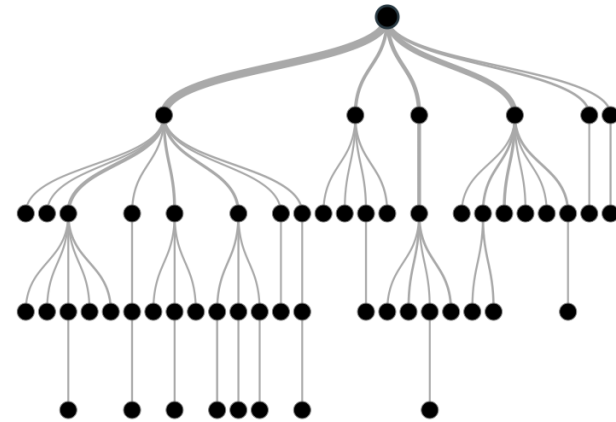


Method



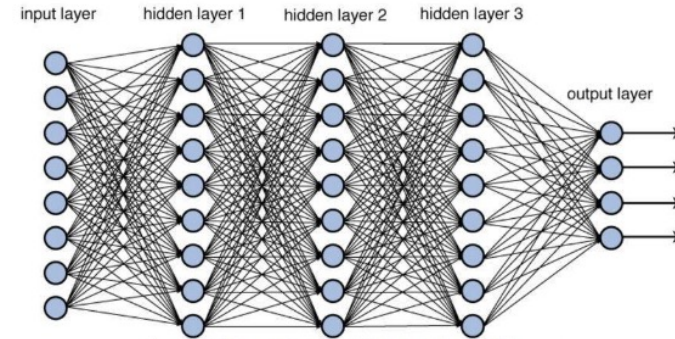
Statistical regression

- It has limitations for modelling complex and non-linear relationships



Conventional machine learning

- Less interpretable
- Requires human-engineered spatiotemporal features to capture the spatial and temporal flow patterns of data

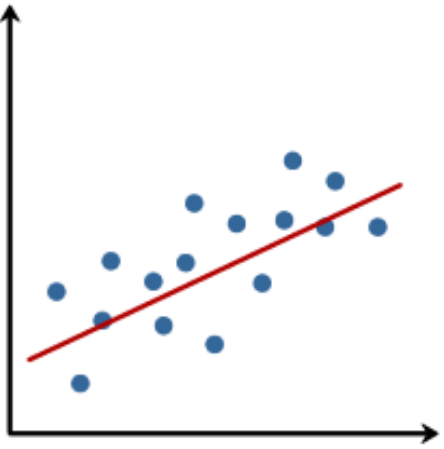


Neural Network

- Automatic learning of spatiotemporal representations from data
- Flexibility to integrate different architectures into hybrid models

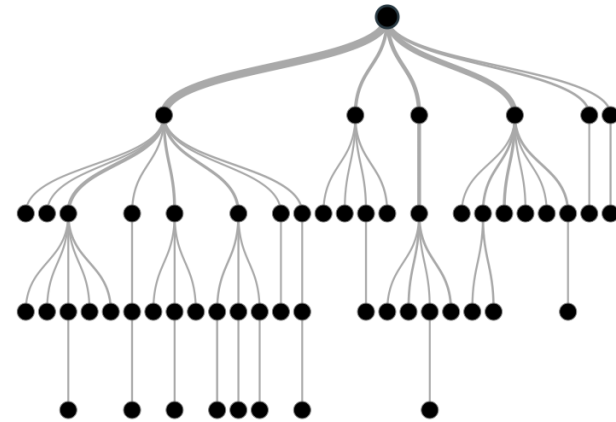


Method



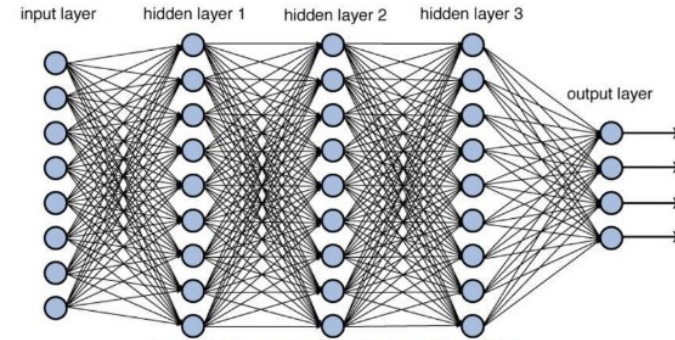
Statistical regression

- It has limitations for modelling complex and non-linear relationships



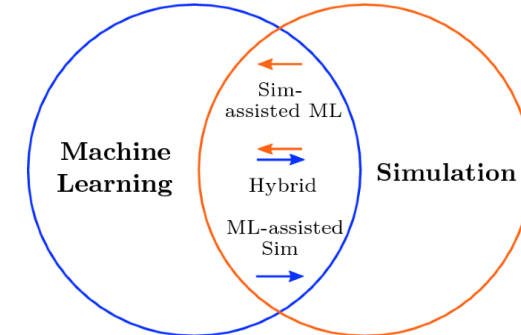
Conventional machine learning

- Less interpretable
- Requires human-engineered spatiotemporal features to capture the spatial and temporal flow patterns of data



Neural Network

- Automatic learning of spatiotemporal representations from data
- Flexibility to integrate different architectures into hybrid models



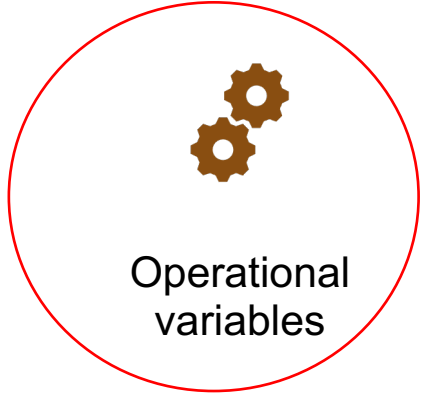
Hybrid model

- Multiple base models with uncorrelated prediction errors

Research question 2

How are selected input variables improving the performance of the train delay prediction model?

Input



Findings

- The **train operation data** greatly influences delays.



Input

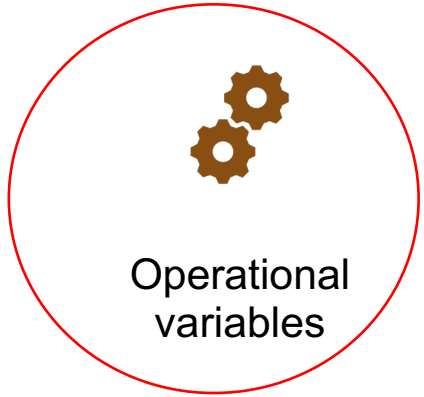


Findings

- The **train operation data** greatly influences delays.
- Other data adds a layer of **adaptability**.



Input



Findings

- The **train operation data** greatly influences delays.
- Other data adds a layer of **adaptability**.
- The **recent observations** from nearby stations or trains are important.



Research question 3

What approaches can enhance
the train delay prediction
model?



Location-conditioned concept

Findings

$$\hat{y} = f(X|i)$$

- Regression models trained **conditionally** on current train location.

where $\hat{y} = (\hat{t}_{i+1}, \hat{t}_{i+2}, \dots, \hat{t}_N)$, denotes the predicted train arrival delays at subsequent stations given current station i . X represents a set of predictor variables encompassing both historical and real-time explanatory factors.

Location-conditioned concept

Findings

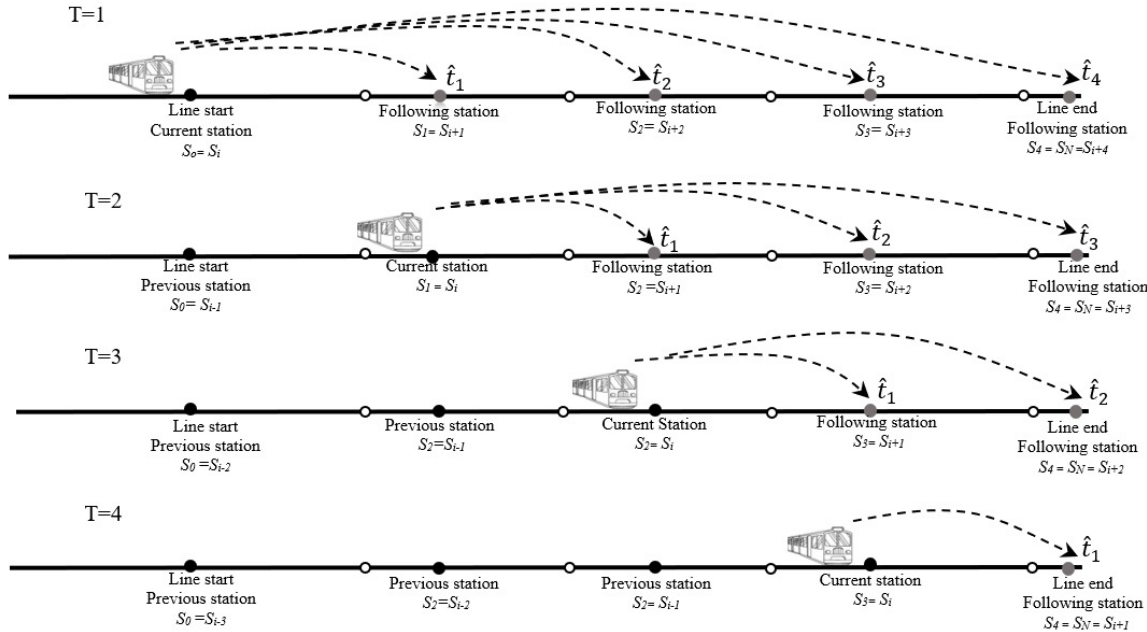
$$\hat{y} = f(X|i)$$

where $\hat{y} = (\hat{t}_{i+1}, \hat{t}_{i+2}, \dots, \hat{t}_N)$, denotes the predicted train arrival delays at subsequent stations given current station i . X represents a set of predictor variables encompassing both historical and real-time explanatory factors.

- Regression models trained **conditionally** on current train location.
- Considers **observable** real-time and historical data.

Multi-output framework

Arrival Times Prediction



- Stop station with observed information
- Stop station with arrival time to be predicted
- Nonstop station

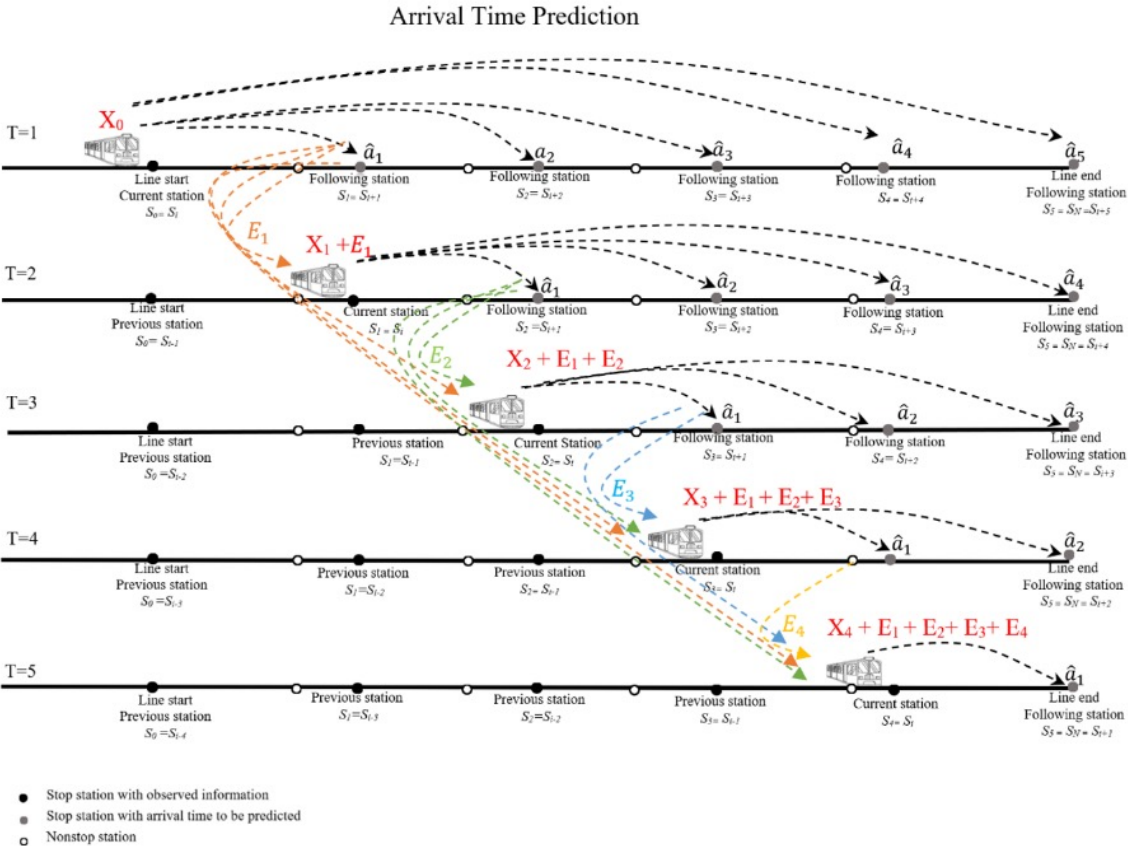
Findings

- Predict arrival delays for multiple downstream stations at arbitrary times.



Error adjustment strategies

Findings

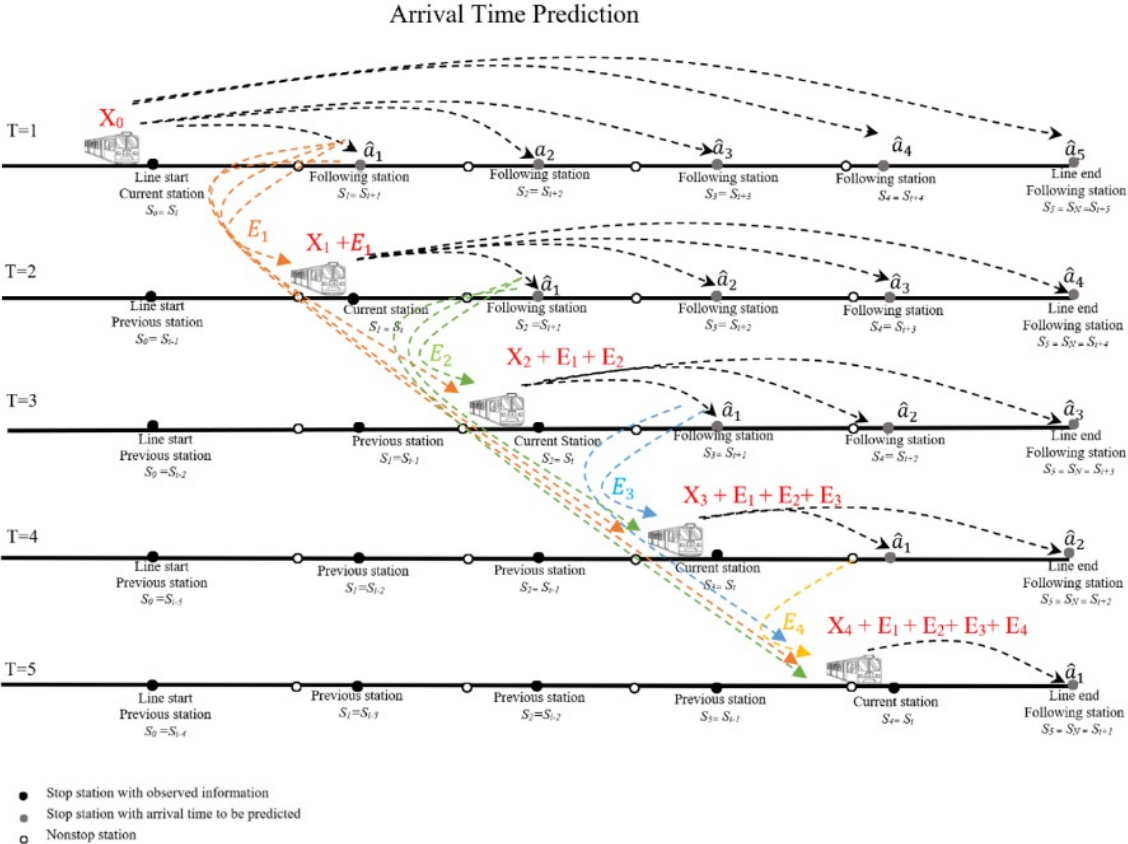


- Use observed information, prediction errors at current and previous stations.

Upstream prediction error correction

Error adjustment strategies

Findings



- Use observed information, prediction errors at current and previous stations.
- Enable the model to constantly adjust itself.

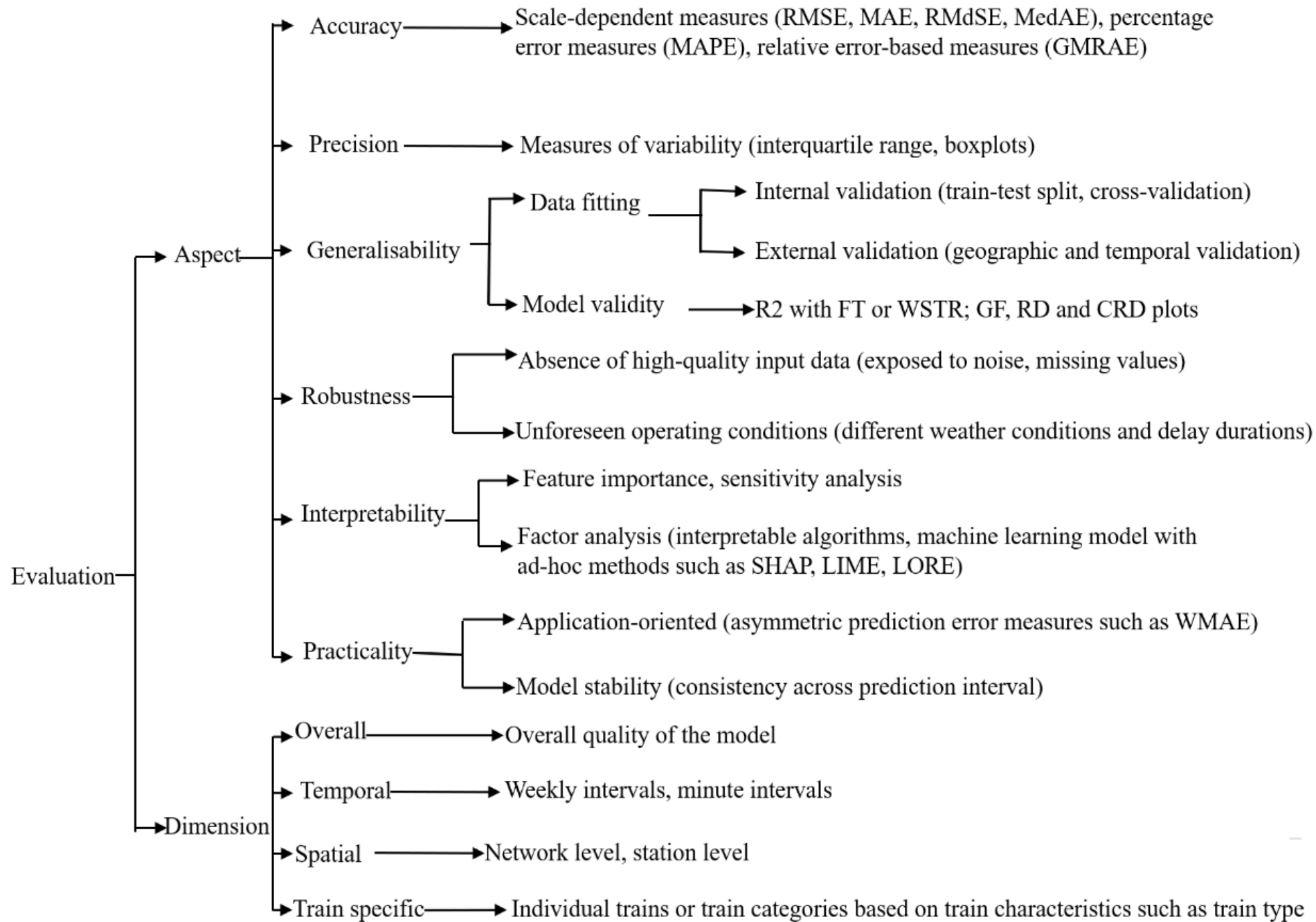
Upstream prediction error correction



Research question 4

How can train delay prediction models be evaluated?

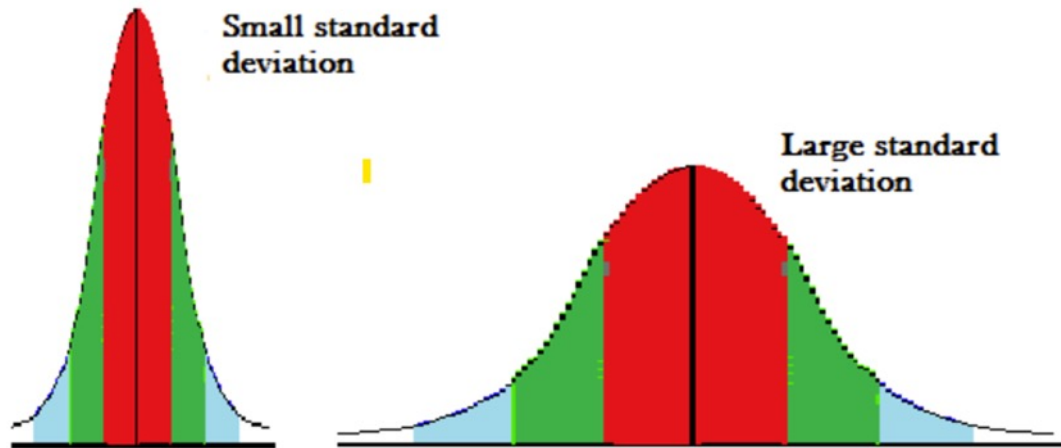




Precision

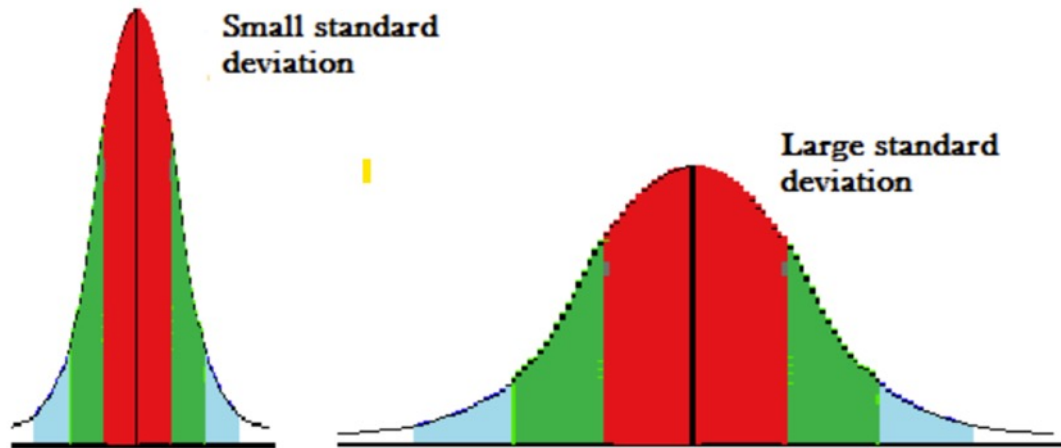
Findings

- Measures **dispersion of prediction error and bias tendencies.**



Statistical variance or the spread of data

Precision

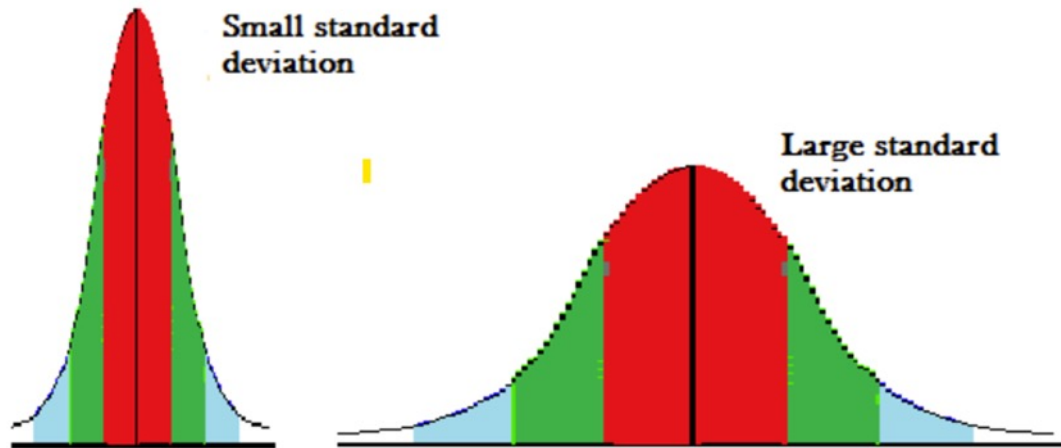


Statistical variance or the spread of data

Findings

- Measures **dispersion of prediction error and bias tendencies.**
- The **narrower ranges** mean more reliable predictions.

Precision



Statistical variance or the spread of data

Findings

- Measures **dispersion of prediction error and bias tendencies**.
- The **narrower ranges** mean more reliable predictions.
- The **interquartile range or boxplot** clarify prediction error uncertainty.

Robustness



(a) Invalid inputs



Robustness



(a) Invalid inputs



(b) Challenging environmental conditions

Robustness



(a) Invalid inputs



(b) Challenging environmental conditions

Findings

- Use datasets with **realistic representative** of real-world application scenarios.

Robustness



(a) Invalid inputs



(b) Challenging environmental conditions

Findings

- Use datasets with **realistic representative** of real-world application scenarios.
- Prevent **purely academic contributions** without real-world industrial use.

Practicality



(a) Application-oriented



Practicality



(a) Application-oriented

Findings

- Tolerance for prediction errors varies depending on the **model's use case**.



Practicality



(a) Application-oriented

Findings

- Tolerance for prediction errors varies depending on the **model's use case**.
- Measure using **asymmetric prediction error measures**



Practicality



(a) Application-oriented

(b) Stability of predictions

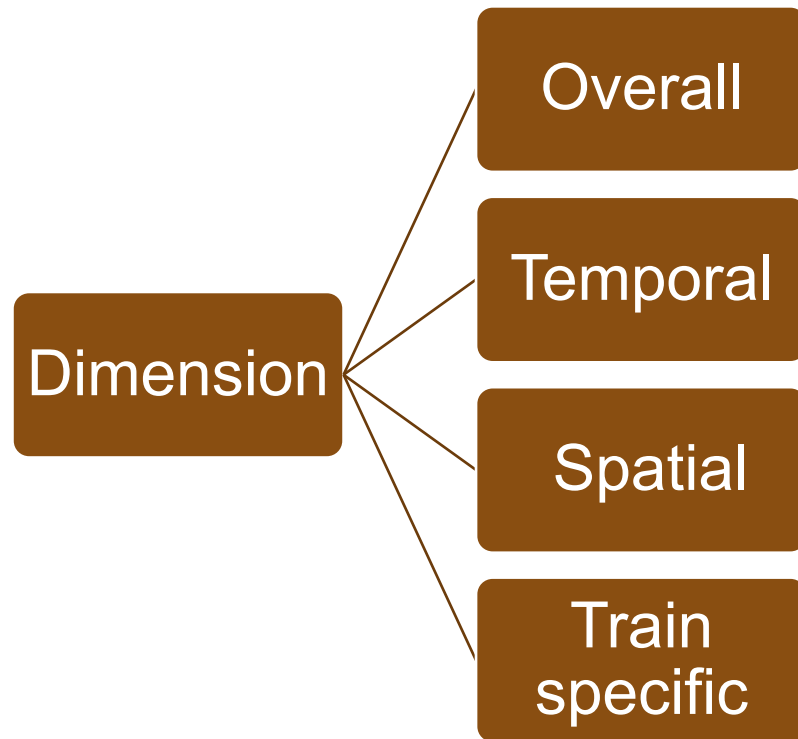
Departures		03:38:07	
Route	Destination	Stop	Time
1	Blackyard leys	R2	~ 5 Min
4A	Wood farm	R7	~ 20 Min
4	Abingdon	R8	~ 40 min
5	Old woodstock	R5	~ 05:10 Pm
7	Aylesbury	R2	~ 05:45 PM
S1	Thornhill park & ride	R6	~06:00 PM
S3	Seacourt park & ride	R1	~06:20 PM

Findings

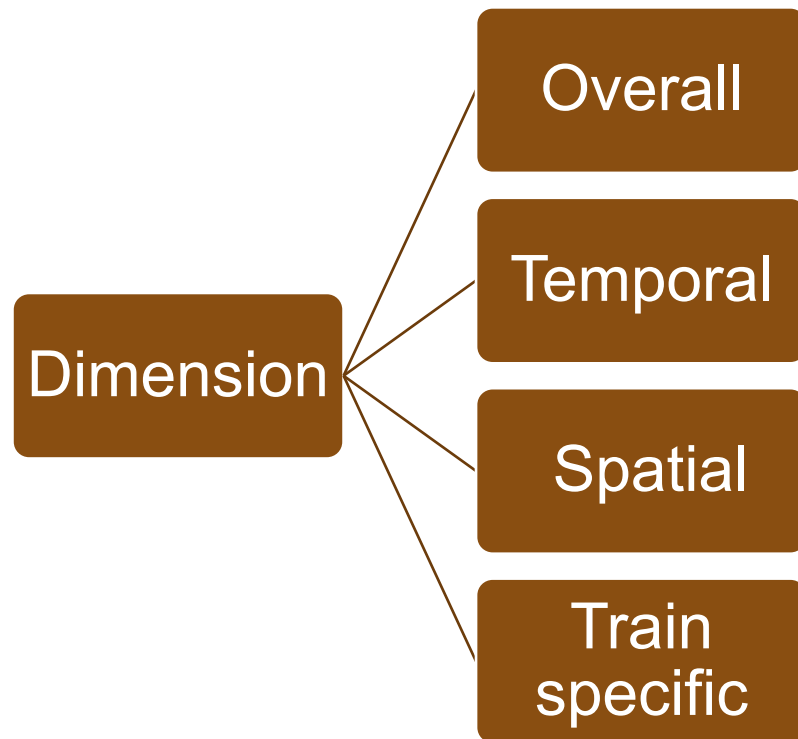
- Tolerance for prediction errors varies depending on the **model's use case**.
- Measure using **asymmetric prediction error measures**
- Assesses the **consistency of the predictions** at each prediction interval .



Dimension



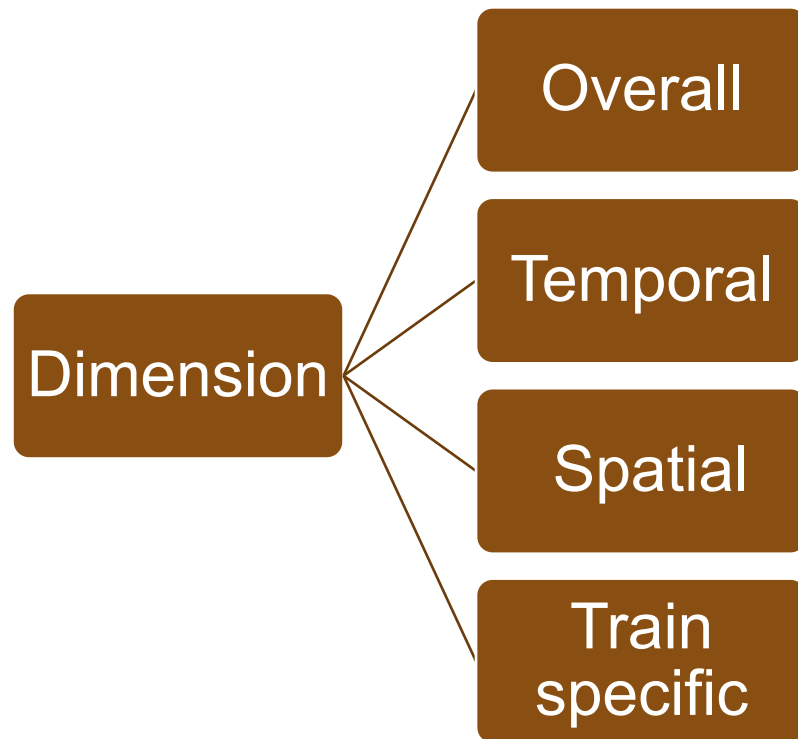
Dimension



Findings

- Overall performance evaluation provides **overview** of the model's quality.

Dimension



Findings

- Overall performance evaluation provides **overview** of the model's quality.
- Detailed evaluations across dimensions uncover **underlying** performance patterns.



Conclusion & Future research

Conclusion

- Use recent data improves train delay prediction model.
 - Introduced location conditional concepts and error adjustment strategies
 - Generate synthetic train events



Conclusion

- Use recent data improves train delay prediction model.
 - Introduced location conditional concepts and error adjustment strategies
 - Generate synthetic train events
- Dynamic multi-output prediction models are crucial for practical applications
 - Introduced line-level multi-output machine learning models
 - Network-level prediction models



Conclusion

- Use recent data improves train delay prediction model.
 - Introduced location conditional concepts and error adjustment strategies
 - Generate synthetic train events
- Dynamic multi-output prediction models are crucial for practical applications
 - Introduced line-level multi-output machine learning models
 - Network-level prediction models
- Evaluate models from various aspects and dimensions
 - Established an evaluation framework
 - Conduct comprehensive case studies

Thank you!!!



LUND
UNIVERSITY



LUND
UNIVERSITY