



**GHENT
UNIVERSITY**

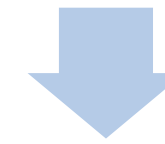
DATABASES

SPARK PROJECT FOR NMBS / 14.12.2017

BELGIAN PUBLIC TRANSPORTATION



Explore



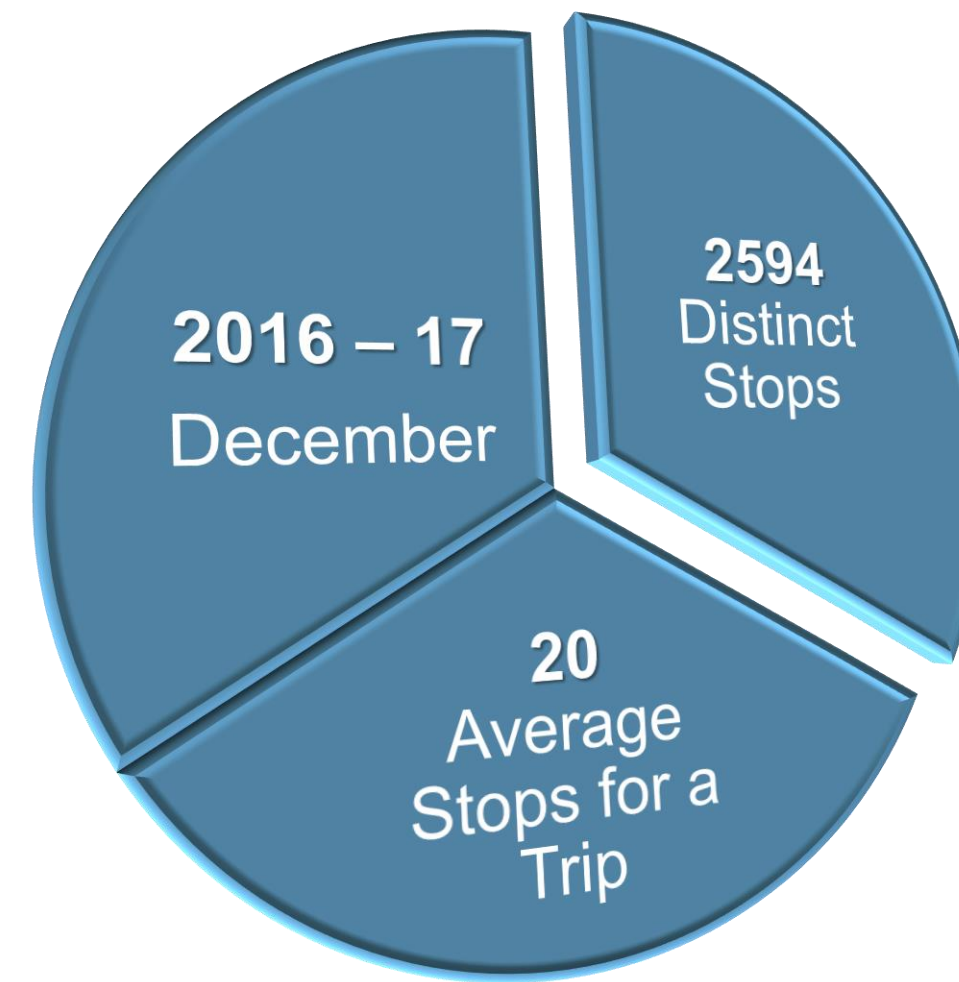
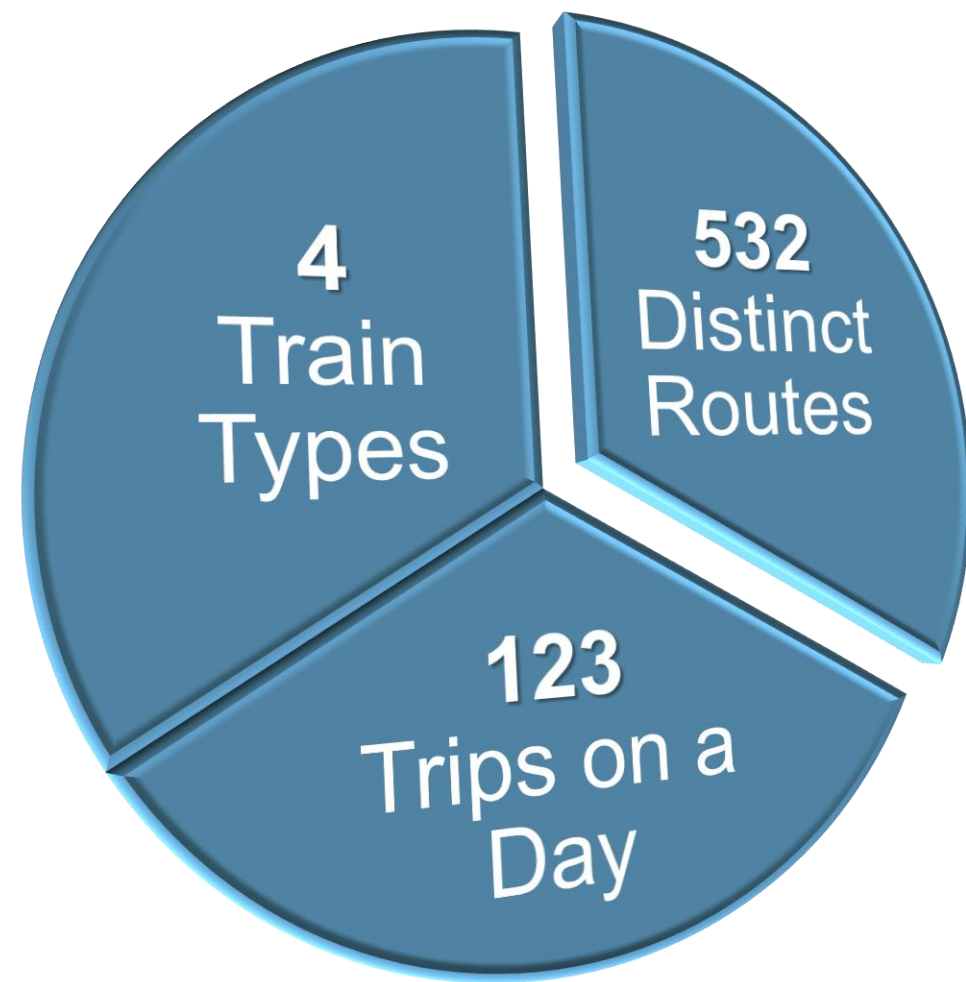
**Identify
Problems**



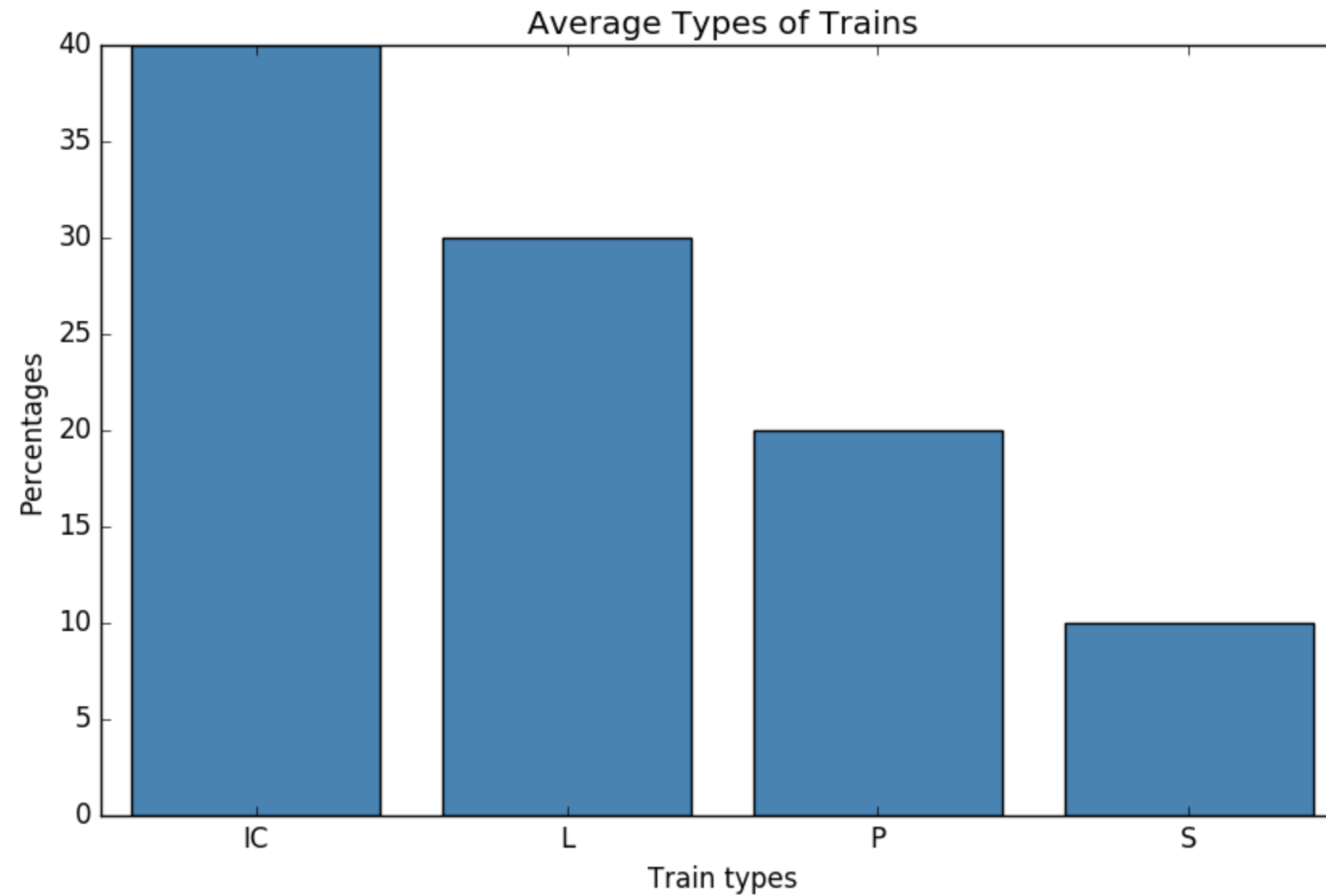
**Suggest
Solutions**

DESCRIPTIVE STATISTICS

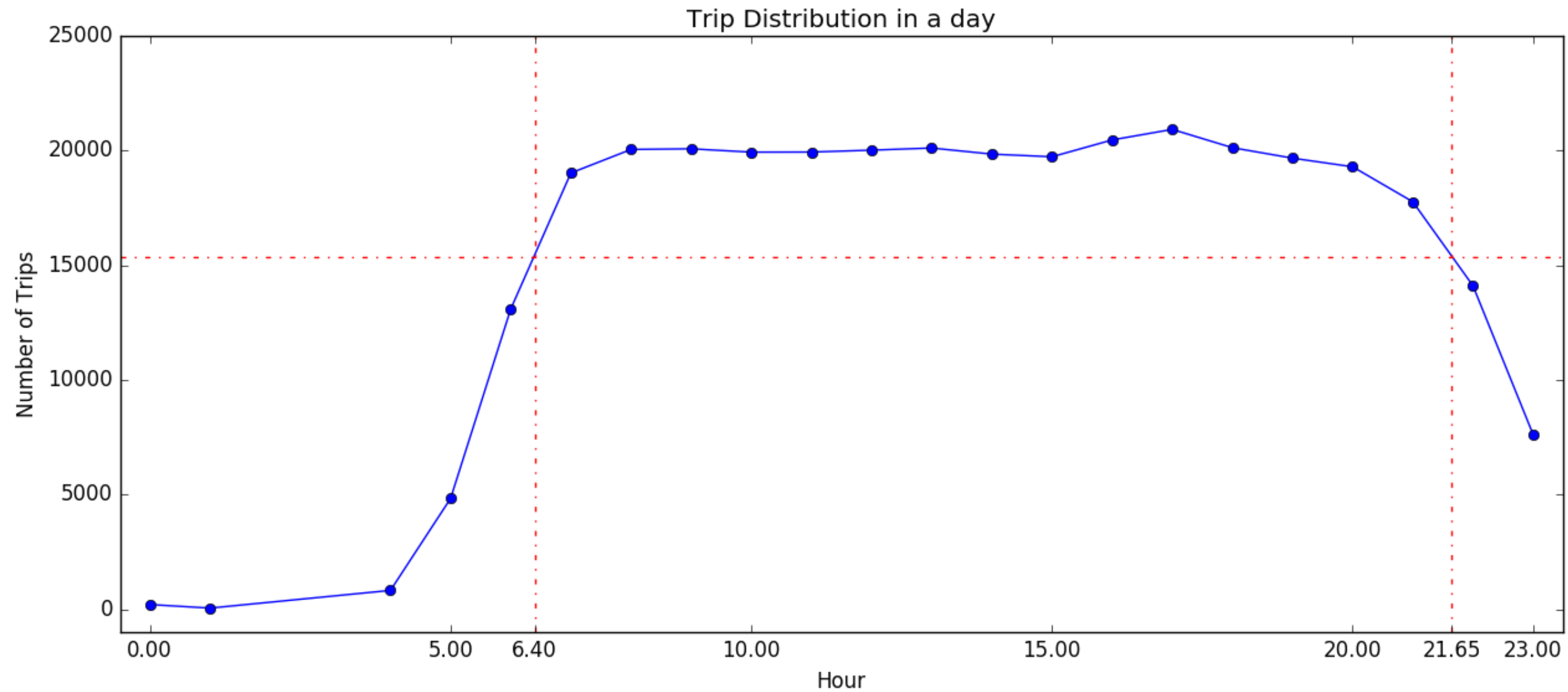
DESCRIPTIVE STATISTICS



DESCRIPTIVE STATISTICS



DESCRIPTIVE STATISTICS

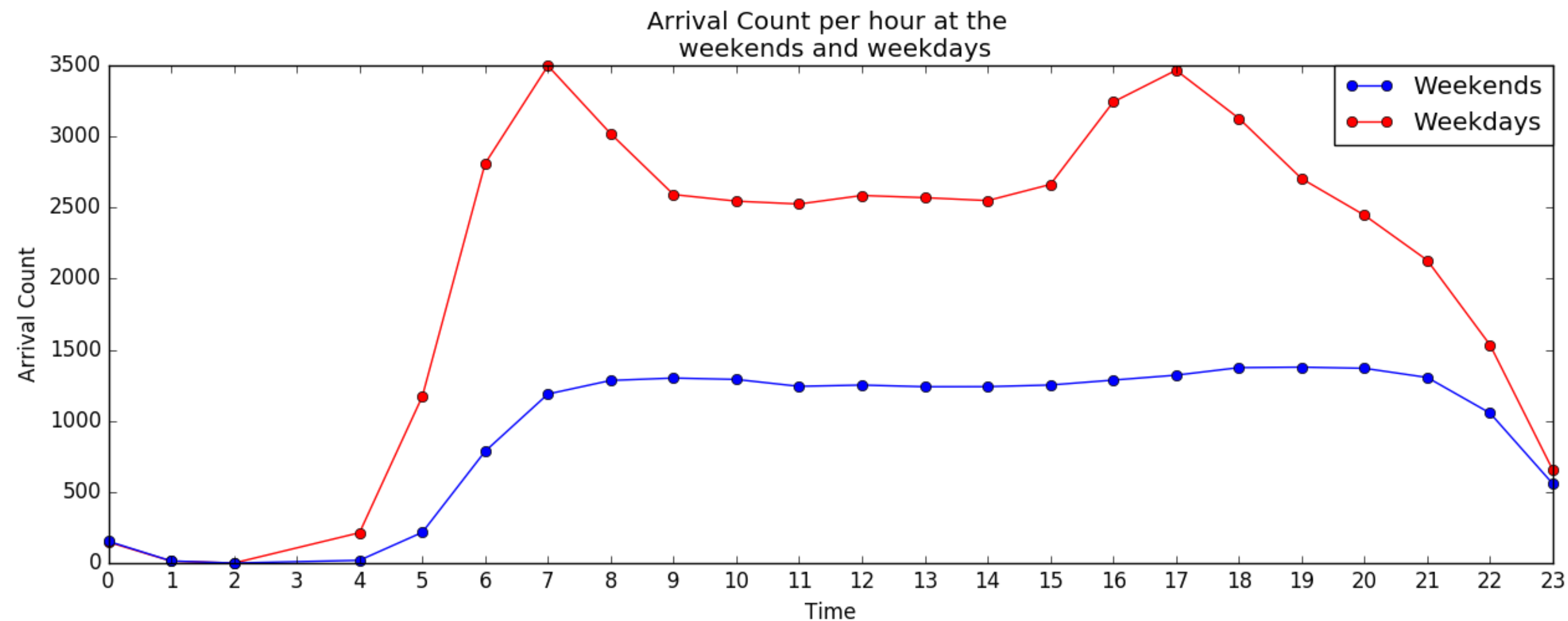


INCOMING – OUTGOING TRAINS

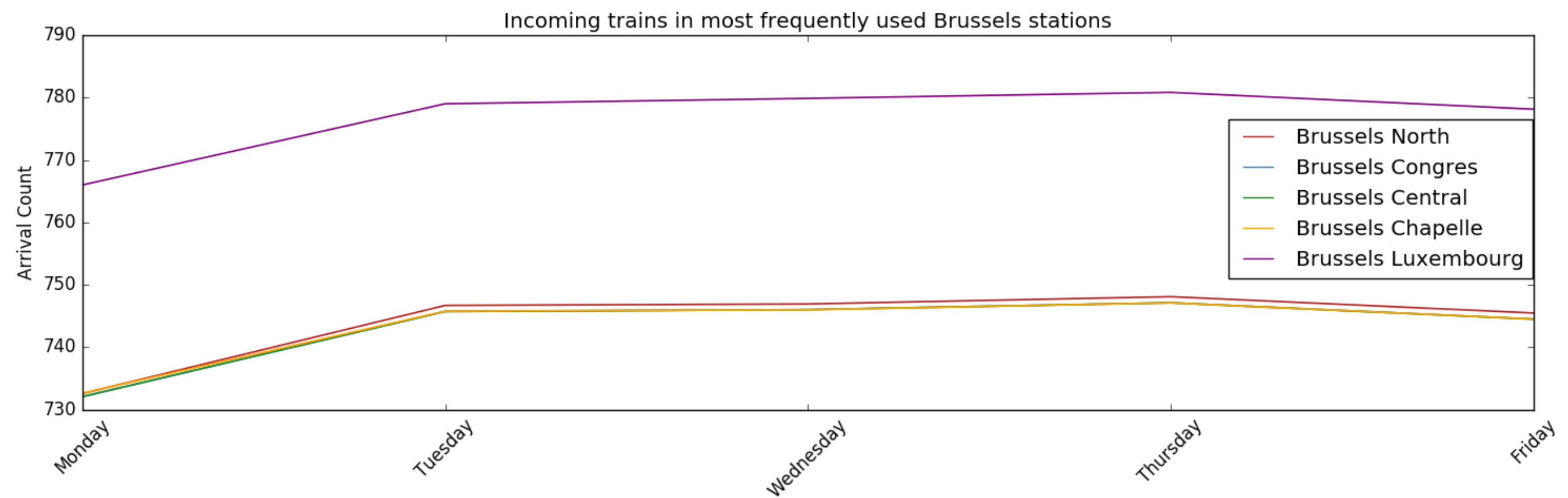
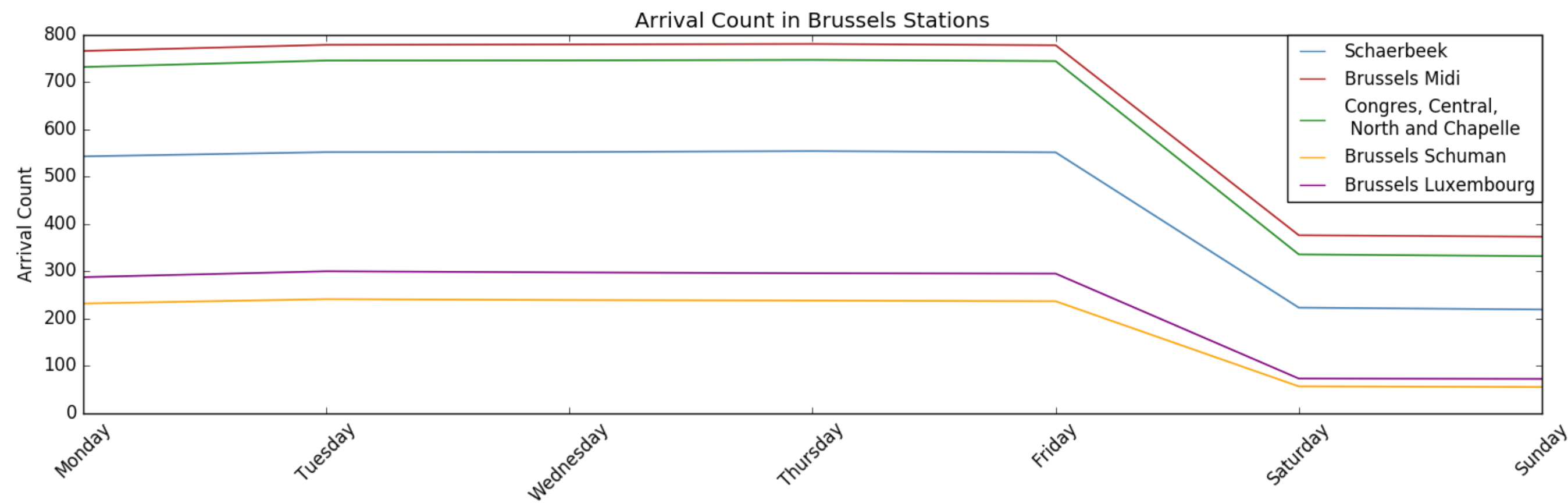
INCOMING – OUTGOING TRAINS



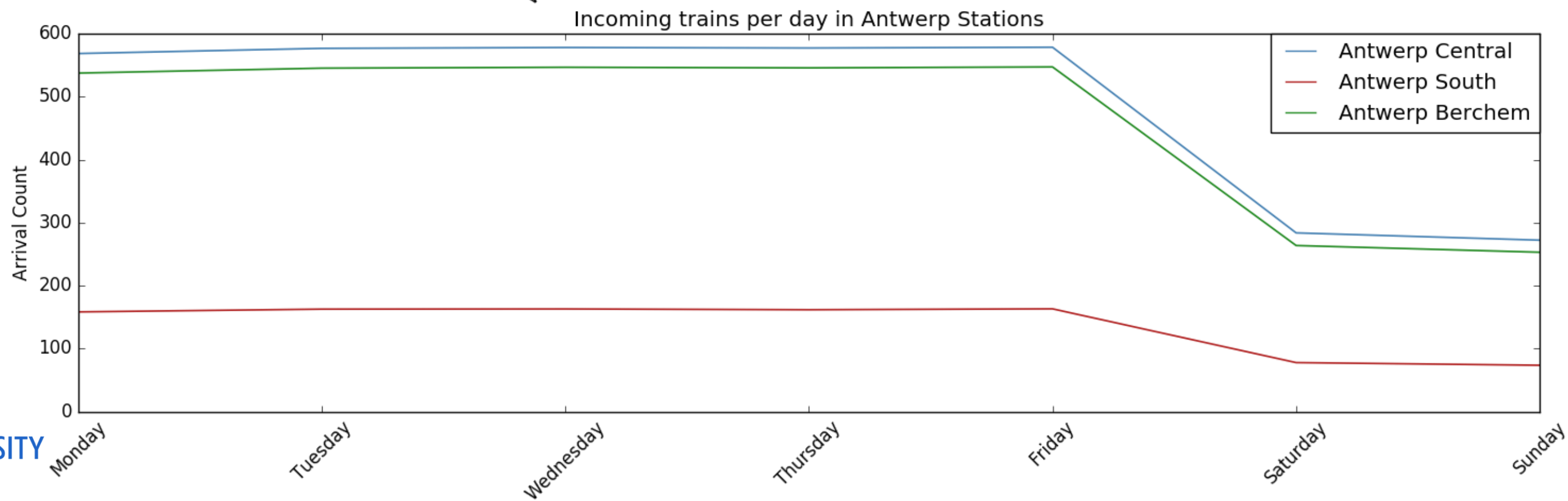
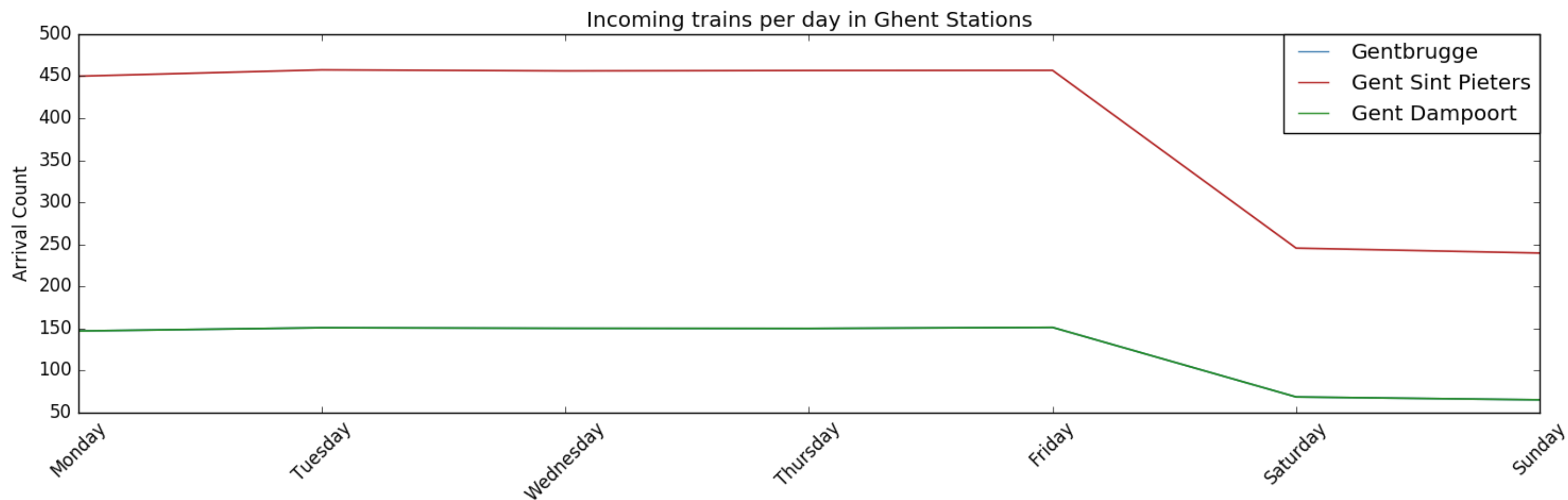
INCOMING – OUTGOING TRAINS



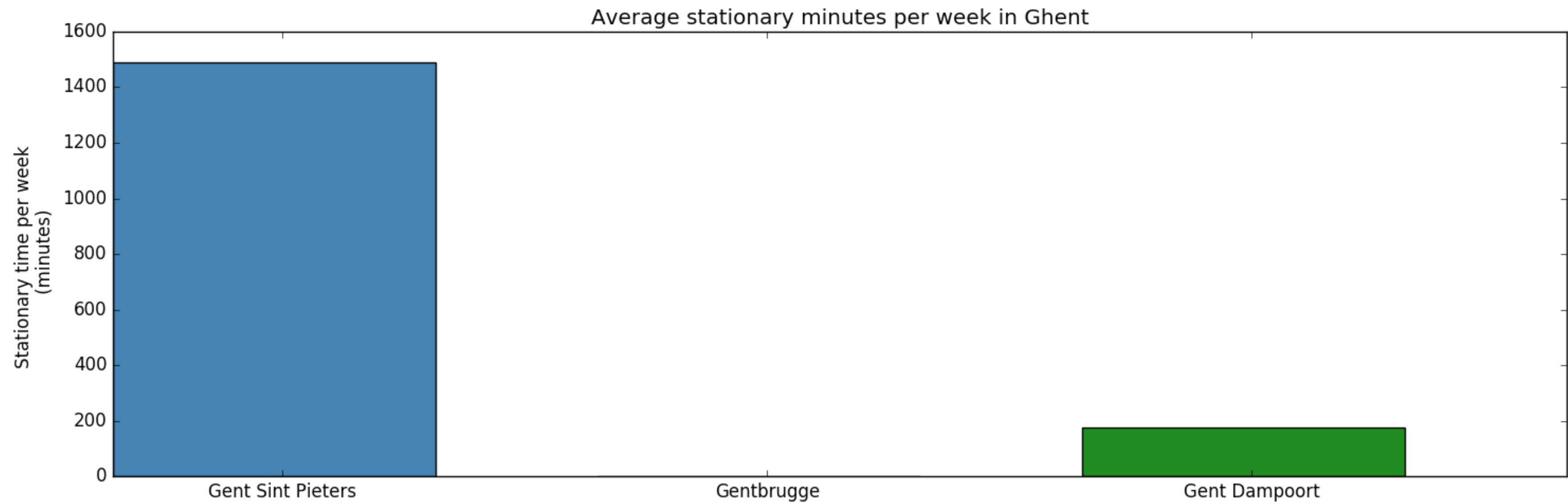
INCOMING – OUTGOING TRAINS



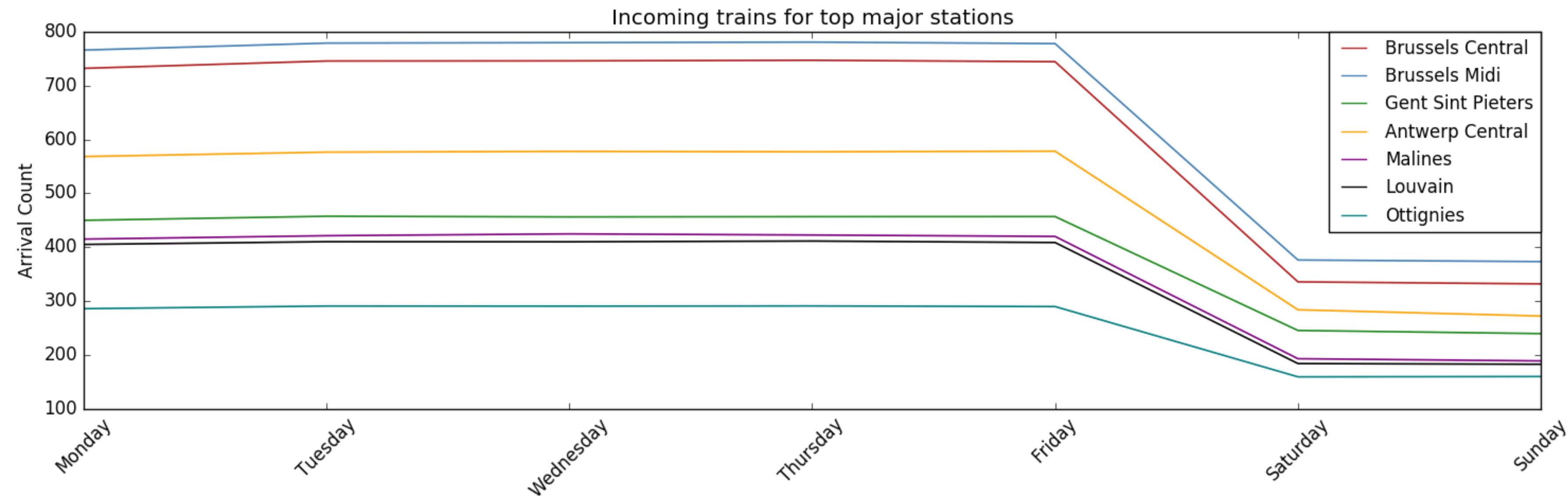
INCOMING - OUTGOING TRAINS



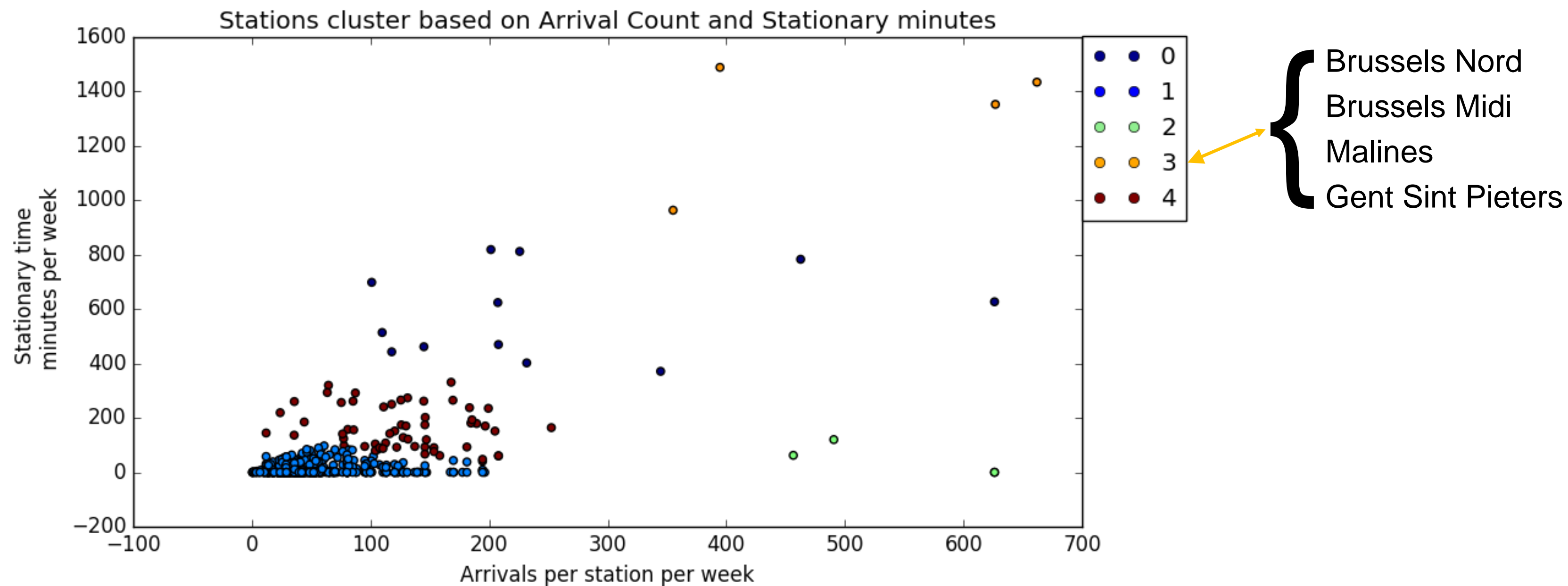
INCOMING – OUTGOING TRAINS



INCOMING - OUTGOING TRAINS



INCOMING – OUTGOING TRAINS



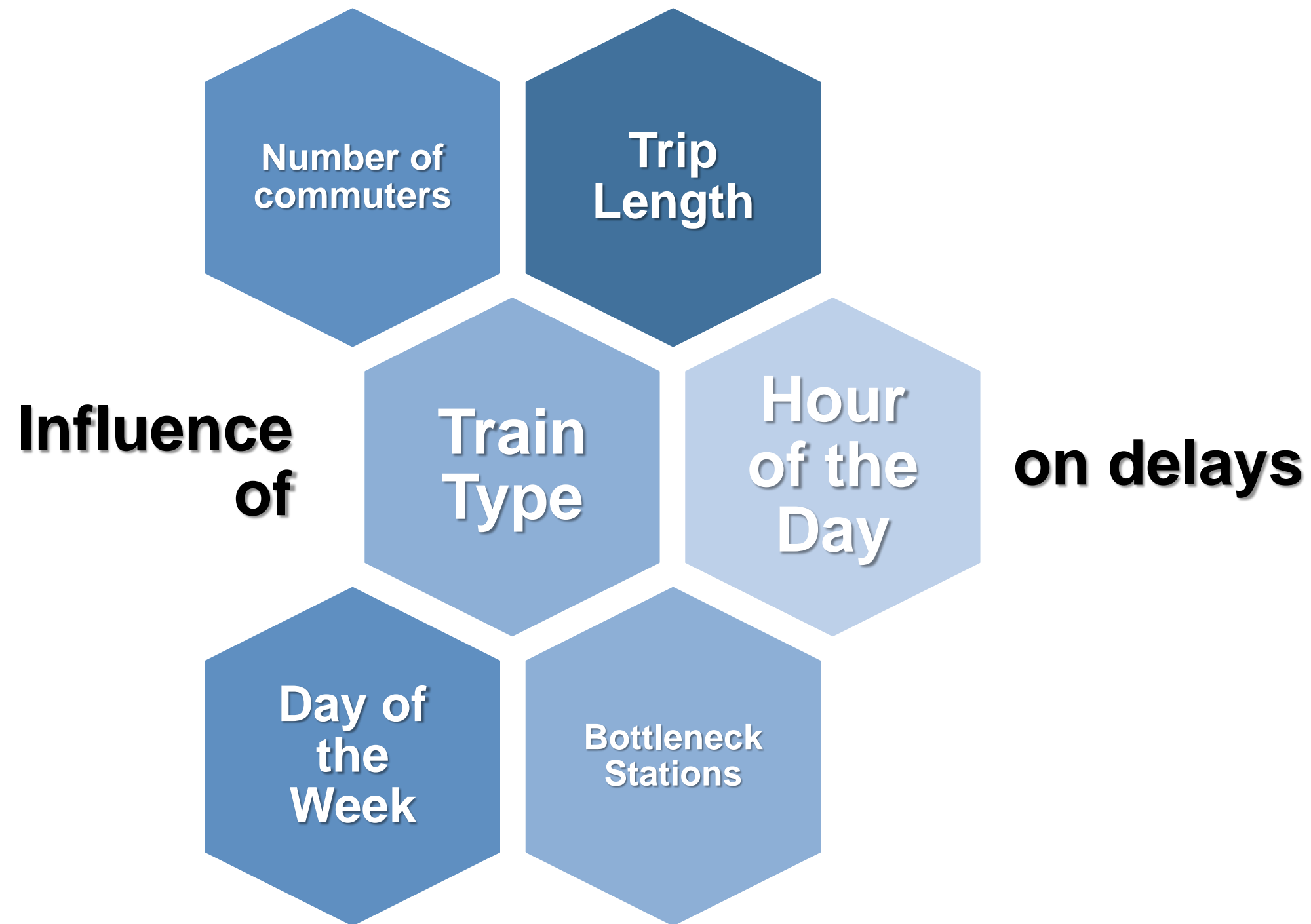
ANALYSIS OF DELAYS

ANALYSIS OF DELAYS

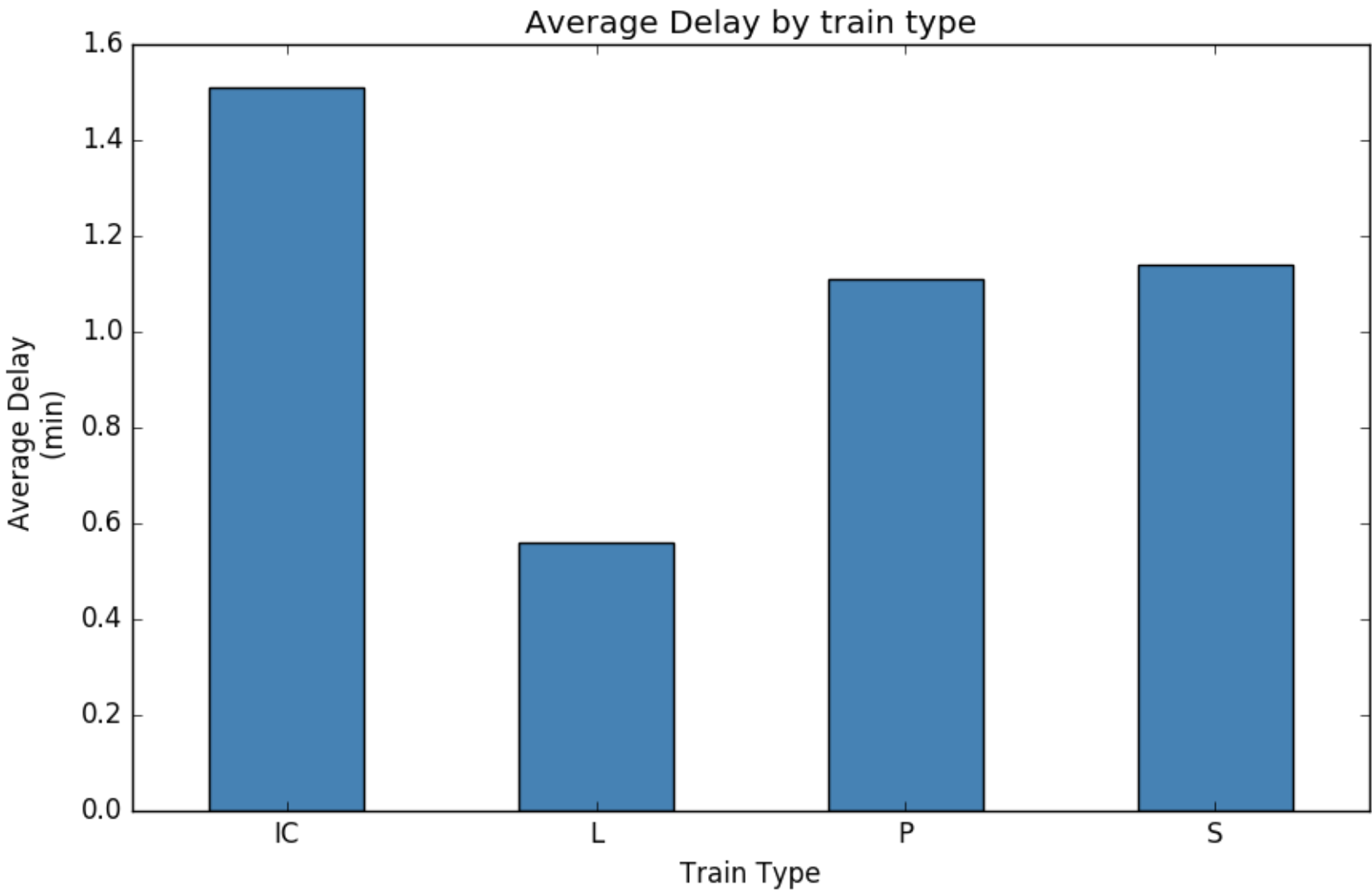
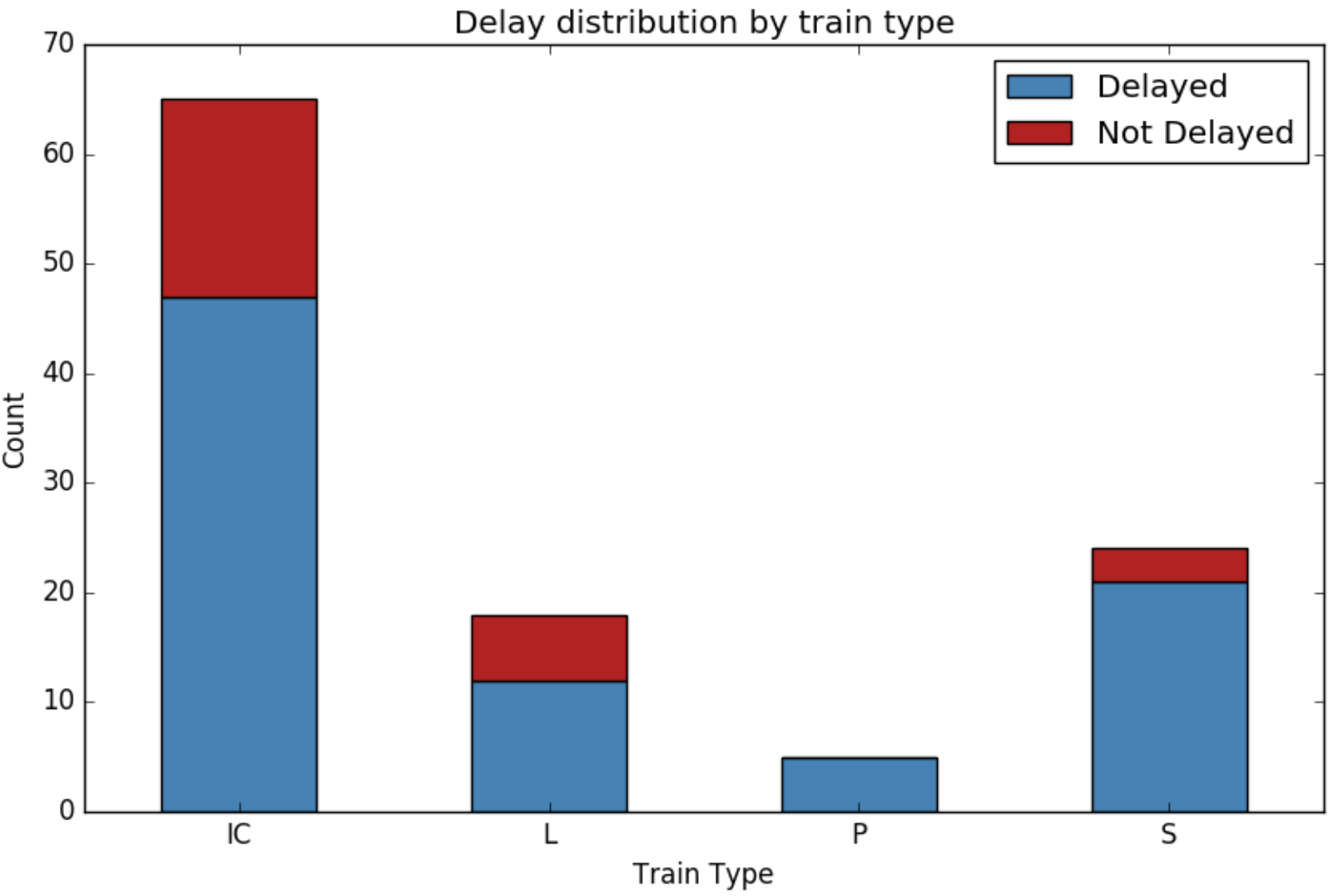


DEPART	TRAIN	RailTime VOIE	RETARD
K-SCHAARBEEK	P	4	+0H35
PE-C. BINCHE	IR	16	+1H00
ORTRIJK	IC	12	+0H27
	P	20	+0H21
			+0H15

ANALYSIS OF DELAYS

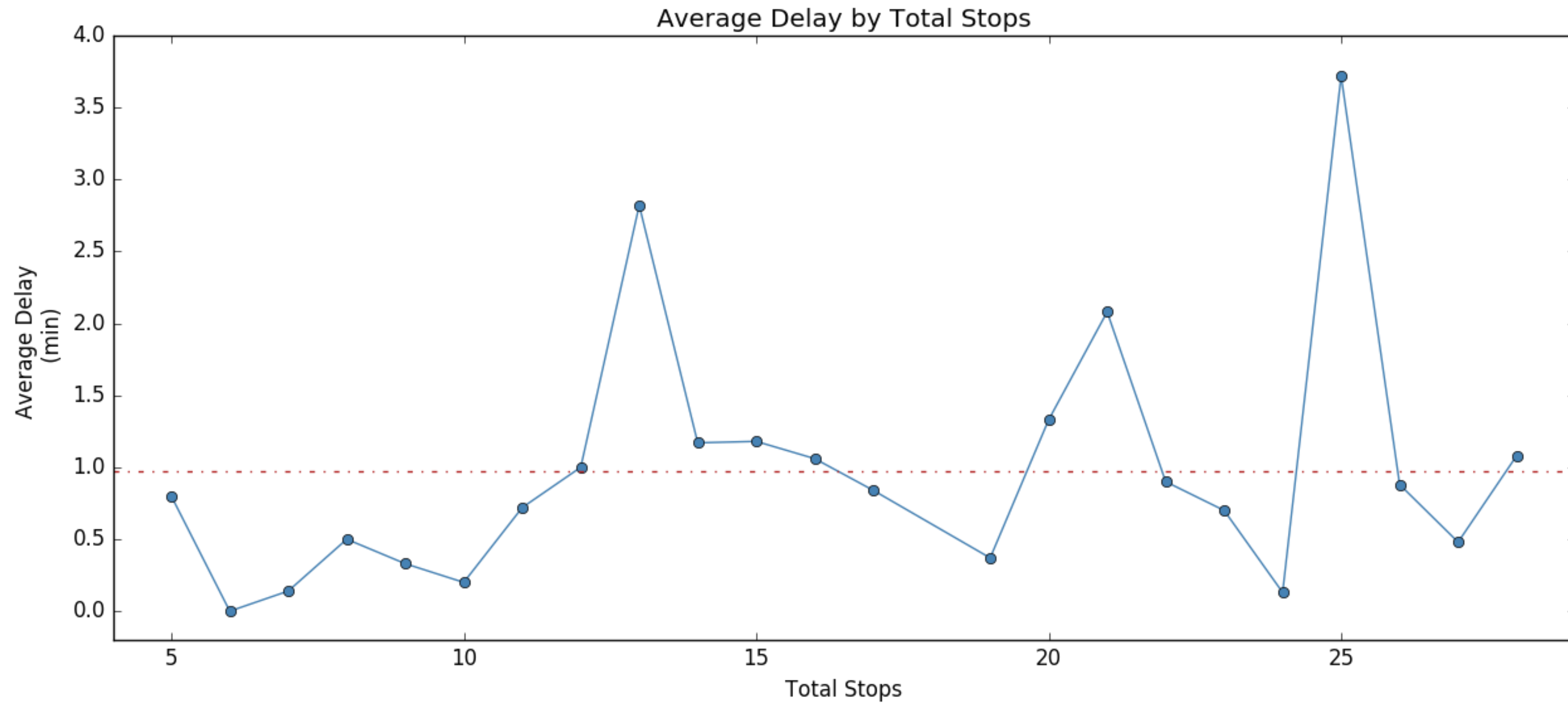


ANALYSIS OF DELAYS

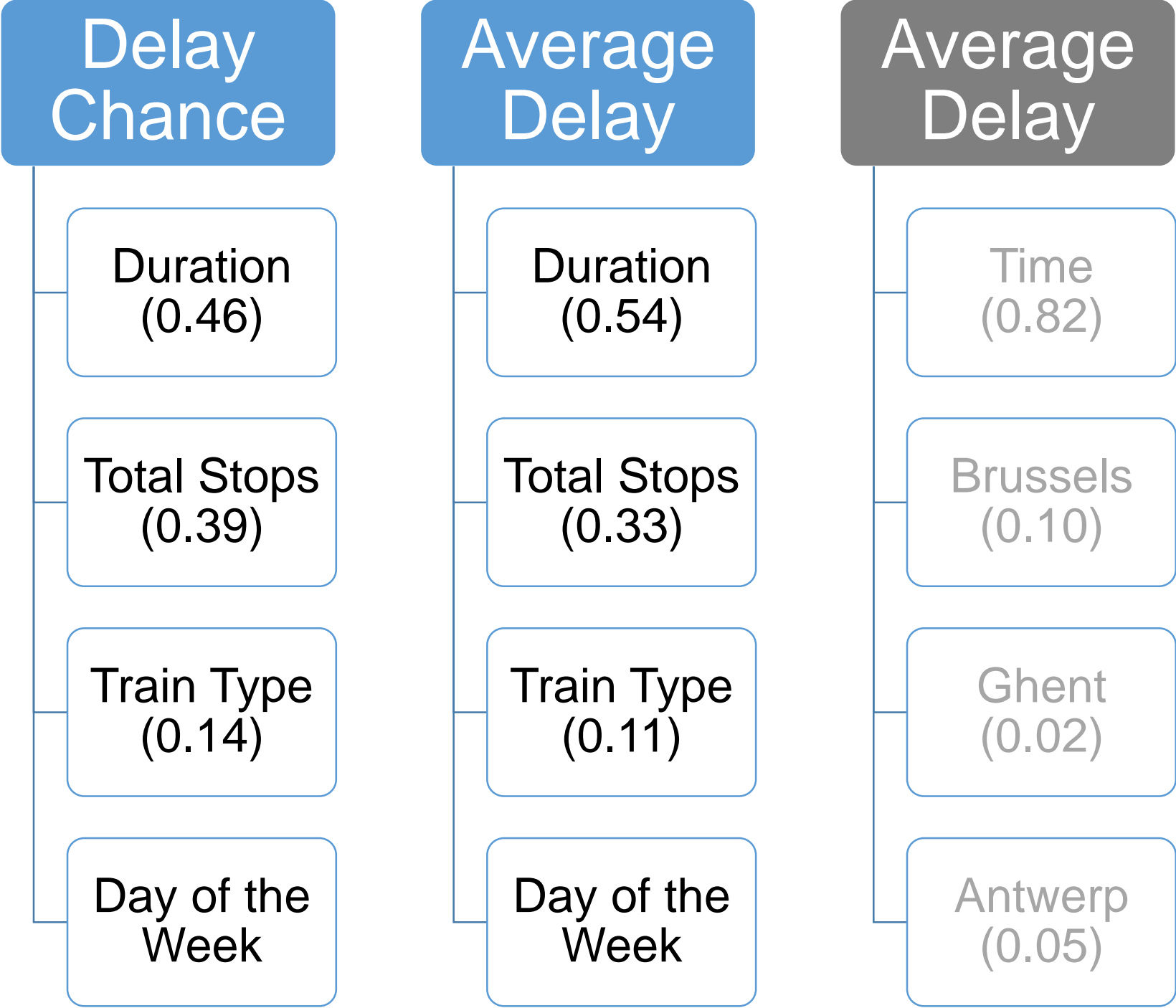


* Not delayed: average delay time ≤ 0.25 min

ANALYSIS OF DELAYS



RANDOM FOREST



Random Forest
Classification
AUC = 0.55
Accuracy = 69%

Random Forest
Regression

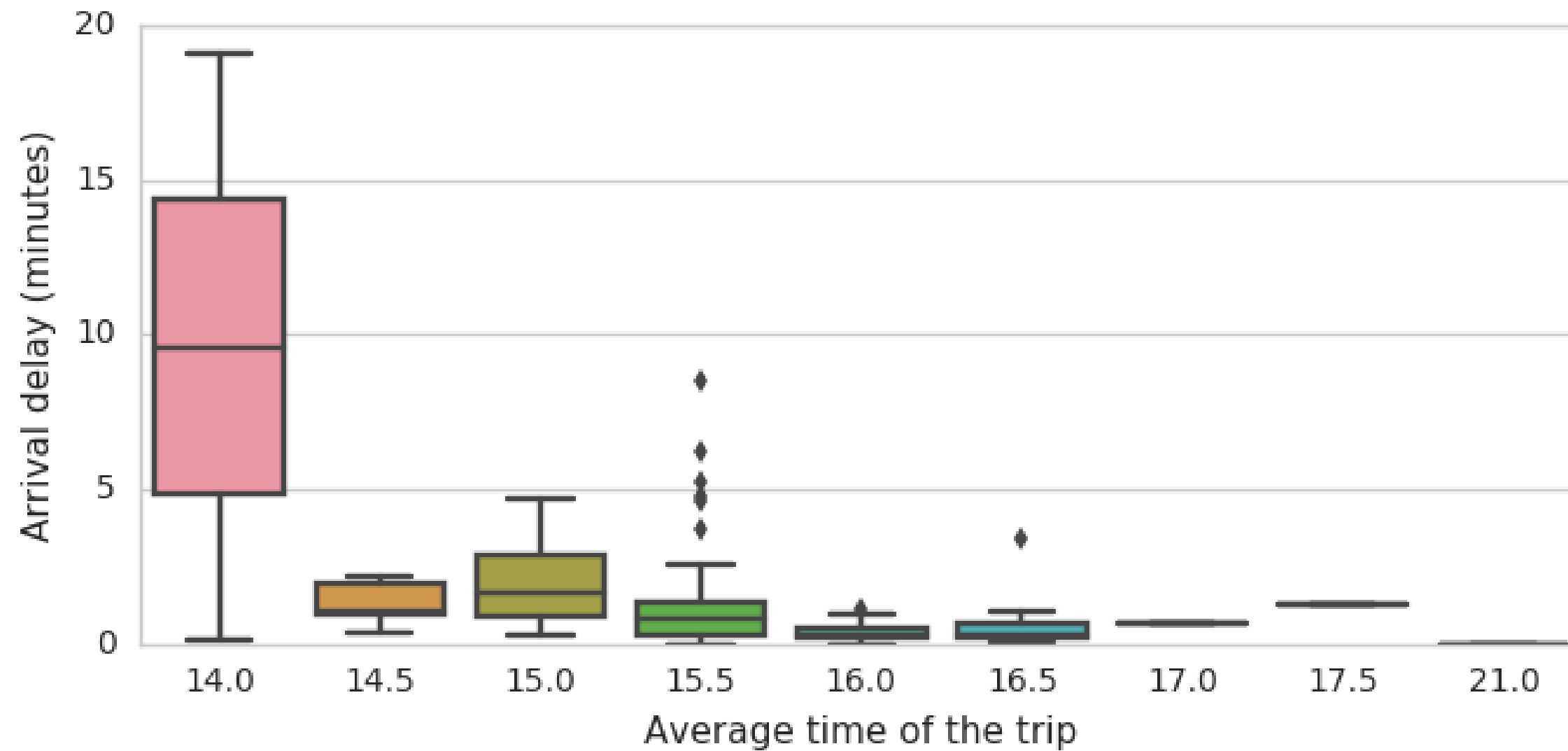
Random Forest
Regression

INFLUENCE OF TIME OF DAY



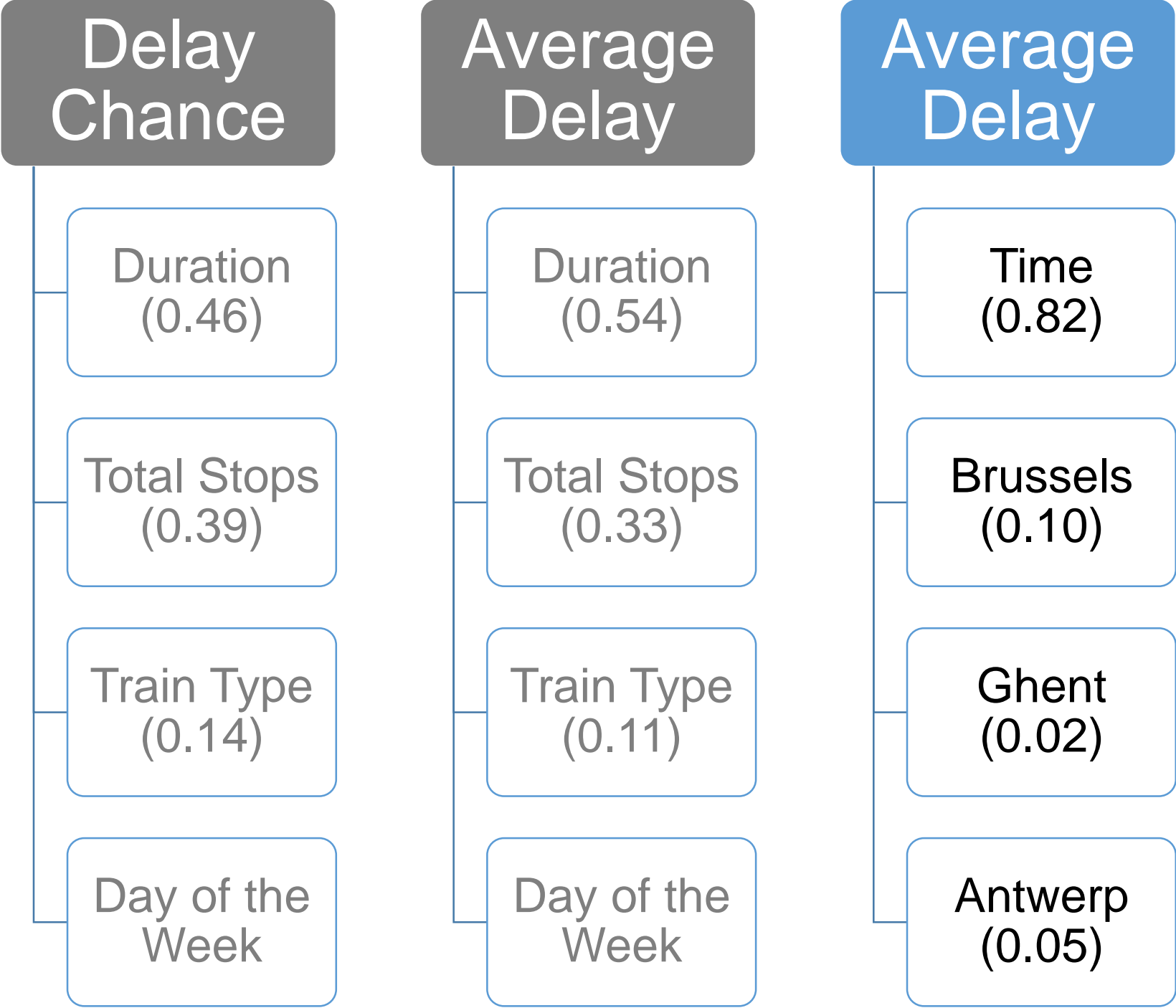
- Is the delay time longer at certain moments?
- Include possible bottleneck cities (Brussels, Ghent, ...)

INFLUENCE OF TIME OF DAY



- Divide 24h day in categories of half an hour each
- Dummies for each train: whether it passed in Brussels, Ghent or Antwerp
- Perform random forest regression

RANDOM FOREST



Random Forest
Classification
AUC = 0.55
Accuracy = 69%

Random Forest
Regression

Random Forest
Regression

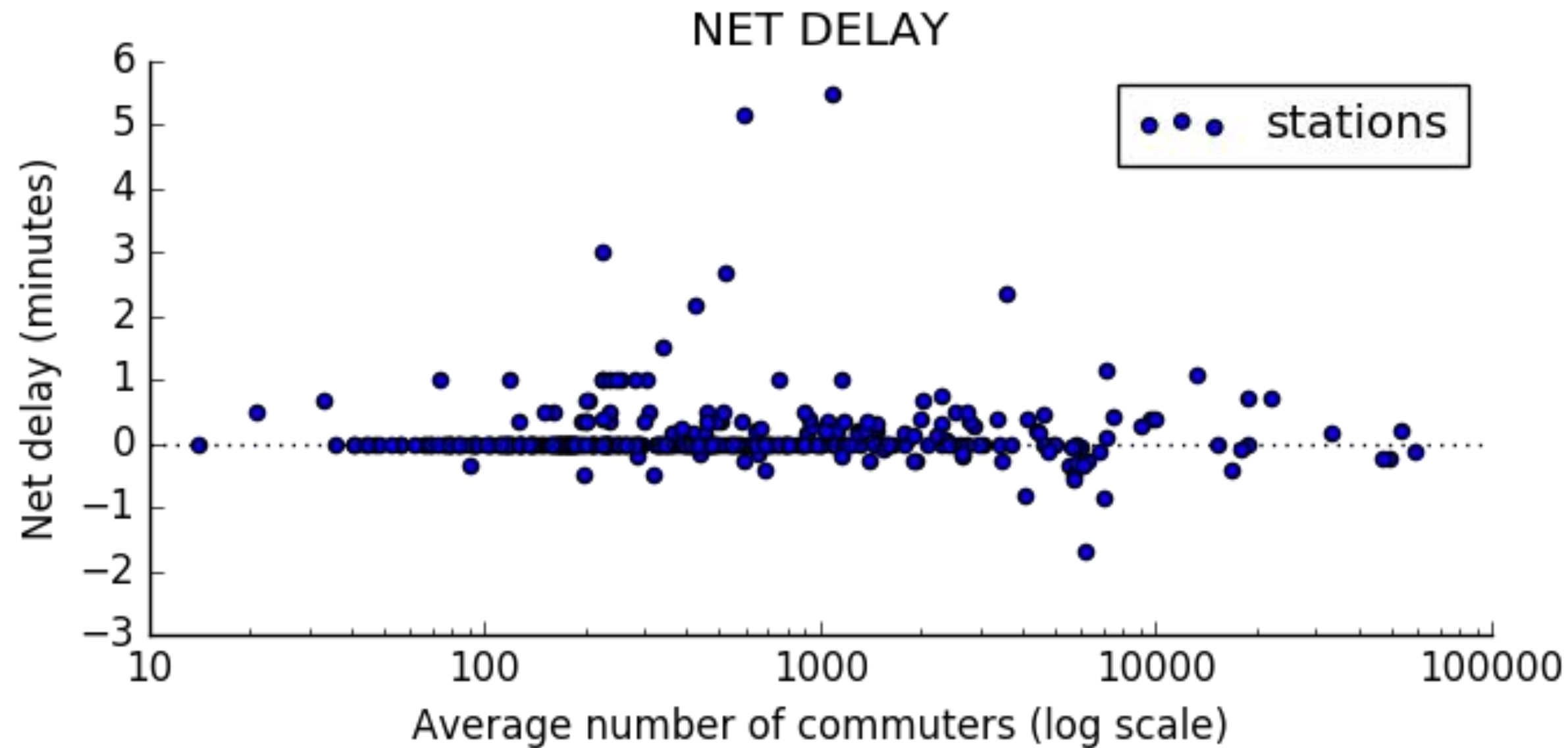
INFLUENCE OF COMMUTERS



- Are trains more easily delayed in busier stations?

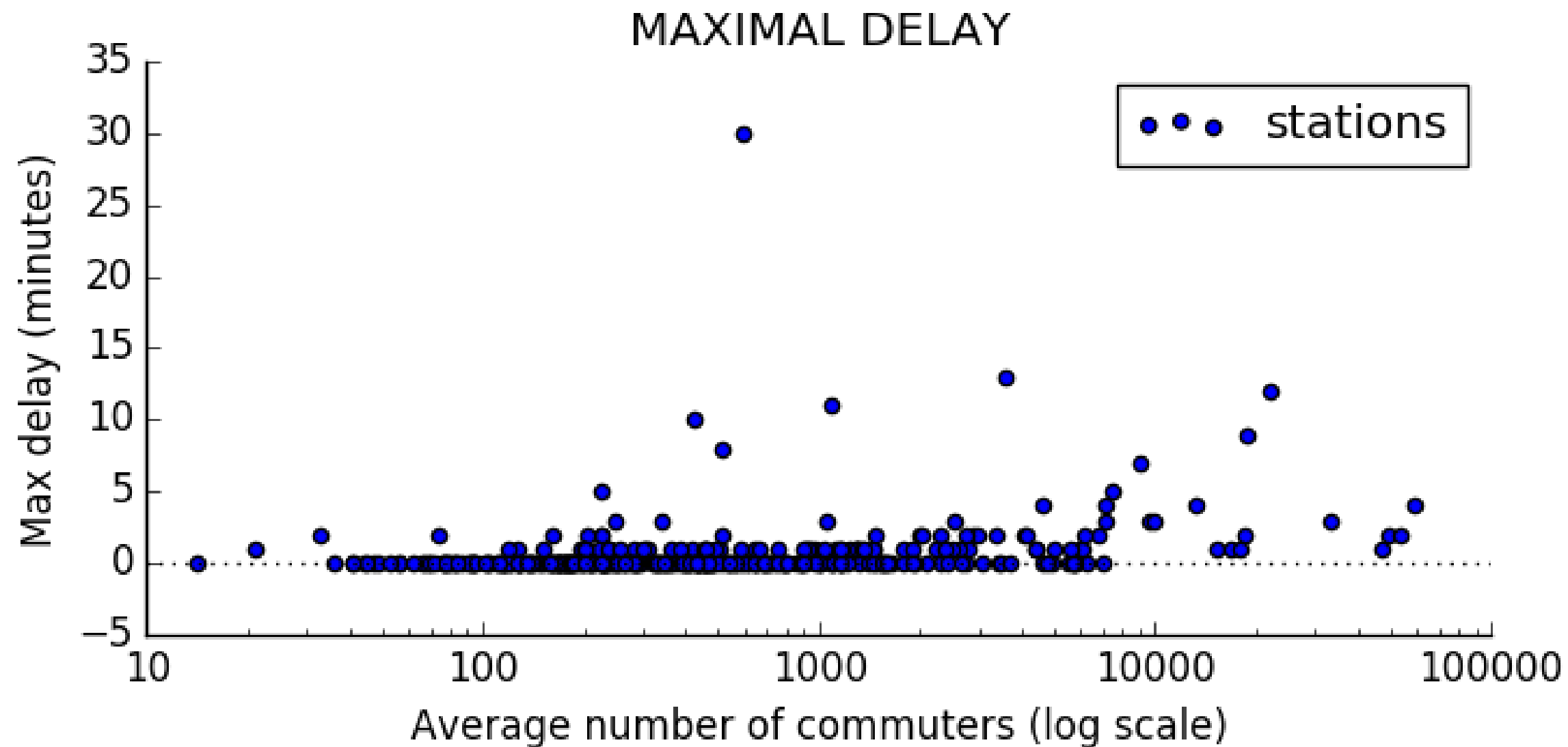
**Extra dataset: average number of commuters per station on weekday*

INFLUENCE OF COMMUTERS



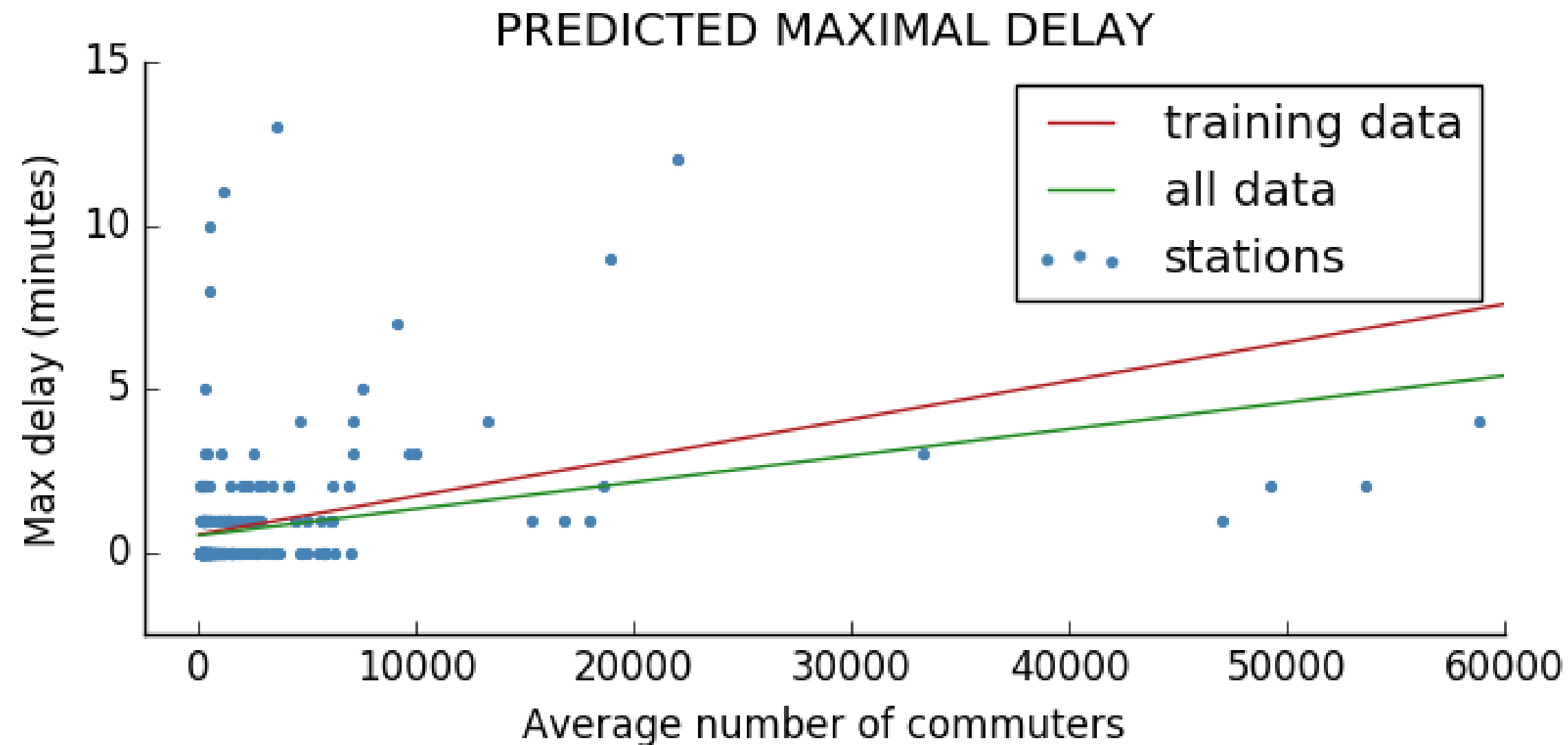
Net delay: difference of departure delay and arrival delay

INFLUENCE OF COMMUTERS



Max delay: the maximal netto delay per station in the data set

INFLUENCE OF COMMUTERS



- Simple linear regression
- Slope has significant p -value!

RECOMMENDATIONS

- More trains with less stops.
- Incentive not to travel during rush hour.



THANK YOU

SPARK PROJECT FOR NMBS / 14.12.2017