

# Estimation of Run Times in a Freight Rail Transportation Network

by

**Kunal Bonsra**

Master of Business Administration, 2010  
University of Iowa  
Bachelor of Science in Marine Engineering, 2003  
Birla Institute of Technology and Science

and

**Joseph Harbolovic**

Bachelor of Science in Chemistry, 2008  
Bachelor of Science in Computer Sciences, 2007  
The University of Texas at Austin

ARCHIVES

Submitted to the Engineering Systems Division in Partial Fulfillment of the  
Requirements for the Degree of

**Master of Engineering in Logistics**

at the

**Massachusetts Institute of Technology**

June 2012

© 2012 Kunal Bonsra and Joseph Harbolovic. All rights reserved.

The author hereby grants to MIT permission to reproduce and to distribute publicly paper and electronic copies of this document in whole or in part.

Signature of Authors .....

Master of Engineering in Logistics Program, Engineering Systems Division  
May 11, 2012

Certified by ....

Dr. Basak Kalkanci  
Postdoctoral Associate, Center for Transportation and Logistics  
Thesis Supervisor

Prof. Eva M. Ponce Cueto  
Visiting Professor, Center for Transportation and Logistics  
Thesis Supervisor

Accepted by .....

Prof. Yossi Sheffi  
Professor, Engineering Systems Division  
Professor, Civil and Environmental Engineering Department  
Director, Center for Transportation and Logistics

# **Estimation of Run Times in a Freight Rail Transportation Network**

by

Kunal Bonsra and Joseph Harbolovic

Submitted to the Engineering Systems Division in Partial Fulfillment of the  
Requirements for the Degree of Master of Engineering in Logistics

## **Abstract**

The objective of this thesis is to improve the accuracy of individual freight train run time predictions defined as the time between departure from an origin node to arrival at a destination node not including yard time. A correlation analysis is conducted to identify explanatory variables that capture predictable sources of delay and influence run times for use in a regression model. A regression model is proposed utilizing the following explanatory variables: rolling historical average, congestion window, meets, passes, overtakes, direction, arrival headway, and departure headway to predict train run times. The performance of the proposed regression model is compared against a baseline simple historical averaging technique for a two year period of actual train operational data. The proposed regression model, though subject to specific limitations, offers substantial improvements in accuracy over the baseline technique and is recommended as justifying further exploration by the railroad to ultimately enable more accurate train schedules with subsequent improvements in railroad capacity, customer service, and asset utilization.

Thesis Supervisor(s): Dr. Basak Kalkanci and Prof. Eva M. Ponce Cueto  
Title(s): Postdoctoral Associate and Visiting Professor

## Acknowledgements

We would like to thank the following individuals for their contributions to our thesis:

- Our thesis sponsors
  - Through their experience and knowledge, we have gained a great deal of knowledge about the freight railroad industry.
- Dr. Basak Kalkanci and Prof. Eva M. Ponce Cueto
  - Their efforts and long hours spent in guiding, discussing and reviewing the results were critical to success of this research.
- Joseph - My dear wife Amanda for getting me to Boston and MIT
- Kunal – My family and friends for always supporting me

## Table of Contents

Abstract.....	2
Acknowledgements.....	3
Table of Contents.....	4
List of Figures.....	5
List of Tables.....	6
Chapter 1: Introduction.....	7
1.2: Research Objective.....	8
1.3 Thesis Overview.....	8
Chapter 2: Literature Review.....	9
2.1: Railroad Capacity.....	9
2.2: Railroad Scheduling.....	10
2.3: Delays.....	13
2.3.1: Sources of Predictable Delays.....	14
2.3.2: Current approaches to estimate delays.....	16
Chapter 3: Methodology and Data Collection.....	18
3.1 Introduction.....	18
3.1.1. Potential Explanatory Variables.....	19
3.2: Data Collection.....	20
3.2.1 Description of Historical Run Data.....	21
Chapter 4: Data Analysis and Findings.....	25
4.1 Analysis.....	26
4.2 Calculated Explanatory Variables.....	26
4.3 Correlation Analysis.....	29
4.4 Regression Model.....	29
4.5 Findings.....	32
Chapter 5: Comparison of Methods and Limitations.....	33
5.1: Comparison of Methods.....	33
5.2: Robustness.....	34
5.2.1 North Dakota.....	34
5.2.2. Oregon.....	39
5.3: Limitations.....	44
Chapter 6: Conclusion and Future Direction.....	46
6.1 Conclusion.....	46
6.2 Future Direction.....	47
References.....	49

## List of Figures

Figure 2-1 Scheduling Process.....	11
Figure 3-1 Train volume on the Missouri track segment from Jan. 2010 – Jan. 2012.....	22
Figure 3-2 Mean run times on the Missouri track segment from Jan. 2010 – Jan. 2012 .....	23
Figure 3-3 Directional dependency of run times on the Missouri track segment .....	24
Figure 5-1 Train volume on the North Dakota segment from Jan. 2010 – Jan. 2012 .....	35
Figure 5-2 Mean run times on the North Dakota segment from Jan. 2010 – Jan. 2012.....	36
Figure 5-3 Directional dependency of run times on the North Dakota segment.....	36
Figure 5-4 Train volume on the Oregon segment from Jan. 2010 – Jan. 2012.....	40
Figure 5-5 Mean run times on the Oregon segment from Jan. 2010 – Jan. 2012.....	41
Figure 5-6 Directional dependency of run times on the Oregon segment.....	41

## List of Tables

Table 3-1 Mean run times for the Missouri track segment .....	23
Table 4-1 Correlation coefficients for potential explanatory variables.....	29
Table 4-2 Regression model(s) summary of fit for Missouri segment.....	30
Table 4-3 Regression model summary for low priority trains on Missouri segment.....	30
Table 4-4 Regression model summary for medium priority trains on Missouri segment.....	31
Table 4-5 Regression model summary for high priority trains on Missouri segment.....	31
Table 5-1 Regression model performance for Missouri segment .....	33
Table 5-2 Mean run times for the North Dakota track segment.....	35
Table 5-3 Regression model(s) summary of fit for North Dakota segment.....	37
Table 5-4 Regression model summary for low priority trains on North Dakota segment .....	37
Table 5-5 Regression model summary for medium priority trains on North Dakota segment .....	38
Table 5-6 Regression model summary for high priority trains on North Dakota segment .....	38
Table 5-7 Regression model performance for North Dakota segment.....	39
Table 5-8 Mean run times for the Oregon track segment.....	40
Table 5-9 Regression model(s) summary of fit for Oregon segment.....	42
Table 5-10 Regression model summary for medium priority trains on Oregon segment .....	42
Table 5-11 Regression model summary for high priority trains on Oregon segment .....	43
Table 5-12 Regression model performance for Oregon segment.....	43
Table 5-13 Regression model performance for Missouri segment schedule data.....	45

## Chapter 1: Introduction

Railroads are a critical element of transportation infrastructure and subject to growing demand. According to a recent Federal Highway Administration report, total freight movements on rail are projected to rise from an estimated 16.9 billion tons in 2010 to 27.1 billion tons in 2040 — almost a 61 percent increase (AAR, 2011a).

In terms of monetary impact, according to the U.S. Department of Commerce, freight railroads generate nearly \$265 billion in total economic activity each year (AAR, 2011b). This includes direct and indirect effects. IBIS World forecasts, that the industry's revenue will grow at an average annual rate of 3.4% reaching \$90.0 billion by 2016. In 2012 alone, IBIS World projects revenue will increase 2.6% to \$78.2 billion.

In terms of the division of revenues, the industry is still heavily dependent on the mining sector; 45% of all freight carried by rail is coal. Over the last 15 years, railroads have diversified to other lines of revenues such as intermodal transportation; this market has grown strongly and holds a great deal of promise. According to the AAR, intermodal transportation accounted for approximately 22 percent of U.S. rail revenue in 2007. In fact, intermodal traffic has risen from 3 million trailers and containers in 1980, to more than 12 million in 2007 (AAR, 2008).

Intermodal traffic did decline due to the recession in 2008 and 2009, but has since rebounded to more than 11 million units in 2010.

Class 1 railroads are large freight railroads that own the majority of track in North America. According to the AAR, to be classified as Class I, railroads needed to have minimum carrier operating revenues of \$346.8 million, \$359 million, \$401.4 million and \$378.8 million in 2006, 2007, 2008, and 2009. These railroads will face significant capacity constraints due to the long term projected growth of freight traffic. To accommodate this growth, significant investment is

required in capacity and changes need to be implemented in railroads' operational practices. There is much room for improvement. According to a market study conducted by Hertenstein and Kaplan in 1991 and cited in Hallowell and Harker (1998), "a 1% improvement in the reliability of cargo delivery time could yield as much as a 5% revenue increase in several markets". Present operational practices need to be improved by more efficient prediction of the train run times and accurate train scheduling.

## **1.2: Research Objective**

The main goal of this thesis is to develop a model to help the partner railroad, a major Class 1 U.S railroad, improve the accuracy of run time predictions. Increased accuracy when estimating run times allows the partner railroad to drive improvements in capacity utilization and punctuality/reliability of its operations. To achieve this objective the following steps were undertaken:

- Investigation into railroad capacity planning, scheduling and stochastic characteristics of delays.
- Identification of underlying factors that result in predictable delays using statistical analysis.
- Development of an analytical prediction model to predict the trains run times.
- Validation of the resulting model in three independent track segments.

## **1.3 Thesis Overview**

Chapter 2 reviews the literature on state-of-the-art railroad operations, scheduling, and delays. We also discuss the pros and cons of commonly used modeling approaches for train delays.

Chapter 3 discusses the data methodology and the data collection approach that is used to enable the data analysis in Chapter 4. Chapter 4 discusses the results of the statistical modeling of parameters influencing trains delays and train run times. A regression model is presented to enable accurate prediction of train run times based on the partner railroad scheduling information. Chapter 5 compares the regression model performance in predicting accurate train run times against that of a simple historical average baseline technique. Limitations of implementing the model within the existing partner railroad scheduling framework are also discussed. Chapter 6 presents conclusions and areas of further research.

## **Chapter 2: Literature Review**

The literature review is organized in two sections. The first section addresses railroad capacity and scheduling, we discuss tools used to analyze capacity and plan scheduling. The second section gives a general view of how different parameters affect delays. An understanding is developed around the different type of variables that influence delays and how these delays have a negative impact on actual train run times. Finally current methods for analyzing delays and predicting train run times are presented.

### **2.1: Railroad Capacity**

Capacity estimation is a long running challenge for the railroad industry. A key motivation for calculating railroad capacity is to validate the feasibility of current train timetables in the existing railroad network. Additionally, it helps when evaluating the need for new track infrastructure, new network design, and modernization of signaling equipment. (Yuan, 2006)

As stated in (Abril et al., 2008), “The goal of capacity analysis is to determine the maximum number of trains that would be able to operate on a given railway infrastructure, during a specific time interval, given the operational conditions.”. (Abril et al., 2008) continues and identifies a number of different railroad capacity definitions:

- **Theoretical Capacity** – The maximum number of trains that can operate on a track segment in a theoretical perfect environment. Generally, this is an upper bound for capacity. This assumes homogenous conditions: similar trains and equal spacing between trains.
- **Practical Capacity** – A number of basic assumptions of theoretical capacity are relaxed. This is the practical limit of train volume that can be moved on the track segment regularly in a reliable fashion. It can be between 65% and 70% of total theoretical capacity of the line.
- **Used Capacity** – The capacity utilized by the existing railroad schedule.
- **Available Capacity** – This difference between practical and used capacity, indicating additional volume the track segment could handle safely.

## **2.2: Railroad Scheduling**

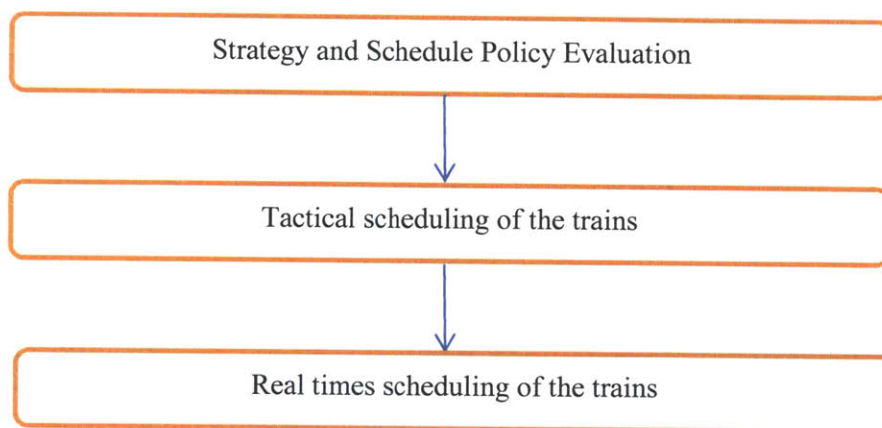
Railroad networks have very high utilization rates with dense traffic patterns. A great deal of effort is needed to create and execute conflict free schedules. Train scheduling requires the development of timings and ordering plans for particular trains based on their routes and demand. Train timetables are created to schedule trains to meet different priorities and other operational constraints. As stated in (D’Ariano et al., 2007), “Timetable development is a complex problem in which a compromise between capacity utilization and timetable robustness has to be provided”. Train timetables typically contain slack time to address routine delays prevalent in the railroad network. (Arcot 2007) presents a useful analysis of slack time:

“Slack time in timetables absorb minor delays and limit the propagation of delays in the network. As the amount of slack in timetables is increased, the train network becomes more stable. The stability of rail networks can be defined as the ability to recover to the original schedule after disruptions to the schedule...However, increasing slack in timetables reduces the number of trains that can be scheduled in the network, which in turn reduces the infrastructure capacity utilization.”

Slack time is incorporated in train run times via dwell times.

Train scheduled run times = Free train run times + Dwell times

The generic steps which are involved in the scheduling process are shown below:



**Figure 2-1 Scheduling Process**

Generally there are two types of scheduling: tactical and operational scheduling (D'Ariano, 2010)

- Tactical scheduling – Requires the creation of a master schedule and is focused on scheduling at the network level. The objective is to satisfy demand while allotting time for maintenance activities and other predictable events. Tactical schedules are generally developed months in advance on a corporate or regional level. Tactical scheduling models usually perform optimization on both train routes and train timings. This approach is most commonly used in Europe.

- **Operational Scheduling** – The operational schedule is created shortly before actual train departures. As a result, sometimes the tactical schedule might be substantially different from the actual schedule followed by operations personnel. North American and Australian railroads generally use this strategy. A draft time table is generally fixed, but train timings are not, which allows for significant flexibility in the system.

### **Railroad Current Scheduling Practices**

North American freight railroads operate multiple trains on a particular track segment at differing speeds. This increases the heterogeneity of the system but also creates multiple conflicts while scheduling the traffic. There are multiple methodologies to address the heterogeneity issue. One such method (Harrod, 2009) schedules a subset of the trains with linear weighted utility values.

### **Dispatch Simulation Software:**

A number of different tools are used by railroads to schedule operations and each railroad has its own philosophies about how to build traffic, (Dingler, 2010) discusses some of the methodologies and tools used by the railroads. Railroads can each utilize the same tool configured differently for their various operational approaches. Class 1 railroads can also develop their own in-house technologies to match their unique specifications. The primary tool used by railroads is the Rail Traffic Controller (RTC) from Berkeley Simulation Software.

### **Rail Traffic Controller (RTC):**

(Dingler, 2010) presents an in-depth discussion of RTC:

“Delay is a measure of level of service and is the primary output from RTC... RTC requires the user to input the infrastructure and traffic into the model. Using the built-in train performance calculator (TPC) and meet-pass logic RTC attempts to dispatch the

traffic in a way similar to an actual dispatcher (Wilson 2010). The current generation of the software resolves conflicts using priority based dispatching; when there are conflicts, the logic seeks alternative routes for the lower priority train (Lai 2008)... since its introduction in 1999 it has become widely accepted by railroads, consultants and government agencies and is the de facto industry standard of the North American railroad industry.”

### **2.3: Delays**

Train run times as discussed earlier are comprised of free train run times and dwell times. Dwell times include slack times that try to accommodate trains delays. Usage of dwell time creates a great deal of inefficiency in the system by significantly increasing actual train run times. The causes of delay that influence actual run times are reviewed in this section. There are two types of delay: predictable and unpredictable. Predictable delays are the scheduled events that are incorporated in run times as a buffer to allow for traffic conflicts. Unpredictable delays are unscheduled events that are random and result in significant instability in the network.

Unscheduled events can be due to flooding, collision, malfunctioning equipment, crew delays, amongst others. A significant portion of the existing research explores identifying predictable delays (Chen & Harker, 1990), (Hallowell & Harker, 1998), (Hallowell & Harker, 1996).

The literature review and this thesis are focused on identifying the magnitude and nature of the variables influencing predictable delay. Unpredictable delay is not analyzed due to its random nature. The literature review for delays is divided into two sections. The first section identifies sources of predictable delays. The second section discusses methodologies to estimate delays and the current industry approaches to estimating capacity.

### 2.3.1: Sources of Predictable Delays

Predictable delays depend on both dynamic and static factors. Dynamic factors are dependent on train characteristics and vary based on individual train type. Static factors are independent of train characteristics and are instead influenced by the track configuration and railroad characteristics.

#### **Dynamic factors:**

Delays are highly variable and depend largely on the volume and type of traffic on the track segment. (Dingler et al., 2010) used a simulation to predict increased traffic results in higher number of meets and passes between trains. A meet is defined as when two trains in opposite directions of travel encounter one another on the same track, resulting in one train stopping for another. A pass is defined as when two trains are moving in same direction and one train stops for other train. Headway is the distance maintained between trains running in the same direction. As a result of meets and passes, train headways are reduced. This reduction in headway can have a cascading effect of propagating delay as trains may be slowed to maintain minimum headway distances. The delay-volume relationship is also dependent on the traffic mix of a particular track segment. (Lai et al., 2010) discusses the relationship between train volume and delay with respect to single track segments:

“On a single-track line, the effect of additional trains on delay is not linear. Instead, the relationship between train volume and delay is exponential, with each train type and traffic mix having its own particular functional relationship (7, 17, 18).”

Apart from volume, train type heterogeneity also impacts delay (Dingler et al., 2009). Different train types have significantly different operating characteristics, generally represented via train velocity, reflecting the varying business needs of each cargo type. Generally with all things

being equal heterogeneous traffic results in greater delays as compared to homogenous traffic. Homogenous networks generally eliminate passes which helps in reducing the network variability.

Additional sources of delay due to heterogeneous traffic are:

1. A faster train trailing a slower train.
2. A train with faster acceleration trailing a slower acceleration train.
3. Lower priority trains waiting for a higher priority train to pass to resolve a meet
4. Lower priority trains slowing down to allow a higher priority train to pass
5. Lower average speeds resulting in an increased number of meets.
  - a. This can be due to lower speeds, lower power, and lower priority or some combination thereof.

Different train characteristics have a variable impact on heterogeneity. Differences in train priority have the most significant influence on delays. Delay is also dependent on the percentage of different train types, or the relative mix of priorities on a particular track segment (Dingler et al., 2009). The impact of the heterogeneity in speed and power-ton ratio is relatively low.

Homogeneity with respect to speed can reduce delay, but will not have significant influence on delays as trains generally are not traveling at their optimum or maximum velocities. With respect to differences in power-to-ton ratios, higher power generally reduces the time required to accelerate after stops. Higher power-to-ton ratios generally have diminishing effect as the maximum power configuration of a train is approached. The maximum change in performance is generally encountered in trains with low power-to-ton ratios.

### **Static factors that influence predictable delays:**

The length of the track segment, siding length and spacing, signal spacing, and percent of the segment single tracked also contribute to train run times (Krueger, 1999). With respect to this thesis and its data analysis, these static factors are considered constant along the particular track segments analyzed and are not pertinent to the analysis.

### **2.3.2: Current approaches to estimate delays**

A number of different methodologies are presently utilized to estimate delays:

1. Simulation Models.
2. Queuing Models.
3. Stochastic Models.
4. Regression Models

Simulation based techniques to estimate delays were used by (Murali, Dessouky, Ordóñez, & Palmer, 2010). Their research included simulation runs to represent train movements with regression analysis conducted on the results of the simulations to predict delays. Their model also tries to predict the effect of network design, network topology, and differing traffic parameters on delay.

Researchers have also proposed models based on queuing theory to estimate delays. Queuing models help to compute the buffer times required to minimize scheduled waiting times. In these models train arrival times and service times are assumed to be independent random variables. (Greenberg, Leachman, & Wolff, 1988) used queuing models on a low speed, single track rail network to predict delays. Poisson distributions were used to model train departures and slow speeds allowed them to model trains with limited headways. (Wendler, 2007) presents an approach to predict scheduled wait times by using a semi-markovian queuing model.

Knock-on delay effects between two trains on a single track have also been explored recently (Carey & Kwiecinski, 1994). In this research stochastic approximations of the relationship between scheduled headways and knock-on delays using a non-linear regression model was conducted. A knock-on delay is the indirect delay suffered by later trains after a primary train gets delayed. Primary train delay creates a cascading effect for other trains, resulting in subsequent trains encountering knock-on delays. Regression models have also been used to evaluate train delays (Flier, Graffagnino, & Nunkesser, 2009). The researchers utilized a greedy step AIC algorithm of Venables and Ripley to create a model that could predict delays. Probability models that provide a realistic estimate of knock-on delays and use track capacity have also been proposed (Yuan & Hansen, 2007). Mean knock-on delays increase when schedule slack time decreases. A switchable dispatching policy for a double-track segment was proposed by (Mu & Dessouky, 2010) suggesting it would be advantageous, potential reductions in delay of up to 65% for fast trains, to allow a fast train to pass a slow train by using opposite direction track if the track is empty. Some researchers also discussed using stochastic models for delay propagation with forecasts of arrival and departure events (Berger & Gebhardt, 2011). In a similar fashion, researchers estimated flight departures using a spline approach (Tu, Ball, & Jank, 2008). This methodology appears relevant and may be applicable to the railroad industry. The researchers grouped delay factors in three major categories: seasonal trend, daily propagation pattern, and random residual. To capture historical trends they utilized a smoothing spline model. They assumed a mixture model for the residuals and estimated mixture components using an expectation maximization algorithm. Of particular relevance to this thesis was the research conducted in (Gorman, 2009). In this research statistically analysis was utilized to identify candidate explanatory variables across a

variety of track segments for a U.S. freight railroad. To understand the impact of heterogeneity in train characteristics, the author classified trains in three priority groups: high, medium and low. Meets, passes and overtakes were identified as primary causes of predictable delay.

## **Chapter 3: Methodology and Data Collection**

### **3.1 Introduction**

Actual train run time is calculated as the time of departure from an origin station to the time of arrival at a destination station, with no yard time included. A variety of different methodologies exist to estimate delays and train run times, such as simulation, queuing, probabilistic and regression modeling techniques. With the dataset available from the partner railroad regression modeling is a suitable technique to analyze the individual impact of the different dynamic variables on both delay and train run times.

Regression uses past history to understand the underlying patterns that have an impact on train run times. Regression analysis is used to estimate the conditional expectation of the dependent variable, in this case, train run times, given different independent variables. In order to identify useful explanatory variables for the historical dataset available from the partner railroad from the list of potential explanatory variables identified in the literature a correlation analysis is conducted. Those variables with sufficient explanatory power will be utilized in the regression model.

From the insights gained a final predictive and a forecasting model is created. The final regression model takes into account the impact and significance of the different independent

variables and predicts the run times based on actual historical data or theoretical train departure schedules.

### **3.1.1. Potential Explanatory Variables**

From the literature review the following parameters were selected as potential explanatory variables:

- Train Volume – The total number of trains operating on the same track segment during a particular train’s transit window.
- Train Priority – Different priorities are assigned to trains depending on their assigned importance with respect to business needs.
- Train Direction – This variable accounts for the direction in which a train is travelling.
- Meets – The number of trains traveling in the opposing direction a particular train encounters during its transit.
- Passes – The number of trains a particular train passes while traveling in the same direction during its transit.
- Overtakes – The number of trains a particular train is overtaken by when traveling in the same direction.
- Arrival Headways –The difference in time between the arrival of the last train at the destination station and the train in question.
- Departure Headways – The difference in time between the departure of the last train at the origin station and the train in question.
- Engine horsepower – Information regarding power-to-ton ratios was unavailable. Trains were assumed to have sufficient power as to not cause delay

- Days and month of the year – Calendar day and month of the year the particular train departed.

The final model will only contain those variables deemed to be statistically significant.

The final regression model will be tested on three different track segments. The predicted values obtained from this model will be compared against actual train run times. To validate the results, we will calculate the root mean square error for our model across all the train segments.

### 3.2: Data Collection

The partner railroad provided a dataset to support this thesis from a 141 mile, 97% single tracked network segment, running East to West in Missouri. This segment was selected out of the entire 31,000 mile network as a sample data set for the following reasons:

- Predominately single tracked, only 3% of the route is double tracked
- Strong directional influence on run times
  - The route has wide changes in elevation, from 800 to 1600 feet
  - Full trains tend to transit from West to East and have measurably longer run times versus empty trains which tend to transit from East to West
- Broad cross-section of train priorities
  - Trains have a unique priority which captures their relative importance and unique operating characteristics. Having a broad mix of priorities captures the heterogeneous nature of freight rail networks

This dataset consists of historical actual network operational data over a two year period from January 2010 to January 2012. This data is herein referred to as the “*Historical Run Data*”.

### 3.2.1 Description of Historical Run Data

The unique key into the data is a combination of a specific train ID and scheduled departure date. For each of the discrete train elements and relevant to the analysis, the data contains information specifying crew route (or direction), priority, total run time, departure time, and arrival time. Not relevant to the analysis but also available were origin and destination station(s) along with station ID(s), line description, division name, subdivision name, responsibility center, supervisor in charge, and discrete delay events by train ID with information specifying delay amount (in minutes), date and time of delay event, and station where the delay was logged. The period of the data is January 2010 through January 2012. A multi-year period was selected to capture any seasonality effects that may be present in the underlying data. Critically the reason for the delay was not deemed reliable and prevented the development of an event-based probabilistic model to estimate train run times.

#### **Train IDs:**

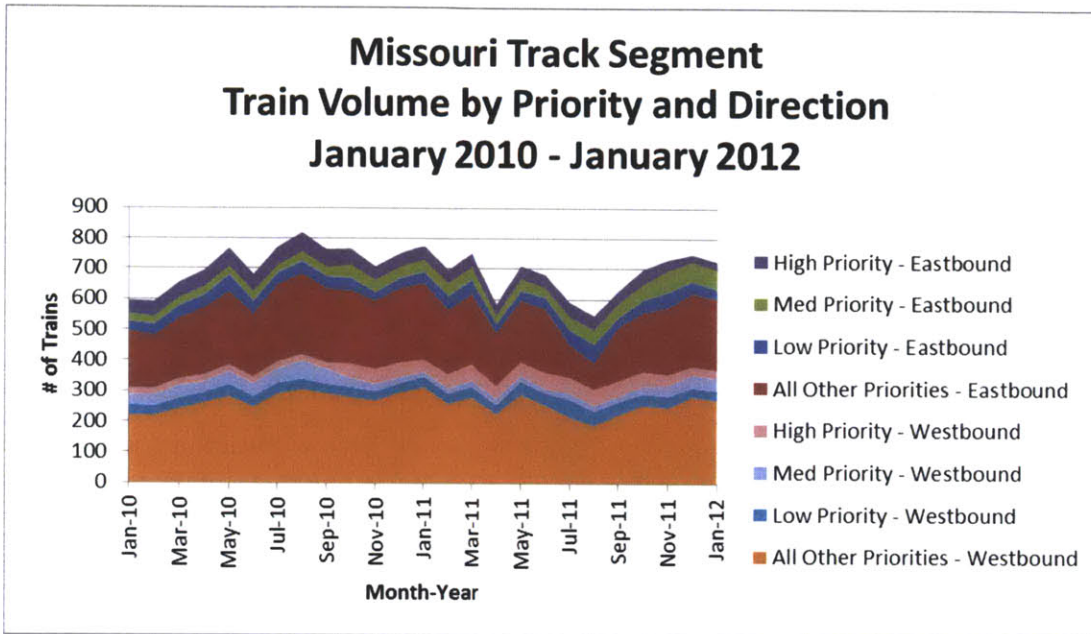
The train ID coupled with a scheduled departure date is the unique identifier in the dataset. A train ID can repeat in the dataset and represents a particular reoccurring train. Within the train ID the following can all be isolated: priority, origin and destination station codes, and daily sequence number (in the event more than one of the train is scheduled per day).

#### **Train Priority:**

The railroad utilizes more than fifteen priority codes to reflect the varying scheduling priorities of the trains within its network. The cargo and loaded status (full or empty) can be derived from the train priorities. The data includes all priority types however this research will focus only on the top three priority trains, designated as low, medium, and high with all other priorities aggregated into the “Other” priority classification. Priorities other than the top three are treated

as equivalent with respect to scheduling preference. The scope of this thesis is limited to prediction for the top three priorities only.

Figure 3-1 shows the train volume transiting the Missouri track segment and the priority and direction mix over the two year historical data period. The volume drop observed in the eastbound traffic in August of 2011 corresponds to specialized track maintenance.



**Figure 3-1 Train volume on the Missouri track segment from Jan. 2010 – Jan. 2012**

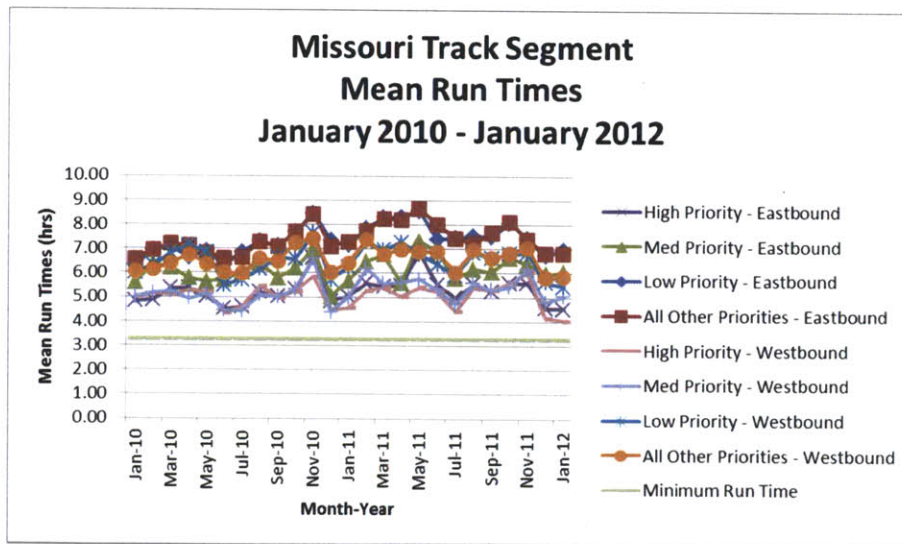
The historical run data captures all trains that operated on the Missouri track segment during the period in question. The minimum run time shown in Figure 3-2 represents the best possible transit speed for the track segment independent of direction and is calculated based on supporting information documenting varying permanent speed restrictions along the track segment.

The mean run times vary significantly by priority and direction as can be observed in Table 3-1.

Priority	Eastbound	Westbound
High	5.21	5.04
Med	6.15	5.19
Low	7.32	6.40
All Other	7.33	6.48

**Table 3-1 Mean run times for the Missouri track segment**

Figure 3-2 shows the mean run times by direction and priority for trains transiting the Missouri track segment over the two year historical data period. While the volume for eastbound traffic was reduced in August 2011 the trains that did transit actually enjoyed an improvement in runtimes notwithstanding the track maintenance that occurred during this period.



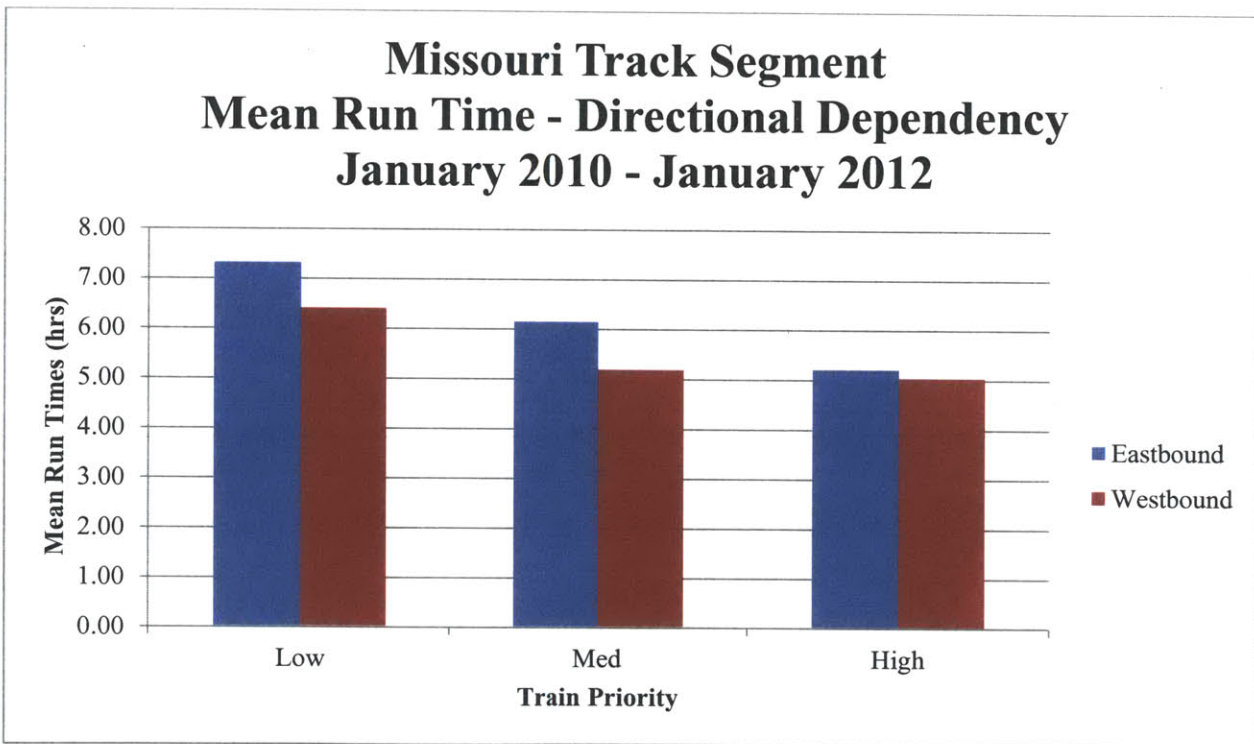
**Figure 3-2 Mean run times on the Missouri track segment from Jan. 2010 – Jan. 2012**

As can be observed in the figures above there are noticeable changes in volume Figure 3-1 without meaningful changes in run times Figure 3-2, with the exception of the track maintenance period in August 2011.

**Crew Routes:**

The crew route specifies the direction of travel of a particular train. The large Class I railroads in the United States tend to be concentrated in specific geographic regions. As such, on the track segment in question loaded trains tend to transit East-bound while empty trains tend to transit West-bound. Operational practice at the railroad dictates that empties are moved at the highest possible velocity though the loaded trains take priority.

Figure 3-3 below shows the dependency of run times on the direction of travel in the Missouri track segment.



**Figure 3-3 Directional dependency of run times on the Missouri track segment**

**Scheduled Departure Date:**

Each train has a scheduled departure date which is derived from a strategic customer-commitment schedule produced several years in advance. The actual departure date of a train

generally occurs within +3 days of the scheduled departure date. This deviation is strongly right-tailed and can reach a maximum of 25 days. As such the use of the scheduled departure date is limited to its role as a component of the unique key.

**Actual Departure Date:**

Each train ID has an actual departure date-time which corresponds to the date and time of departure from the origin stations. This value represents when the train in question has left the station and does not include any yard time.

**Actual Arrival Date:**

Each train ID has an actual arrival date-time which corresponds to the date and time of arrival at the destination station. This value represents when the train in question arrived at the destination station and does not include any yard time.

**Total Run Time:**

The run times referred to in this thesis are the length of time it takes a train to transit the track segment from time of departure, herein referred to as  $T_D$ , at the origin station to time of arrival, herein referred to as  $T_A$  at the destination station. It is important to note the run time predictions do not include yard time at either the origin or destination stations, only the transit time between the stations.

## **Chapter 4: Data Analysis and Findings**

The analysis of the historical run data proceeds with a removal of erroneous records followed by calculation of the potential explanatory variable fields selected in the literature review for use in

a correlation analysis. Those variables found to have significant explanatory power are utilized in the regression model construction.

## **4.1 Analysis**

To identify erroneous records in the dataset trains with total run times outside fixed maximum or minimum values were removed. Any train record with a total train run time in excess of twelve hours was removed. This was done as federal law mandates a crew may operate a train for a period of twelve hours before a different crew must take over. Should a crew change, occur the run time would need to be recalculated.

The minimum run time is calculated for the Missouri track segment based on the length of track, maximum permitted velocity, and any permanent speed restrictions that may be in place. For the Missouri track segment the minimum transit time was 3.28hrs. A logarithm transformation was applied to the train run times, to compensate for the right-tailed nature of the run times distribution. In place of an untransformed run time, the logarithm of the run time is used both as the dependent variable in the regression model and to build a rolling run time average.

## **4.2 Calculated Explanatory Variables**

### **Rolling Run Time Average:**

Past train performance is believed to be a predictor of future train performance. The selection of a prior performance period is somewhat arbitrary and was fixed at 30 days for this analysis.

The rolling run time average captures the influence of train priority on run times. The average of the logarithms of total run time for trains with the same direction and priority from the previous 30 day period, as based on the original departure date and time of the train in question is computed for each train.

**Departure Congestion Window:**

Statistical significance is present for measures of congestion, specifically the number of trains that have or will depart within some time period of the departure of the train in question. The departure congestion window captures the impact of train volume on run times. The strongest explanatory variable for this behavior was found using a twelve hour window, which captures a full day period split evenly before and after train departure. Alternative windows of four and six hours were also calculated but found to have less explanatory power in the correlation analysis. The congestion and thirty day average periods were mutually exclusive.

**Meets:**

The number of meets is defined as the number of trains in the opposing direction a particular train encounters while in transit. This is calculated for a specific train by counting the number of trains in the two year period of historical data that meet the following conditions:

- 1) Must be in opposing directions of travel
- 2) Must have overlapping transit windows (i.e., departure and arrival windows)

**Passes:**

The number of passes is defined as the number of trains a particular train passes traveling in the same direction while in transit. This is calculated for a specific train by counting the number of trains in the two year period of historical data that meet the following conditions:

- 1) Must be in the same direction of travel
- 2) Must have overlapping transit windows (i.e., departure and arrival windows)
- 3) Must have departed before the train in question
- 4) Must have arrived after the train in question

**Overtakes:**

The number of overtakes is defined as the number of trains a particular train is overtaken by when traveling in the same direction. This is calculated for a specific train by counting the number of trains in the two year period of historical data that meet the following conditions:

- 1) Must be in the same direction of travel
- 2) Must have overlapping transit windows (i.e., departure and arrival windows)
- 3) Must have departed after the train in question
- 4) Must have arrived before the train in question

**Direction:**

The direction of the train plays a meaningful role in determining the run time. The business realities of the railroad generally mean that West-bound trains are empty while East-bound trains are loaded. Loaded trains take priority; however empty trains are run at a higher velocity.

Direction is a binary field.

**Arrival Headway:**

The arrival headway is defined as the difference in time between the arrival of the last train at the destination station and the train in question. Arrival headway is used to capture any delay effects from congestion at the destination station.

**Departure Headway:**

The departure headway is defined as the difference in time between the departure of the last train at the origin station and the train in question. Departure headway is used to capture any delay effects from congestion at the origin station.

### Day and Month:

The day of month and month of year are both calculated from the departure date and time to determine any seasonal influence on train run times.

### 4.3 Correlation Analysis

A standard correlation analysis was conducted in the statistical software JMP from SAS Institute, Inc. to quantify the statistical significance of each of the potential explanatory variables in the historical run data. Day and month were found to have little statistical relevance and were removed from consideration for use in the regression model. Table 4-1 shows the results of the correlation analysis listing the potential explanatory variables tested and their respective correlation coefficients. Meets have the highest correlation value, followed closely by the historical rolling average, and then overtakes and passes.

Explanatory Variable	Correlation Coefficient
Rolling Run Time Average	0.49
Departure Congestion Window	0.09
Meets	0.56
Passes	-0.24
Overtakes	0.34
Direction	0.19
Arrival Headway	0.07
Departure Headway	-0.11
Day	-0.01
Month	0.01

Table 4-1 Correlation coefficients for potential explanatory variables

### 4.4 Regression Model

A regression model was constructed in statistical software JMP from SAS Institute, Inc. using the explanatory variables determined to be significant. To capture the distinct operational characteristics of each of the top three priorities, regression models were calculated for each train

priority individually which generated unique term coefficients. Table 4-2 shows the performance of each of these priority specific regression models. Interestingly the regression model explained a greater portion of the variability for the low priority trains.

Summary of Fit	Low Priority	Med Priority	High Priority
RSquare	0.53	0.48	0.37
RSquare Adj	0.53	0.48	0.37
Root Mean Square Error	0.07	0.08	0.08
Mean of Response	0.83	0.74	0.70
Observations	1870	1659	2019

**Table 4-2 Regression model(s) summary of fit for Missouri segment**

Table 4-3 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to low priority trains operating on the Missouri track segment. The number of meets carries the greatest statistical significance for low, medium, and high priority trains on the Missouri segment. The departure congestion window, number of passes, and departure headway all negatively influence run times while all other variables carry a positive sign for all priorities. Interestingly an increase in the volume of trains within the congestion window resulted in reduced run times for all train priorities.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.5065	0.0383	13.23	<0.0001*
Rolling Run Time Average	0.3217	0.0474	6.79	<0.0001*
Departure Congestion Window	-0.0049	0.0005	-10.86	<0.0001*
Meets	0.0229	0.0008	29.1	<0.0001*
Passes	-0.0462	0.0080	-5.78	<0.0001*
Overtakes	0.0718	0.0037	19.54	<0.0001*
Direction	0.0161	0.0044	3.69	0.0002*
Arrival Headway	0.0128	0.0012	11.03	<0.0001*
Departure Headway	-0.0142	0.0012	-11.81	<0.0001*

**Table 4-3 Regression model summary for low priority trains on Missouri segment**

Table 4-4 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to medium priority trains operating on the Missouri track segment. Meets continue to carry the most statistical significance, though to a lesser extent than in the case of low priority trains. The direction variable is a binary variable, with 1 representing an east-bound train and 0 representing a west-bound train.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.4463	0.0395	11.31	<0.0001*
Rolling Run Time Average	0.3175	0.0554	5.73	<0.0001*
Departure Congestion Window	-0.0038	0.0005	-7.06	<0.0001*
Meets	0.0236	0.0010	24.31	<0.0001*
Passes	-0.0367	0.0034	-10.73	<0.0001*
Overtakes	0.0942	0.0071	13.34	<0.0001*
Direction	0.0319	0.0059	5.37	<0.0001*
Arrival Headway	0.0129	0.0012	11.09	<0.0001*
Departure Headway	-0.0135	0.0011	-11.72	<0.0001*

**Table 4-4 Regression model summary for medium priority trains on Missouri segment**

Table 4-5 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to high priority trains operating on the Missouri track segment. As would be expected given the limited difference between directional mean run times for the high priority trains, the direction variable has less statistical relevance in this instance.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.4298	0.0323	13.29	<0.0001*
Rolling Run Time Average	0.3744	0.0453	8.26	<0.0001*
Departure Congestion Window	-0.0051	0.0005	-10.87	<0.0001*
Meets	0.0240	0.0009	27.57	<0.0001*
Passes	-0.0211	0.0028	-7.57	<0.0001*
Overtakes	0.1146	0.0210	5.45	<0.0001*
Direction	0.0110	0.0039	2.81	0.0050*
Arrival Headway	0.0103	0.0011	9.67	<0.0001*
Departure Headway	-0.0098	0.0011	-8.59	<0.0001*

**Table 4-5 Regression model summary for high priority trains on Missouri segment**

The explanatory variables show differing statistical significance across the three different priorities adding support to the development of models specific to priority in place of a single model for a particular track segment.

## 4.5 Findings

The regression models explain a reasonable level of variability in the train run times utilizing a set of explanatory variables that represent a broad cross section of historical operational performance and predictable sources of delay. The usage of train priority specific regression models captures the operational differences between the priorities and more accurately estimates train run times. Interestingly, the departure congestion window, or the volume of trains departing within +/- 12 hours of a particular train, for all three priorities carried a negative sign. This would seem to indicate, that as track segment volume increases, individual run times are reduced. The source of this relationship was not explored but may justify further research. The model explains the greatest percentage of the variability in the low priority trains with gradually reducing R-square values as the train priority increases. The historical run time distributions of the lower priority trains tend to have more variability which appears to be explained to a greater extent with the explanatory variables being utilized than the relatively tight high priority train run time distributions. The impact, on model explanatory power and accuracy, of removing one-off events from the historical data, such as flooding or specialized track maintenance, was not evaluated. It may be worthwhile to investigate periods of relative calm with respect to unpredictable delay causing events to get an improved sense of the model's true capabilities.

## Chapter 5: Comparison of Methods and Limitations

### 5.1: Comparison of Methods

The regression model will be utilized to estimate train run times on a train ID specific basis for all low, medium, and high priority trains in the historical run data.

These predictions were compared with the actual run times and those estimated run times found to be within +/- one hour of the actual run time were deemed “accurate”. A baseline run time prediction methodology of using a simple thirty day average, by priority and direction, to estimate run times represents a base case estimate and stand-in for the existing prediction methodology. A similar comparison of these base case estimates against the actual run times was conducted with those base case estimates found to be within +/- one hour of the actual run time deemed “accurate”. Table 5-1 shows the accuracy and percent improvement over the baseline simple historical average technique of each of the priority specific regression models. The greatest run time prediction accuracy, 73.97%, is generated by the new model for the high priority trains. This represents a substantial improvement in accuracy of 20.95% over the baseline simple average technique. The greatest improvement in accuracy over the baseline technique, 36.79%, occurred with the low priority trains. Accuracy tended to increase when historical run time variability decreased.

	Low Priority	Med Priority	High Priority
<b># of Observations</b>	1870	1659	2019
<b>Regression Model Accuracy</b>	62.83%	66.43%	73.97%
<b>Simple Average Accuracy</b>	45.94%	51.60%	61.16%
<b>% Improvement</b>	36.79%	28.74%	20.95%

Table 5-1 Regression model performance for Missouri segment

## 5.2: Robustness

To evaluate the robustness of the regression model the partner railroad provided two additional datasets containing historical run data from two alternative track segments:

- 1) A 200 mile, 84% single tracked network segment in North Dakota with a minimum transit time of 3.36 hours.
- 2) A 222 mile, 91% single-tracked network segment in Oregon with a minimum transit time of 4.07 hours.

These segments were selected as additional sample sets for the following reasons:

- Predominantly single tracked
- Directional influence on run times
- Broad cross-section of train priorities with unique mixes relative to each other and the Missouri segment
- Geographical distinct areas.

The additional datasets both consist of historical actual network operational data over a two year period from January 2010 to January 2012 with identical data fields to that of the Missouri segment historical run data. The regression model was recreated for each of the two track segments utilizing the same calculated explanatory variables to generate term coefficients specific to the low, medium, and high priority trains operating on each segment.

### 5.2.1 North Dakota

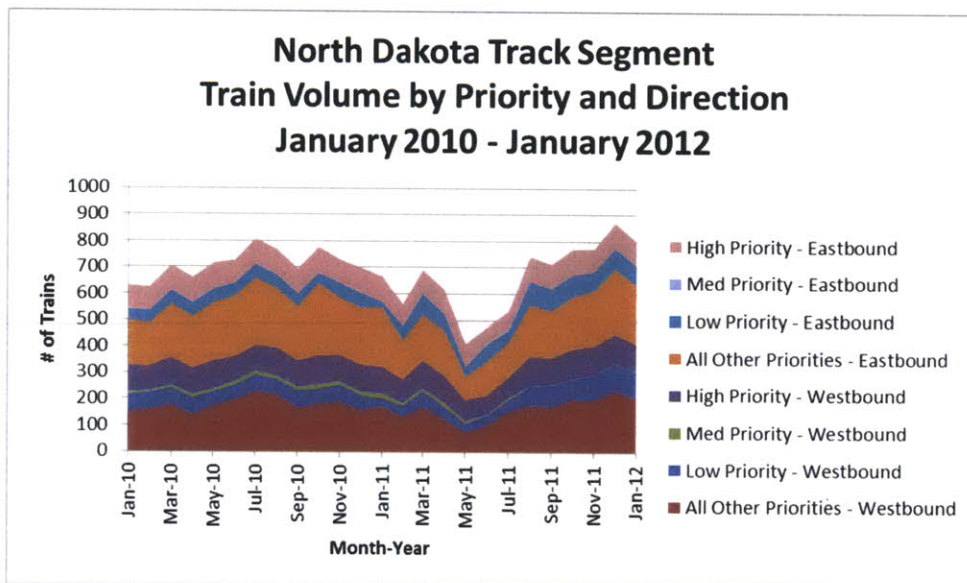
As Figures 5-1 and 5-2 show below there is a noticeable change in volume around the May 2011 timeframe with a substantial increase in run times, this was due to severe flooding. Directional dependencies are largely absent from the low and high priority trains in the North Dakota

segment, Figure 5-3. The mean run times vary somewhat by priority and direction as can be observed in Table 5-2.

Priority	Eastbound	Westbound
High	5.84	5.84
Med	6.90	8.63
Low	8.08	8.19
All Other	8.07	8.26

**Table 5-2 Mean run times for the North Dakota track segment**

Figure 5-1 shows the train volume transiting the North Dakota track segment and the priority and direction mix over the two year historical data period. The volume drop observed for all priorities and both directions in May of 2011 corresponds to severe flooding.



**Figure 5-1 Train volume on the North Dakota segment from Jan. 2010 – Jan. 2012**

Figure 5-2 shows the mean run times by direction and priority for trains transiting the North Dakota track segment over the two year historical data period. There was a general rise in run times around the May 2011 flooding event with a substantial spike in east and west bound high priority traffic.

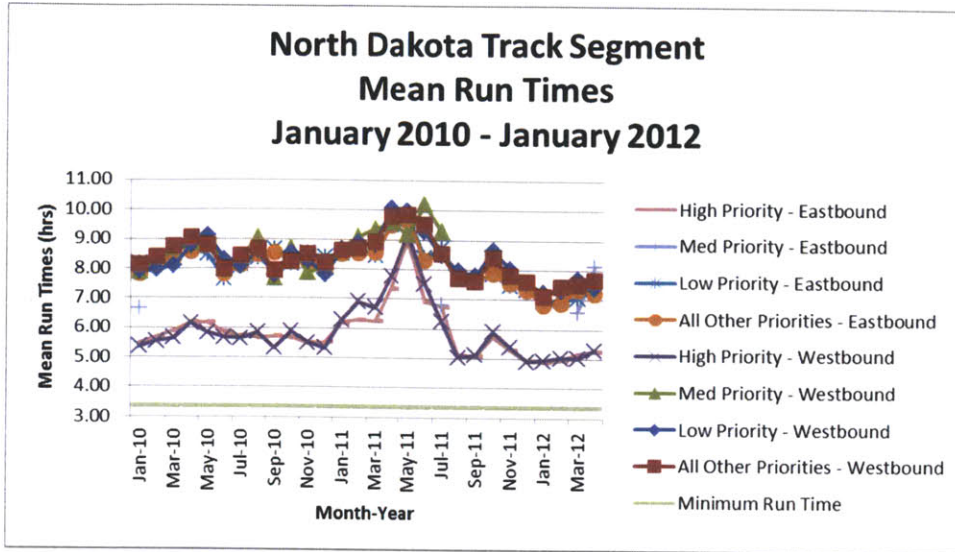


Figure 5-2 Mean run times on the North Dakota segment from Jan. 2010 – Jan. 2012

Figure 5-3 below shows the dependency of run times on the direction of travel in the North Dakota track segment.

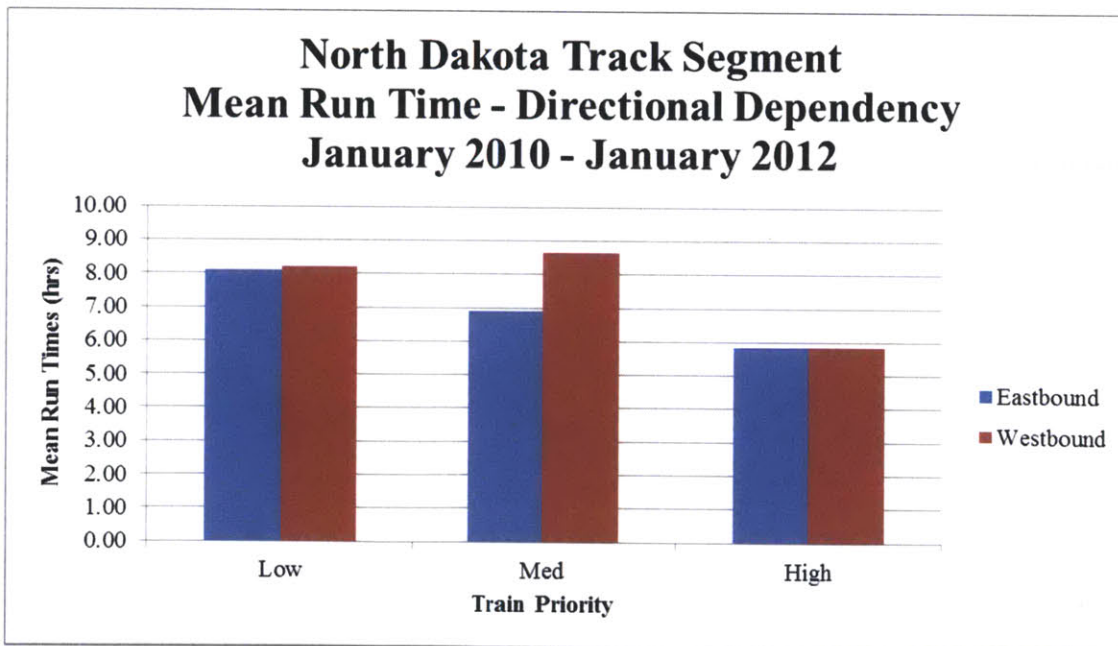


Figure 5-3 Directional dependency of run times on the North Dakota segment

Table 5-3 shows the performance of each of these priority specific regression models. The R-square value for the medium priority model exhibited an unexpected drop. The relatively low number of observations for the medium priority model would dictate an additional investigation to fully evaluate the model against this priority train type operating on the North Dakota track segment.

Summary of Fit	Low Priority	Med Priority	High Priority
<b>RSquare</b>	0.45	0.39	0.47
<b>RSquare Adj</b>	0.45	0.37	0.47
<b>Root Mean Square Error</b>	0.06	0.05	0.07
<b>Mean of Response</b>	0.91	0.93	0.76
<b>Observations</b>	2962	172	4632

**Table 5-3 Regression model(s) summary of fit for North Dakota segment**

Table 5-4 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to low priority trains operating on the North Dakota track segment. As was observed with all train priorities on the Missouri track segment, meets has the greatest statistical significance for low priority trains. Unlike Missouri, the overtake variable for low priority trains on the North Dakota track segment is relatively close in statistical significance to meets.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.2058	0.0324	6.35	<0.0001*
Rolling Run Time Average	0.7655	0.0334	22.9	<0.0001*
Departure Congestion Window	-0.0040	0.0003	-15.02	<0.0001*
Meets	0.0122	0.0005	26.03	<0.0001*
Passes	-0.0300	0.0039	-7.72	<0.0001*
Overtakes	0.0453	0.0018	25.43	<0.0001*
Direction	-0.0015	0.0020	-0.72	0.4745
Arrival Headway	0.0105	0.0006	16.3	<0.0001*
Departure Headway	-0.0118	0.0006	-18.38	<0.0001*

**Table 5-4 Regression model summary for low priority trains on North Dakota segment**

Table 5-5 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to medium priority trains operating on the North Dakota track segment. Due to the low number of observations, the statistical values for medium priority trains are viewed with some skepticism.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.5276	0.1407	3.75	0.0002*
Rolling Run Time Average	0.4769	0.1413	3.37	0.0009*
Departure Congestion Window	-0.0067	0.0012	-5.64	<0.0001*
Meets	0.0143	0.0019	7.51	<0.0001*
Passes	-0.0225	0.0119	-1.89	0.0608
Overtakes	0.0229	0.0079	2.91	0.0041*
Direction	0.0000	0.0000		
Arrival Headway	0.0041	0.0025	1.64	0.1038
Departure Headway	-0.0053	0.0021	-2.45	0.0154*

**Table 5-5 Regression model summary for medium priority trains on North Dakota segment**

Table 5-6 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to high priority trains operating on the North Dakota track segment.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.2598	0.0165	15.75	<0.0001*
Rolling Run Time Average	0.7196	0.0185	38.8	<0.0001*
Departure Congestion Window	-0.0052	0.0002	-20.9	<0.0001*
Meets	0.0131	0.0005	28.28	<0.0001*
Passes	-0.0128	0.0013	-9.58	<0.0001*
Overtakes	0.0714	0.0076	9.45	<0.0001*
Direction	0.0023	0.0020	1.16	0.2463
Arrival Headway	0.0060	0.0005	11.19	<0.0001*
Departure Headway	-0.0047	0.0005	-8.67	<0.0001*

**Table 5-6 Regression model summary for high priority trains on North Dakota segment**

Table 5-7 shows the accuracy and percent improvement over the baseline simple historical average technique of each of these priority specific regression models. In the case of the high

and low priority trains the regression model results in meaningful improvements over the baseline simple average technique. As was observed on the Missouri track segment the greatest improvement in accuracy occurs with the low priority trains.

	Low Priority	Med Priority	High Priority
<b># of Observations</b>	2962	172	4632
<b>Regression Model Accuracy</b>	66.70%	61.63%	76.32%
<b>Simple Average Accuracy</b>	54.12%	55.81%	71.61%
<b>% Improvement</b>	23.25%	10.42%	6.57%

**Table 5-7 Regression model performance for North Dakota segment**

The relatively low number of observations for the medium trains would necessitate further study to confirm the improvement finding. As was observed on the Missouri track segment for high priority trains, the historical run means are effectively equivalent for high and low priority trains on the North Dakota segment which reduces the statistical relevance of the direction variable.

### 5.2.2. Oregon

As Figure 5-4 shows below the top three priorities represent a much smaller percentage of overall traffic relative to the Missouri and North Dakota segments. The run times as shown in Figure 5-5 are relatively stable throughout the period. Directional dependencies are significant for low and medium priority trains in the Oregon segment, Figure 5-6. With the exception of the high priority trains, the mean run times vary significantly by priority and direction as can be observed in Table 5-8.

Priority	Eastbound	Westbound
High	5.67	5.72
Med	6.38	7.29
Low	6.18	8.10
All Other	7.78	8.13

Table 5-8 Mean run times for the Oregon track segment

Figure 5-4 shows the train volume transiting the Oregon track segment and the priority and direction mix over the two year historical data period.

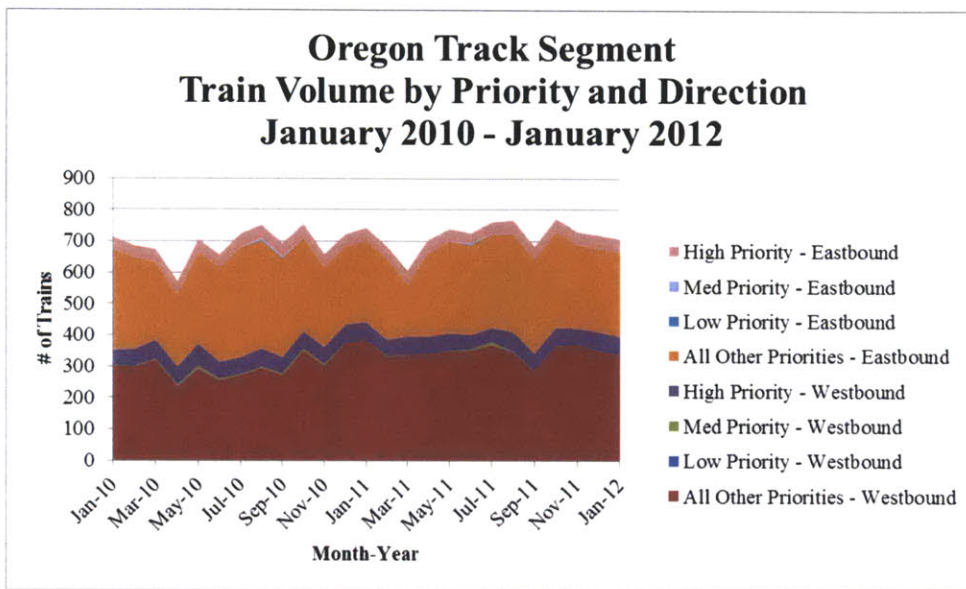


Figure 5-4 Train volume on the Oregon segment from Jan. 2010 – Jan. 2012

Figure 5-5 shows the mean run times by direction and priority for trains transiting the Oregon track segment over the two year historical data period. Due to the limited number of observations for the medium priority category in both directions there are some breaks in the averages during the period.

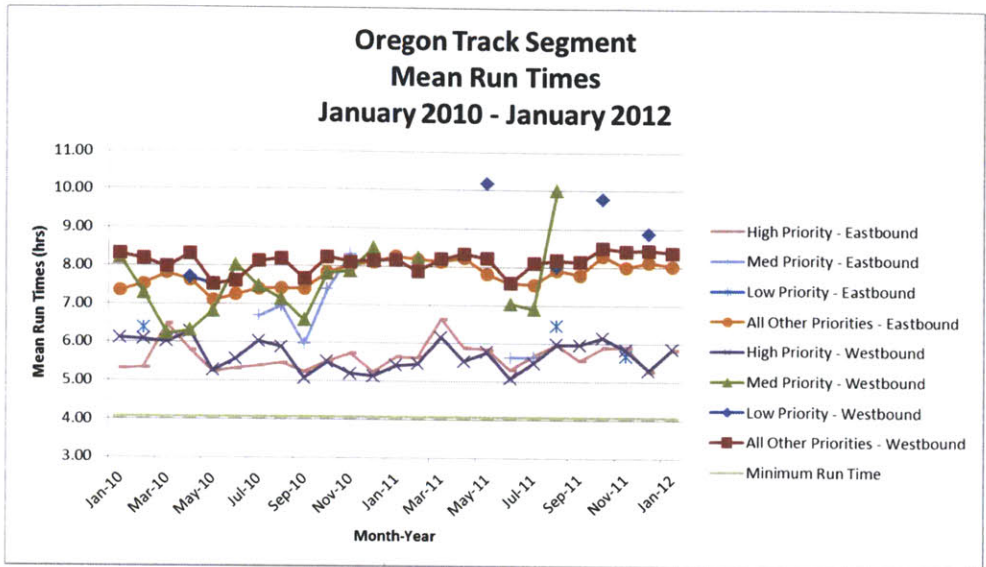


Figure 5-5 Mean run times on the Oregon segment from Jan. 2010 – Jan. 2012

Figure 5-6 below shows the dependency of run times on the direction of travel in the Oregon track segment.

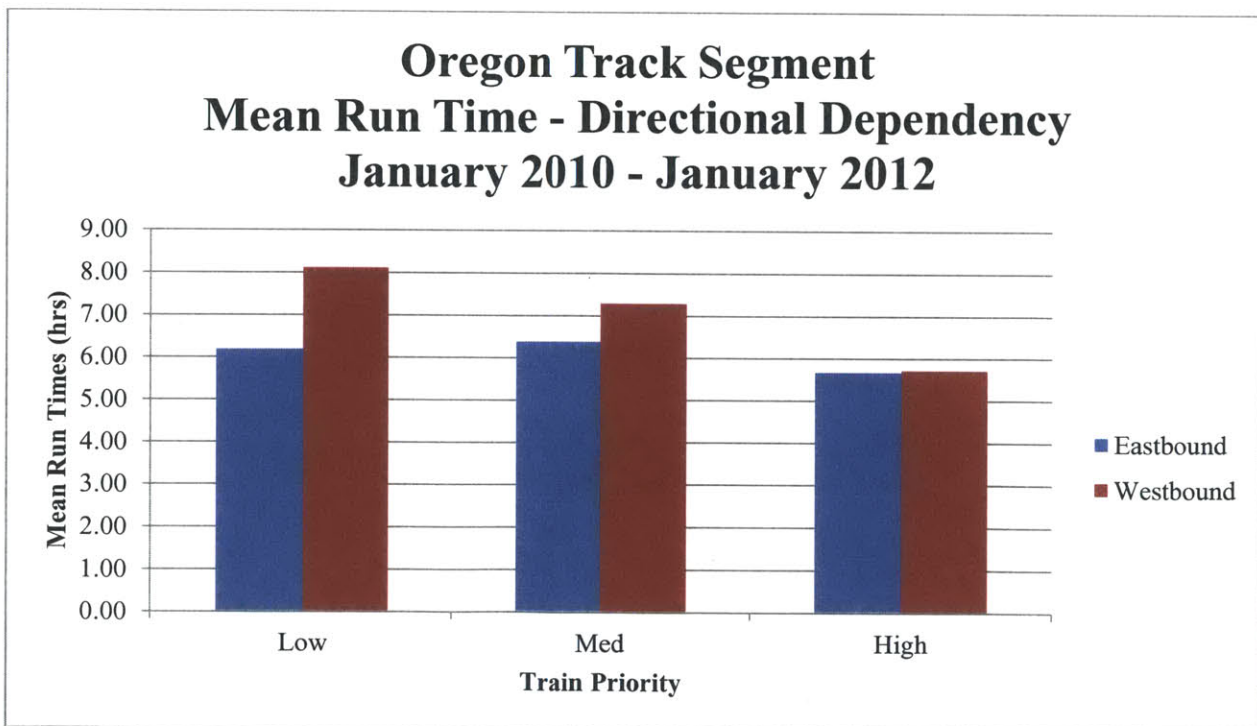


Figure 5-6 Directional dependency of run times on the Oregon segment

Table 5-9 shows the performance of each of these priority specific regression models.

Summary of Fit	Low Priority	Med Priority	High Priority
<b>RSquare</b>	Insufficient Observations	0.46	0.26
<b>RSquare Adj</b>		0.40	0.26
<b>Root Mean Square Error</b>		0.07	0.07
<b>Mean of Response</b>		0.84	0.75
<b>Observations</b>		87	2377

**Table 5-9 Regression model(s) summary of fit for Oregon segment**

Table 5-10 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to medium priority trains operating on the Oregon track segment. Due to the limited number of observations for the medium priority trains further investigation is likely needed to validate the regression model’s performance for this train priority on the Oregon track segment.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.8666	0.1485	5.84	<0.0001*
Rolling Run Time Average	-0.1667	0.1802	-0.92	0.3578
Departure Congestion Window	0.0006	0.0026	0.25	0.8069
Meets	0.0104	0.0044	2.39	0.0194*
Passes	-0.0724	0.0195	-3.72	0.0004*
Overtakes	0.0797	0.0203	3.92	0.0002*
Direction	0.0419	0.0179	2.34	0.0221*
Arrival Headway	0.0219	0.0083	2.64	0.0100*
Departure Headway	-0.0188	0.0071	-2.65	0.0097*

**Table 5-10 Regression model summary for medium priority trains on Oregon segment**

Table 5-11 shows the coefficient, standard error, t-Ratio, and p-value for all explanatory variables in the regression model specific to high priority trains operating on the North Dakota track segment. As was observed on the Missouri and North Dakota routes, meets has the highest statistical significance.

Term	Coefficient	Std. Error	t Ratio	Prob> t
Intercept	0.4763	0.0397	12.01	<0.0001*
Rolling Run Time Average	0.3669	0.0516	7.11	<0.0001*
Departure Congestion Window	-0.0035	0.0004	-8.61	<0.0001*
Meets	0.0127	0.0007	17.85	<0.0001*
Passes	-0.0158	0.0018	-8.98	<0.0001*
Overtakes	0.1383	0.0119	11.57	<0.0001*
Direction	0.0148	0.0029	5.14	<0.0001*
Arrival Headway	0.0056	0.0007	7.99	<0.0001*
Departure Headway	-0.0032	0.0008	-3.95	<0.0001*

**Table 5-11 Regression model summary for high priority trains on Oregon segment**

Table 5-12 shows the accuracy and percent improvement over the baseline simple historical average technique of each of these priority specific regression models. The relatively low number of observations for the medium priority trains would necessitate further study to evaluate the observed accuracy gains. The high priority train model has sufficient observations and mirrors the accuracy and percent improvement over the baseline technique of the model for high priority trains on the North Dakota track segment.

	Low Priority	Med Priority	High Priority
<b># of Observations</b>	Insufficient Observations	87	2377
<b>Regression Model Accuracy</b>		67.82%	77.62%
<b>Simple Average Accuracy</b>		57.47%	71.73%
<b>% Improvement</b>		18.00%	8.21%

**Table 5-12 Regression model performance for Oregon segment**

The regression model when applied to historical data from the Oregon segment shows similar gains in accuracy to that of the Missouri and North Dakota segments for the high priority trains. Due to the lack of observations for the low and medium priority trains further study would be needed to confirm the accuracy of the model when applied to the Oregon segment and these priorities.

### 5.3: Limitations

The regression model utilizes explanatory variables that depend upon accurate operational data. With respect to the historical data the variability is 0%, however with forward-looking operational schedule data there is variability present in the form of trains departing later or earlier than their scheduled departure time or a different train priority departs than the one scheduled.

The partner railroad provided an additional data set consisting of forward looking schedule information for a two month period from February 2012 to April 2012 for the Missouri track segment. This data is herein referred to as the “*Active Schedule Data*”. The basic data element in each set is a specific train ID and scheduled departure date.

The schedule dataset and historical dataset are mutually exclusive with respect to period. As is the case with the historical run data the basic element in the data set is a specific train ID and a scheduled departure data. For each of the discrete train elements and relevant to the analysis, the data contains information specifying estimated crew route (or direction), estimated priority, and an estimated departure time.

Critically an estimated arrival time is not available and must be estimated based on the estimated departure time. The methodology used was add the mean run time from the two year period of historical data specific to the train priority and direction to the estimated departure time to arrive at a calculated estimated arrival time. Further study is likely needed to arrive at an improved arrival time estimation methodology.

The explanatory variables are all calculated based on the schedule data to derive the needed inputs for the regression model. The existing coefficients from the historical run data analysis are used to make run time estimations for the trains present in the active schedule data.

Utilizing the same validation scheme as the historical run data, a baseline run time prediction methodology of using a simple thirty day average, by priority and direction, to estimate run times will represent a base case estimate and stand-in as the existing prediction methodology. A similar comparison of these base case estimates against the actual run times will be conducted with those base case estimates found to be within +/- one hour of the actual run time deemed “accurate”. Table 5-13 shows the accuracy and percent change relative to the baseline simple historical average technique of each of these priority specific regression models.

	Low Priority	Med Priority	High Priority
<b># of Observations</b>	122	230	75
<b>Regression Model Accuracy</b>	47.15%	56.52%	77.33%
<b>Simple Average Accuracy</b>	51.22%	47.83%	92.00%
<b>% Change</b>	-7.94%	18.18%	-15.94%

**Table 5-13 Regression model performance for Missouri segment schedule data**

As shown in Table 5-13, the regression model results in deteriorated performance relative to the baseline simple average technique. Independent of the low number of observations there is a clear failure of the model when utilized against an unreliable schedule. It was found that 89% of the trains scheduled to depart either did not depart within a one hour window (+/- thirty min) of the scheduled departure time or a different train was departed in the scheduled slot. This could result in miscalculated values for congestion, meets, passes, overtakes, arrival headway, and departure headway or making run time predictions for a high priority train when in fact a low priority train departed.

## Chapter 6: Conclusion and Future Direction

### 6.1 Conclusion

The results from the proposed regression models across three geographically distinct regions and three train priorities show a substantial improvement in accuracy when predicting freight train run times over a baseline simple historical averaging technique. The regression model coefficients vary across the low, medium, and high priority trains within a particular track segment which suggests support for the development of priority specific regression models versus a single model for a specific track segment. The model is dependent on having sufficient historical observations in order to correctly model the explanatory variable relationships. Those situations where a low number of observations were collected would need additional study to confirm the applicability of the model to those specific priority and track segments. The proposed approach offers substantial accuracy improvements of 36.79%, 28.74%, 20.95% for the low, medium, and high priority trains respectively over baseline simple historical averaging techniques in a representative single tracked segment.

The explanatory variables selected for use in the regression model represent a broad cross-section of known delay sources and are statistically significant in predicting train run times. Meets represent the highest statistical significance which intuitively mirrors the expected operational dynamics. Attention should be focused on train conflict resolution: meets, passes, and overtakes, as they weigh heavily on train run times. It is clear the model shows strong promise in being able to improve the accuracy of freight train run time predictions. Accurate train scheduling information is needed however to maximize the model's accuracy.

Surprisingly increased train volume in the +/- 12 hour window around a particular train's departure time resulted in decreased run times. The cause of this behavior was not investigated but may suggest operational efficiency improvements to be had as train volume increases.

In most scenarios with respect to predictable delay causing events we find substantial support for the usefulness of meets, passes, and overtakes in estimating train run times.

Due to the substantial improvements of up to 36.79% in run time accuracy that were generated with the new models, we recommend further exploration by the railroad into the proposed solution as warranted as they can ultimately produce more accurate train schedules, with subsequent improvements in railroad capacity, customer service, and asset utilization.

## 6.2 Future Direction

It is important to note that this thesis only discusses run time prediction with respect to transit time. It does not attempt to address: time spent at stations, time spent loading/unloading, and the time requirements of executing a crew change. These factors have significant influence on total train run times, which are a combination of transit and station time.

There is likely some potential for usage of a delay event-based probabilistic model which takes into account the probability of a particular delay event occurring and the estimated quantity of delay from such an event. As noted earlier in this thesis the lack of reliable event-based delay data precludes exploration of this model as a potential solution.

Several of the explanatory variables are dependent on a train's transit window. The scheduled data available from the partner railroad only includes an estimated time of departure necessitating the need to estimate arrival times. It is unclear if the methodology employed in this thesis of utilizing historical mean run times by train priority and train direction to calculate estimated arrival times is the most suitable technique. Further studies should be conducted to

evaluate alternative estimated arrival time's calculation methodologies. The usefulness for applying this model for active run time prediction is constrained by the lack of accurate schedule information. Further studies should be conducted with the partner railroad to identify a more accurate source of information. This research and methodology could likely be utilized to predict transit times in other single degree of freedom transportation systems.

## References

- AAR. (2008). Overview of America's freight railroads. *Association of American Railroads*, (May)
- AAR. (2011a). America needs more rail capacity. *Association of American Railroads*, (October)
- AAR. (2011b). The economic impact of the America's freight railroads. *Association of American Railroads*, (October)
- Abril, M., Barber, F., Ingolotti, L., Salido, M. A., Tormos, P., & Lova, A. (2008). An assessment of railroad capacity. *Transportation Research: Part E*, 44(5), 774-806. doi:10.1016/j.tre.2007.04.001
- Arcot, V.C. (2007). *Modeling uncertainty in rail freight operations: implications for service reliability*. (Master's thesis). Retrieved from the University of Maryland Digital Repository (DRUM)
- Berger, A., & Gebhardt, A. (2011). Stochastic delay prediction in large train networks.
- Carey, M., & Kwiecinski, A. (1994). Stochastic approximation to the effects of headways on knock-on delays of trains. *Transportation Research: Part B*, 28B(4), 251.
- Chen, B., & Harker, P. T. (1990). Two moments estimation of the delay on single-track rail lines with scheduled traffic. *Transportation Science*, 24(4), 261-275.
- D'Ariano, A. (2010). *Improving real-time train dispatching performance: optimization models and algorithms for re-timing, re-ordering and local re-routing*. (Doctoral dissertation). Retrieved from 4OR: A Quarterly Journal of Operations Research, Vol. 8, No 4, DOI: 10.1007/s10288-010-0131-y, 429-432
- D'Ariano, A., Pranzo, M., & Hansen, I.A. (2007). Conflict resolution and train speed coordination for solving real-time timetable perturbations. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 8, No. 2, 208-222.
- Dingler, M., Lai, Y-C., & Barkan, C. (2009). Impact of train type heterogeneity on single-track railroad capacity. *Transportation Research Record: Journal of the Transportation Research Board*, 2117, 41-49.

- Dingler, M., Lai, Y-C., & Barkan, C. (2009). Impact of operational practices on rail line capacity: A simulation analysis. *Proceedings of the American Railway Engineering and Maintenance of Way Association Annual Conference 2009*
- Dingler, M., Koenig, A., Sogin, S., & Barkan, C. (2010). Determining the causes of train delay. *Proceedings of the American Railway Engineering and Maintenance of Way Association Annual Conference 2010*
- Dingler, M. (2010). *The impact of operational strategies and new technologies on railroad capacity*. (Master's thesis). Retrieved from The University of Illinois at Urbana-Champaign Digital Environment for Access to Learning and Scholarship (IDEALS).
- Flier, H., Graffagnino, T., & Nunkesser, M. (2009). Scheduling additional trains on dense corridors. *SEA 2009 Proceedings of the 8<sup>th</sup> International Symposium on Experimental Algorithms*, 149-160
- Greenberg, B. S., Leachman, R. C., & Wolff, R. W. (1988). Predicting dispatching delays on a low speed, single track railroad. *Transportation Science*, 22(1), 31.
- Hallowell, S. F., & Harker, P. T. (1998). Predicting on-time performance in scheduled railroad operations: methodology and application to train scheduling. *Transportation Research, Part A (Policy and Practice)*, 32A(4), 279-95. doi:10.1016/S0965-8564(97)00009-8
- Hallowell, S. F., & Harker, P. T. (1996). Predicting on-time line-haul performance in scheduled railroad operations. *Transportation Science*, 30(4), 364.
- Harrod, S. (2009). Capacity factors of a mixed speed railroad network. *Transportation Research Part E*, 45, 830-841.
- Krueger, H. (1999). Parametric modeling in rail capacity planning. *Proceedings of the 1999 Winter Simulation Conference, Phoenix, AZ*.
- Lai, Y-C., Dingler, M., Hsu, C-E., & Chiang, P-C. (2010). Optimizing train network routing with heterogeneous traffic. *Transportation Research Board of the National Academies*, 2159, 69-76.
- Michael, F. G. (2009). Statistical estimation of railroad congestion delay. *Transportation Research Part E: Logistics and Transportation Review*, 45(3), 446-456. doi:10.1016/j.tre.2008.08.004

- Mu, S., & Dessouky, M. Efficient dispatching rules on double tracks with heterogeneous train traffic.  
(This is a working paper)
- Murali, P., Dessouky, M., Ordóñez, F., & Palmer, K. (2010). A delay estimation technique for single and double-track railroads. *Transportation Research: Part E*, 46(4), 483-495.  
doi:10.1016/j.tre.2009.04.016
- Setar, L. (2011). On track: higher production and spending will bolster demand for rail transportation. *IB*, 48211(Rail Transportation in US)
- Tu, Y., Ball, M. O., & Jank, W. S. (2008). Estimating flight departure delay distributions—A statistical approach with long-term trend and short-term pattern. *Journal of the American Statistical Association*, 103(481), 112-125. doi:10.1198/016214507000000257
- Wendler, E. (2007). The scheduled waiting time on railroad lines. *Transportation Research: Part B*, 41(2), 148-158. doi:10.1016/j.trb.2006.02.009
- Yuan, J. (2006). *Stochastic modelling of train delays and delay propagation in stations*. (Doctoral dissertation). Retrieved from Delft University of Technology Institutional Repository
- Yuan, J., & Hansen, I. A. (2007). Optimizing capacity utilization of stations by estimating knock-on train delays. *Transportation Research: Part B*, 41(2), 202-217. doi:10.1016/j.trb.2006.02.004