

Who Will Take the Light Rail?

Jordan Aron (Biostats)
Kwangho Baek (Civil Engineering)
Minh Nguyen (Data Science)

Structure

1.Introduction

2.Data Overview

3.Methodology

4.Conclusion and Discussion





Introduction

Travel Mode Choice Theory



Which Travel Mode Would you Take?

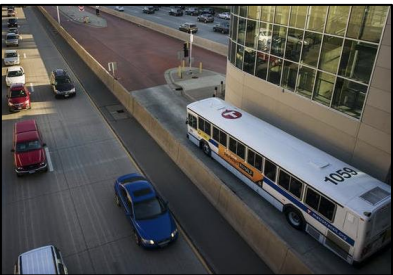


A function of

- (1) socio-demographic attributes
- (2) travel purpose (work/school/...)
- (3) origin/destination location
- (4) path attributes (cost/transit supply)



The most important thing is how a traveler **perceives** each transport mode for a given circumstances-> **latent**



(Metro Transit, Star Tribune)

LRT Ridership Matters: Machine Learning Approach?

\$715.3 M
Blue Line

\$957.0 M
Green Line

$$\frac{\text{Benefit}}{\text{Cost}} > 1?$$

Inaccurate ridership prediction: tremendous social cost



Develop a mode prediction model, unlike the conventional one, which can handle large number of demographic variables that predicts the number and attributes of potential LRT riders.

Who Will Take the Light Rail?



Data

Data Source

2016 Transit On-Board Survey (OBS 2016) from Metro Transit, MN

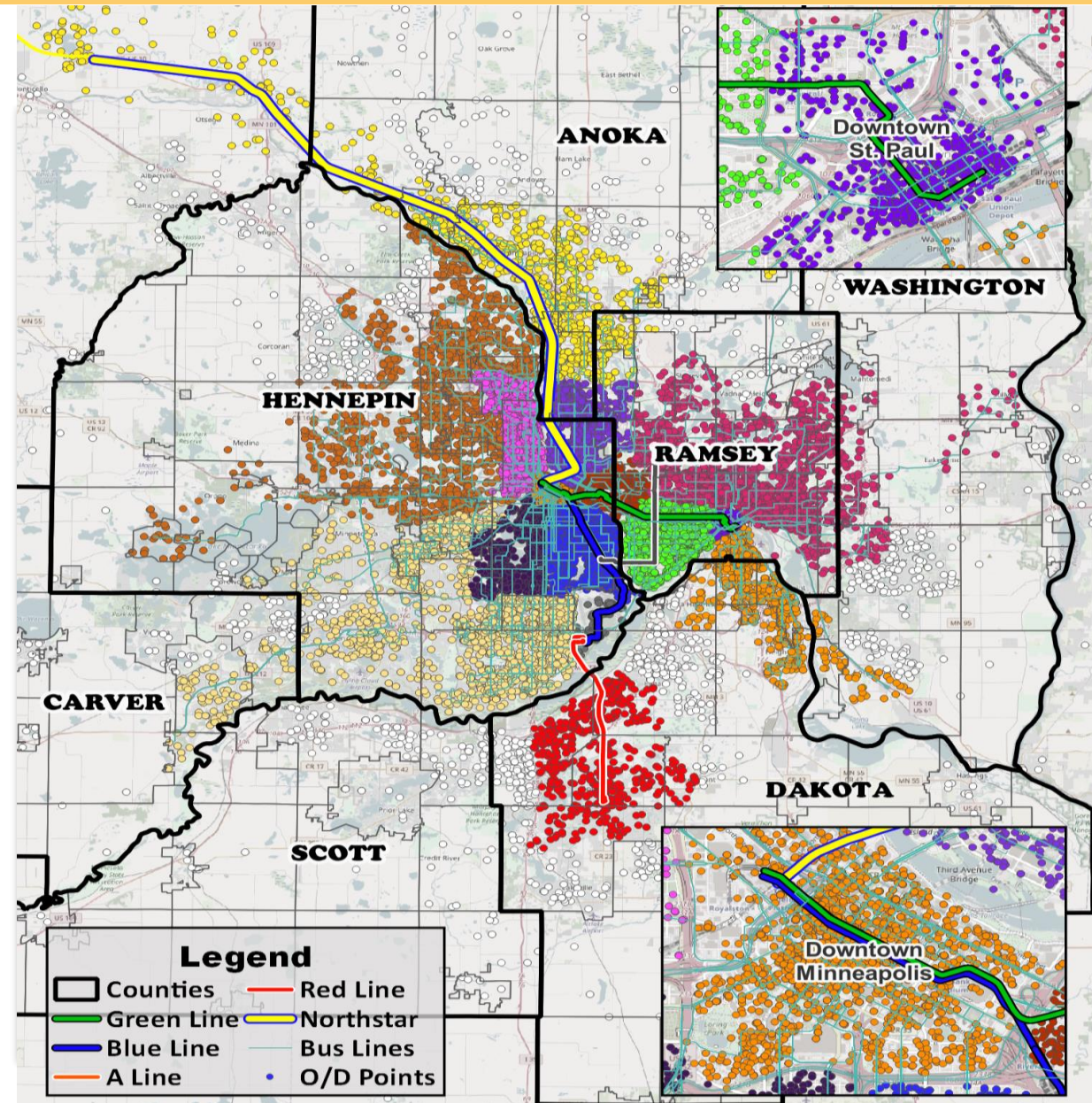
30,491
Observations
(Passenger
Count)

174
Demographic
& Travel-related
Attributes

Weights
Available for
adjusting
sampling error

- Travel Start Time
- Gender
- Age Group
- Access Mode
- Egress Mode
- Payment Method
- Transfers
- Income
- Employment or studentship
- Disability
- Ethnicity
- Home Language
- Fare Subsidized?
- Car Availability
- OD Coordinates
- OD Place Type
- ...

OD: Origin and Destination



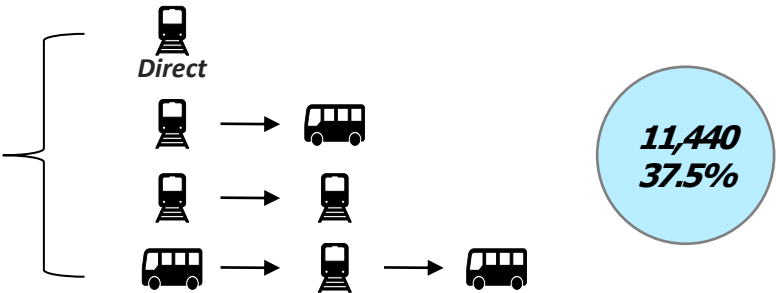
Binary Target Class Label

The Label

N=30,491

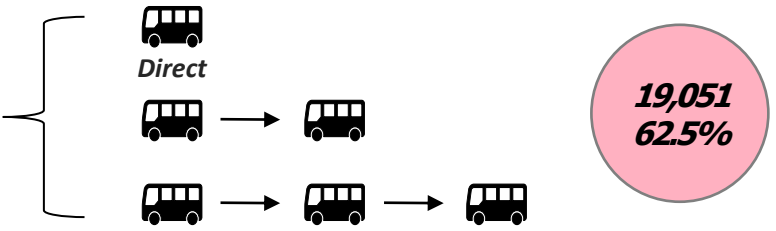
:LRT (Light Rail Transit) :Bus → :Transfer

isLRT=1



If a passenger utilized at least one LRT route at any time during their travel

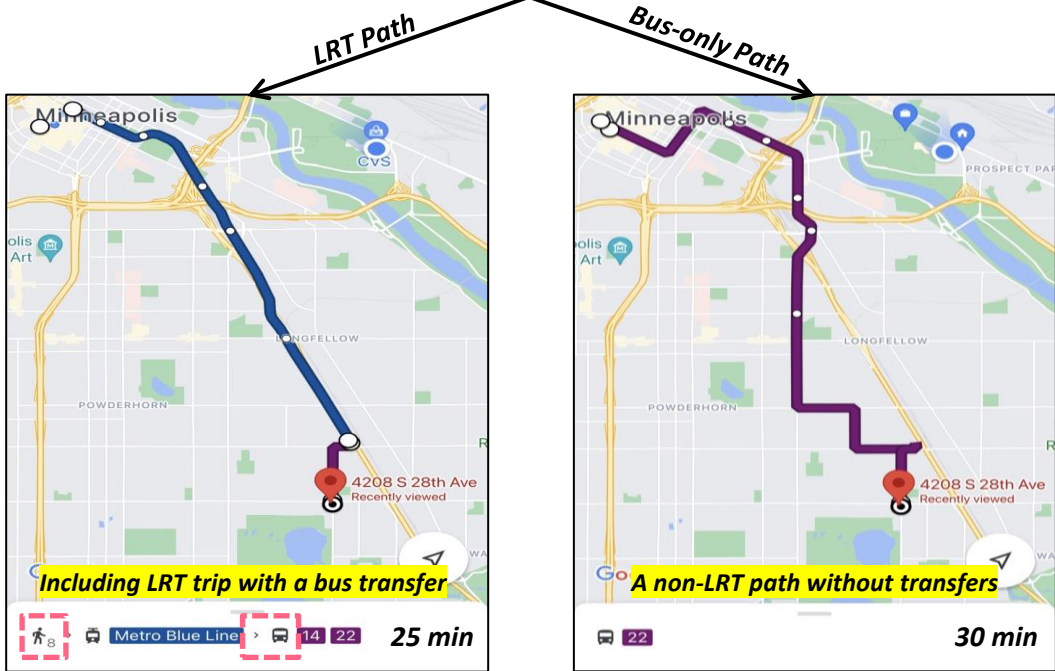
isLRT=0



If a passenger only utilized bus route(s) during their travel

*The number and order of transfers of each path do not affect the classification

Example



isLRT=1

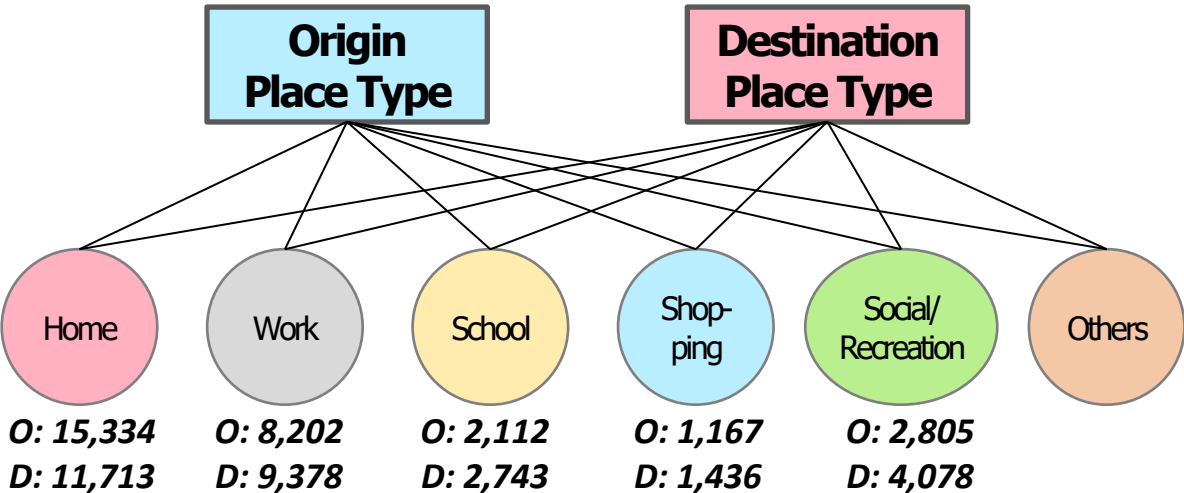
isLRT=0



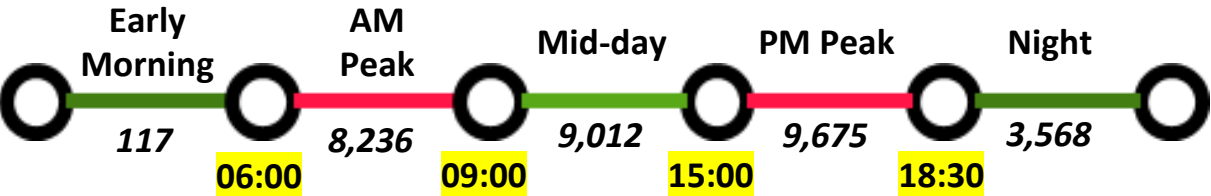
Predictors

Place Type (6 levels)

N=30,491

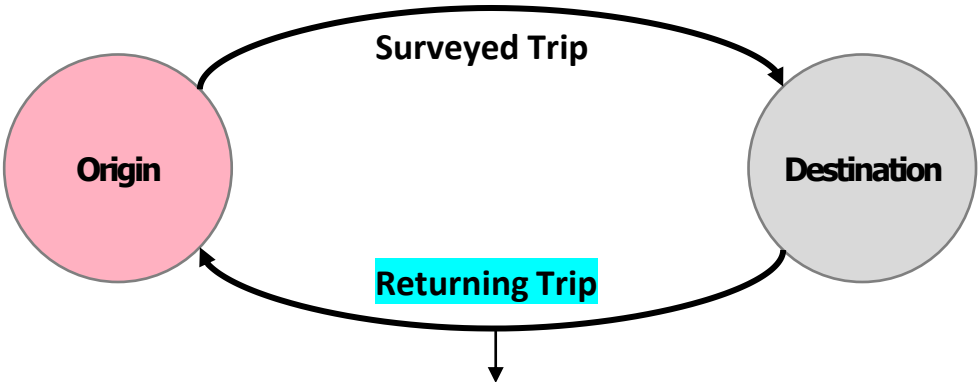


Time Period (5 levels)



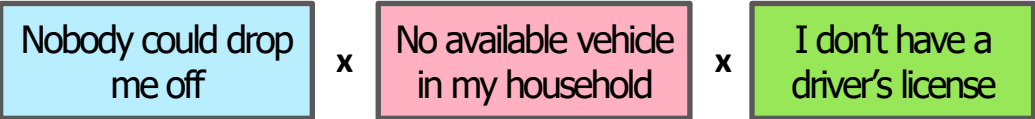
Opposite trip planned? (binary)

N=30,491



If a respondent said this trip was planned: 1 (19,879), otherwise: 0 (10,612)

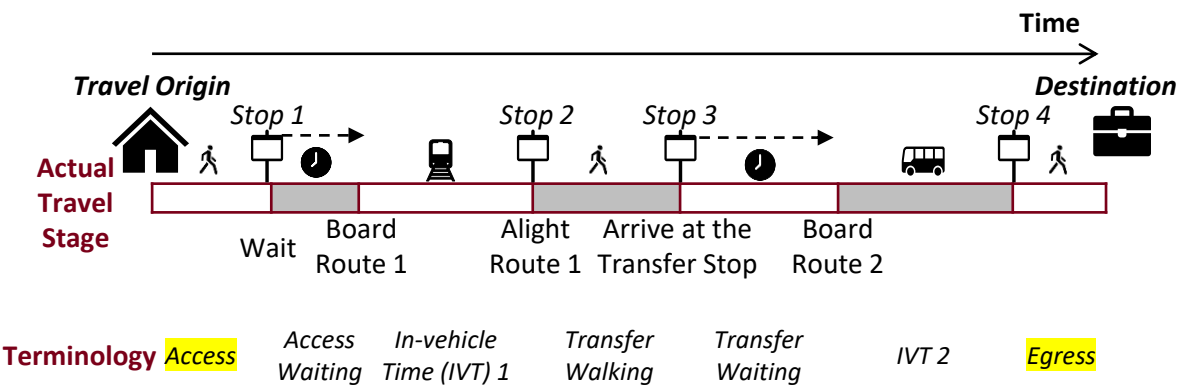
Vehicle Availability (binary)



If the joint condition is true: 0 (22,014), otherwise: 1 (13,747)

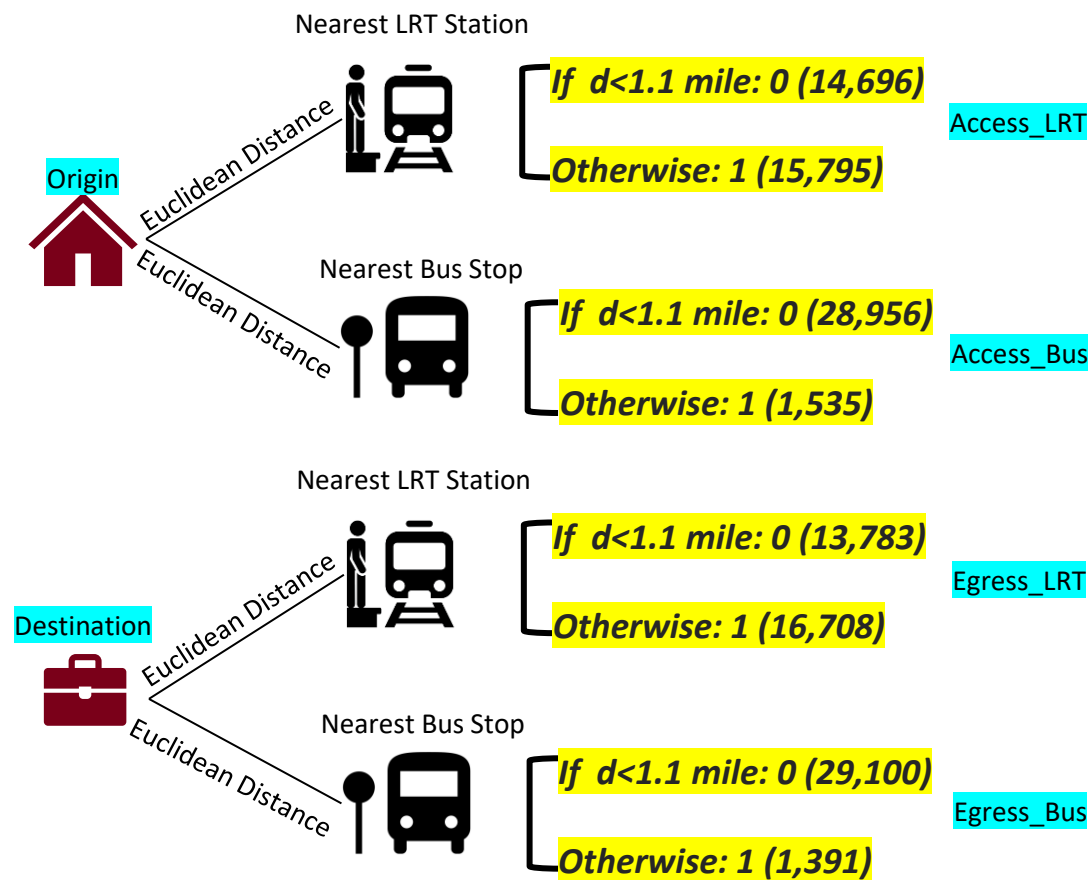
Predictors

Motorized Access/Egress Mode and Transfer (binary) N=30,491



If access/egress modes are “on foot”, “bike”, “wheelchair” :0, otherwise 1

Distances (four binary variables) N=30,491

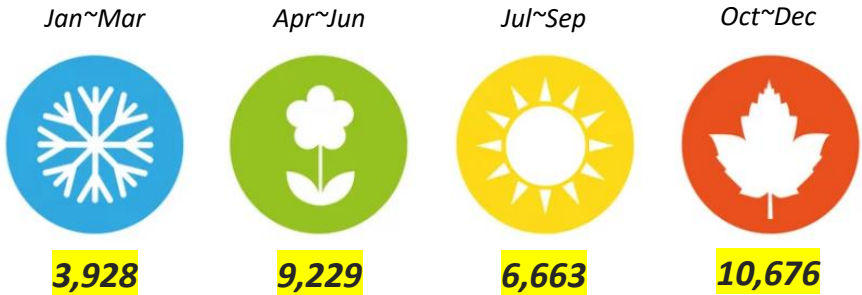
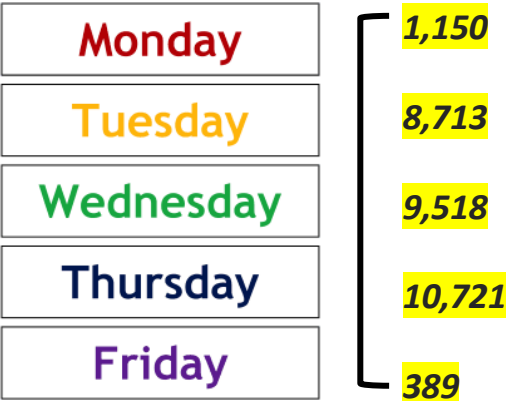


1.1 mile: 95% cutoff value of transit traveler's walking access distance

Predictors

Day/Season

N=30,491



Others (Self-Explanatory)

N=30,491

Category	Level	Count
Gender	Female	17,088
	Male	18,673
Age Group	Age over 44	10,889
	Age 25-44	15,588
	Age under 25	9,284
Ethnicity	White	22,079
	Non-white	13,682
Annual Income (in 2016 USD)	Over \$60K	16,814
	\$15K~\$60K	14,381
	Under \$15K	4,566
Student/ Employment Status	Student	10,364
	Employed	22,402
	None of the above	2,995
Disability	Disabled	3,837
	Not Disabled	31,924

Who Will Take the Light Rail?



First 10 Rows of the Data

	ORIGIN_PLACE_TYPE	DESTIN_PLACE_TYPE	TIME_PERIOD	TRIP_IN_OPPOSITE_DIR
1	Social Visit / Community / Religious / Personal	Social Visit / Community / Religious / Personal	Midday	No
2	Your HOME	Recreation / Sightseeing / Restaurant	Midday	No
3	Work	Your HOME	Midday	Yes
4	Your HOME	Recreation / Sightseeing / Restaurant	PM Peak	Yes
5	Doctor / Clinic / Hospital (non-work)	Social Visit / Community / Religious / Personal	PM Peak	No
6	Work	Shopping	Midday	Yes
7	Your HOME	Work	Midday	Yes
8	Work	Your HOME	PM Peak	No
9	Social Visit / Community / Religious / Personal	Your HOME	PM Peak	No
10	Doctor / Clinic / Hospital (non-work)	Social Visit / Community / Religious / Personal	PM Peak	Yes

	GENDER	Access	Egress	Vehicle	pay	income	status	disable	ages
1	Male	Non-Motorized	Non-Motorized	Not Available	Free Ride or Pass Used	\$15K~\$60K	Others	Disabled	Over 44
2	Male	Non-Motorized	Non-Motorized	Not Available	Cash	Under \$15K	Others	Disabled	Over 44
3	Female	Non-Motorized	Non-Motorized	Not Available	Go-to Card or Mobile	Under \$15K	Employed	Disabled	Under 25
4	Male	Motorized	Non-Motorized	Available	Go-to Card or Mobile	\$15K~\$60K	Employed	Disabled	Over 44
5	Female	Non-Motorized	Non-Motorized	Available	Free Ride or Pass Used	\$15K~\$60K	Student	Disabled	Age 25-44
6	Female	Non-Motorized	Non-Motorized	Not Available	Cash	Under \$15K	Student	Disabled	Age 25-44
7	Female	Non-Motorized	Non-Motorized	Available	Free Ride or Pass Used	\$15K~\$60K	Others	Disabled	Under 25
8	Female	Non-Motorized	Non-Motorized	Not Available	Free Ride or Pass Used	Under \$15K	Employed	Disabled	Age 25-44
9	Male	Motorized	Non-Motorized	Available	Go-to Card or Mobile	Under \$15K	Others	Disabled	Age 25-44
10	Female	Non-Motorized	Non-Motorized	Not Available	Go-to Card or Mobile	Under \$15K	Others	Disabled	Over 44

	eng	HH	iswhite	Label	DAY	SEASON	Orig_Rail_Meter_Disc	Dest_Rail_Meter_Disc	Orig_Bus_Meter_Disc	Dest_Bus_Meter_Disc
1	Others	1	Yes	LRT	Monday	Spring	1	0	0	0
2	English	1	No	LRT	Monday	Spring	0	1	0	1
3	English	More than 2	No	Bus	Monday	Spring	0	1	0	0
4	English	1	Yes	LRT	Monday	Spring	1	0	1	0
5	English	1	Yes	LRT	Monday	Spring	0	0	0	0
6	English	2	Yes	LRT	Tuesday	Spring	0	0	0	0
7	Others	More than 2	No	LRT	Tuesday	Spring	0	0	0	0
8	English	1	Yes	LRT	Tuesday	Spring	1	0	0	0
9	English	1	No	LRT	Tuesday	Spring	0	0	0	0
10	English	1	Yes	LRT	Tuesday	Spring	0	1	0	0

P=22 * N=30,491

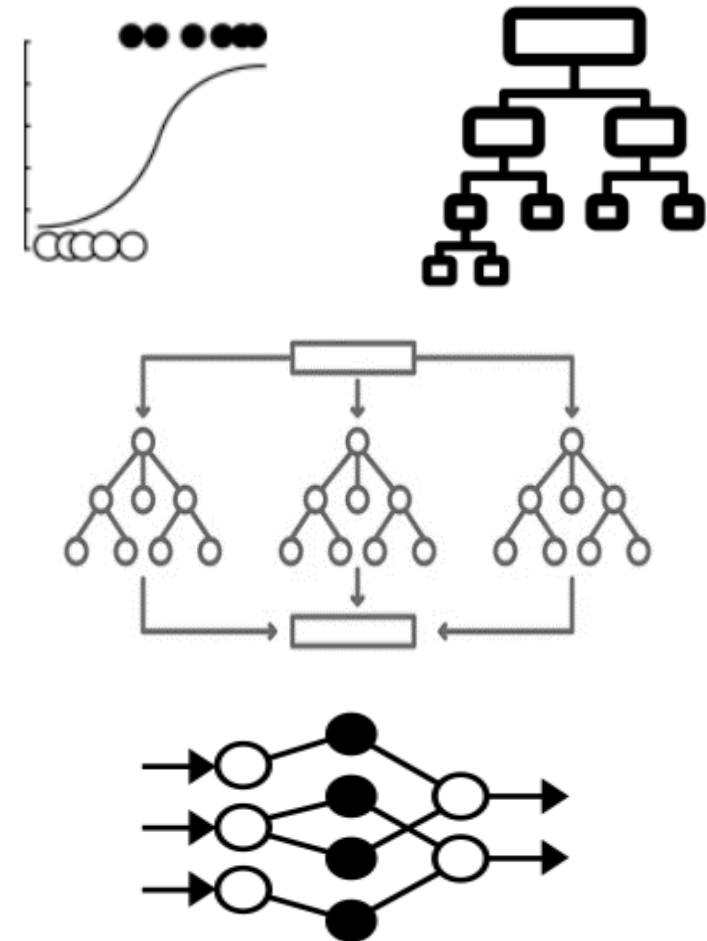
Methods Overview

1. Implemented models:

- Logistic regression
- Decision trees
- Random forest
- Neural network

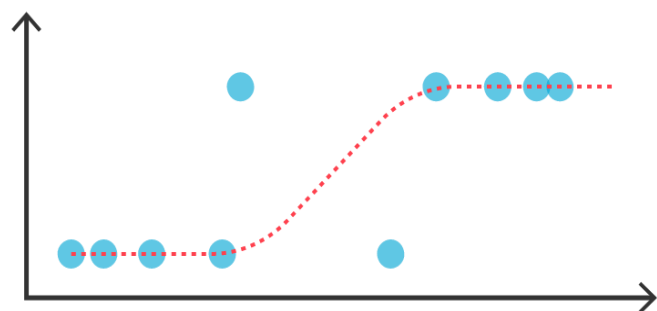
2. Training and Testing Process

- Original dataset: 30,491 observations
- 80:20 training testing data split (24,393 vs 6,098)



Logistic Regression

Logistic Regression Model



1. Used **all** (22) available predictors
2. **Test error rate: 0.249**
3. Simple to implement
4. Not doing so well in categorizing LRT users
5. Baseline model

Results

Accuracy : 0.7506
95% CI : (0.7395, 0.7614)
No Information Rate : 0.6378
P-Value [Acc > NIR] : < 2.2e-16

Sensitivity : 0.5378
Specificity : 0.8714
Pos Pred Value : 0.7038
Neg Pred Value : 0.7685
Prevalence : 0.3622
Detection Rate : 0.1948
Detection Prevalence : 0.2768
Balanced Accuracy : 0.7046

Kappa : 0.4312
McNemar's Test P-Value : < 2.2e-16

6,098 Tests	Pred Bus	Pred LRT
True Bus	3,389 (87%)	500 (13%)
True LRT	1,021 (46%)	1,188 (54%)

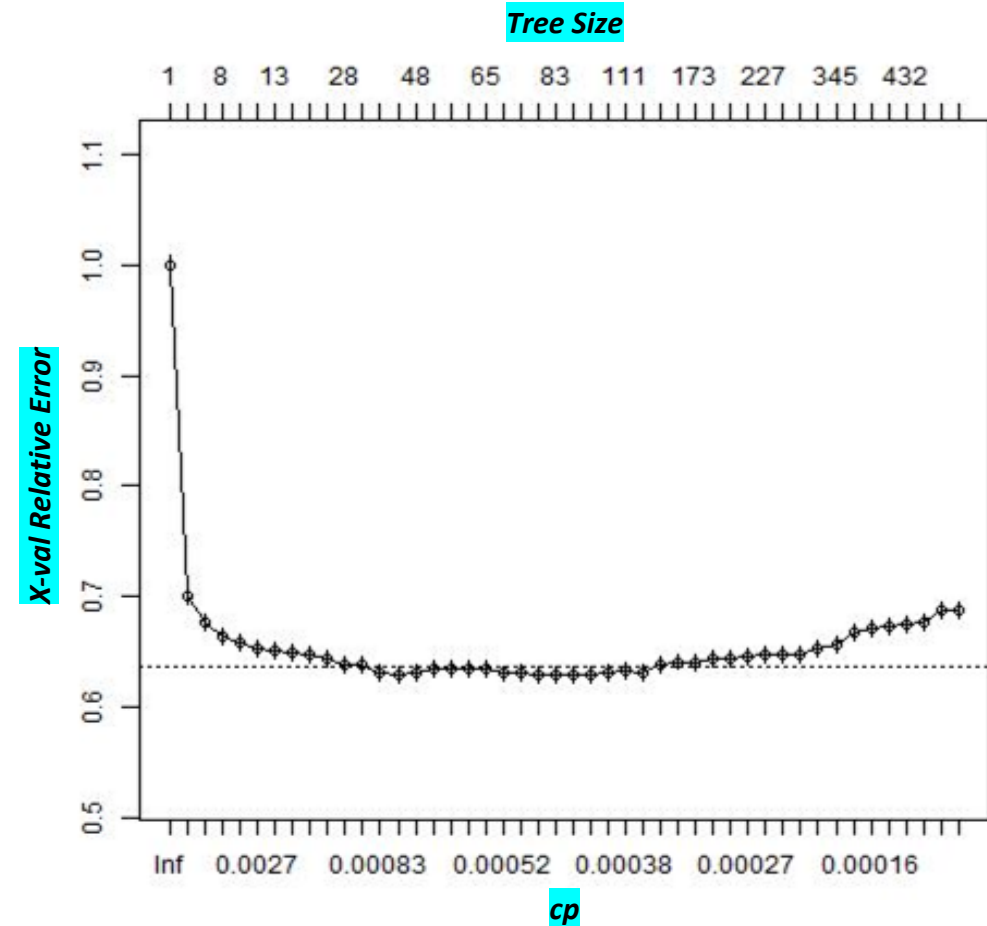
Decision Tree (Baseline Tree)

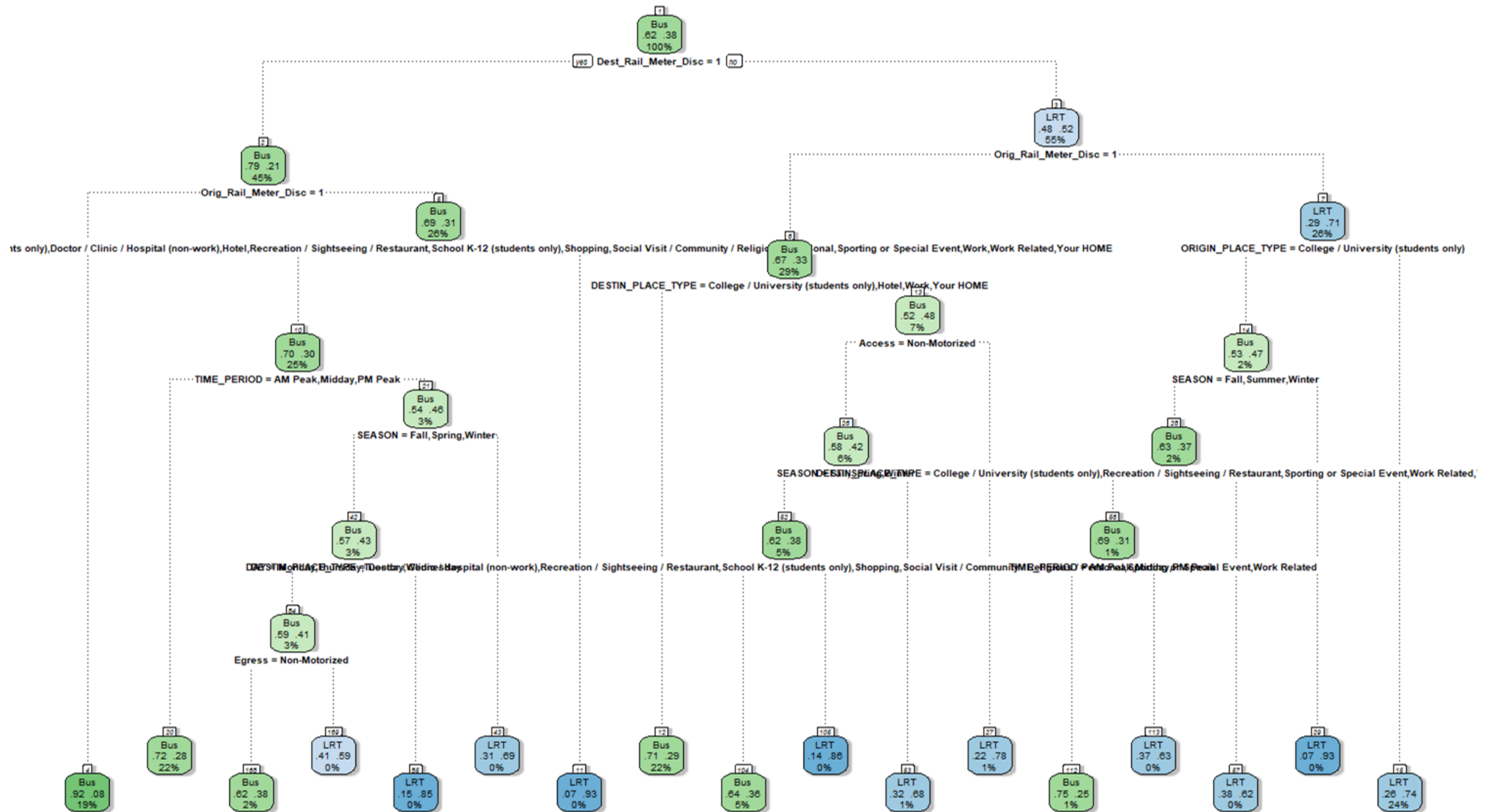
■ Baseline Decision Tree

1. Baseline tree was grown using all available predictors
2. 10 fold CV was implemented for tuning parameters
3. **Test error rate: 0.255**

6,098 Tests	Pred Bus	Pred LRT
True Bus	3,190 (82%)	699 (18%)
True LRT	850 (39%)	1,350 (61%)

■ Relative Error





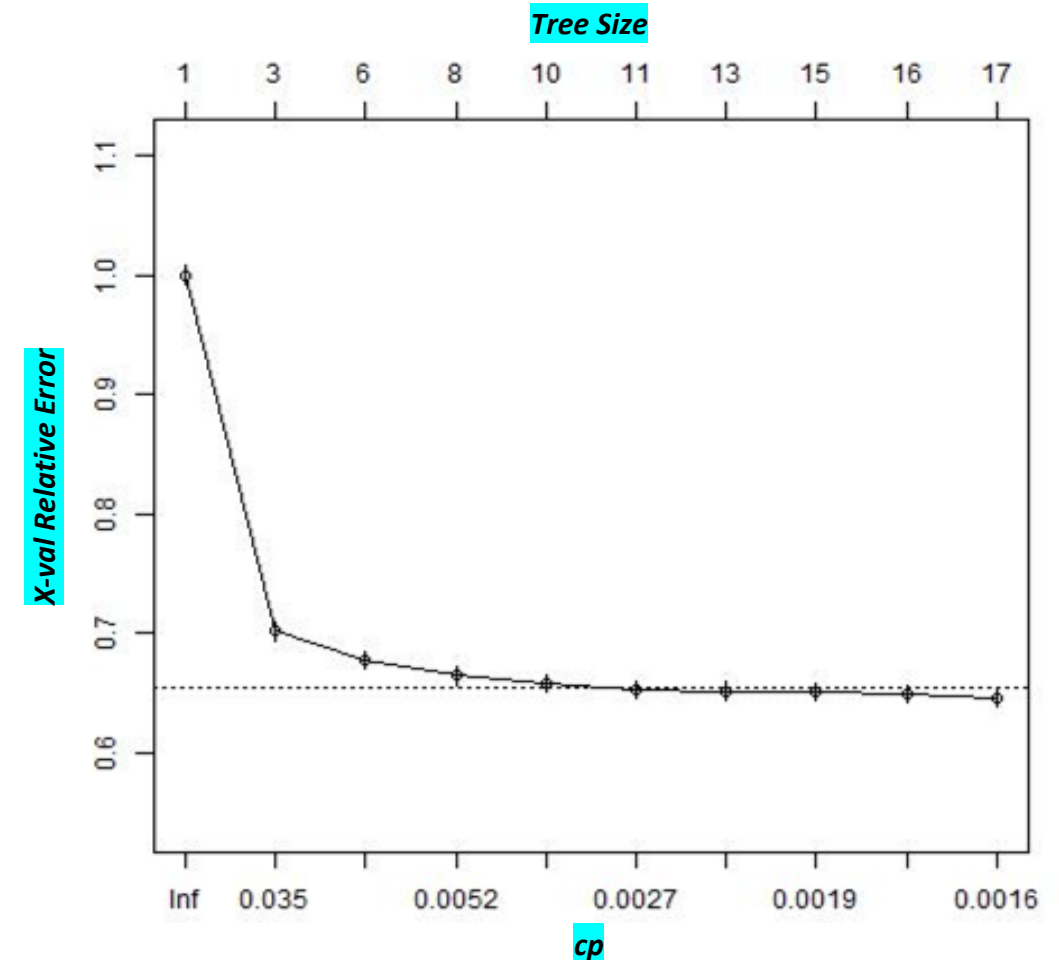
Pruned Tree

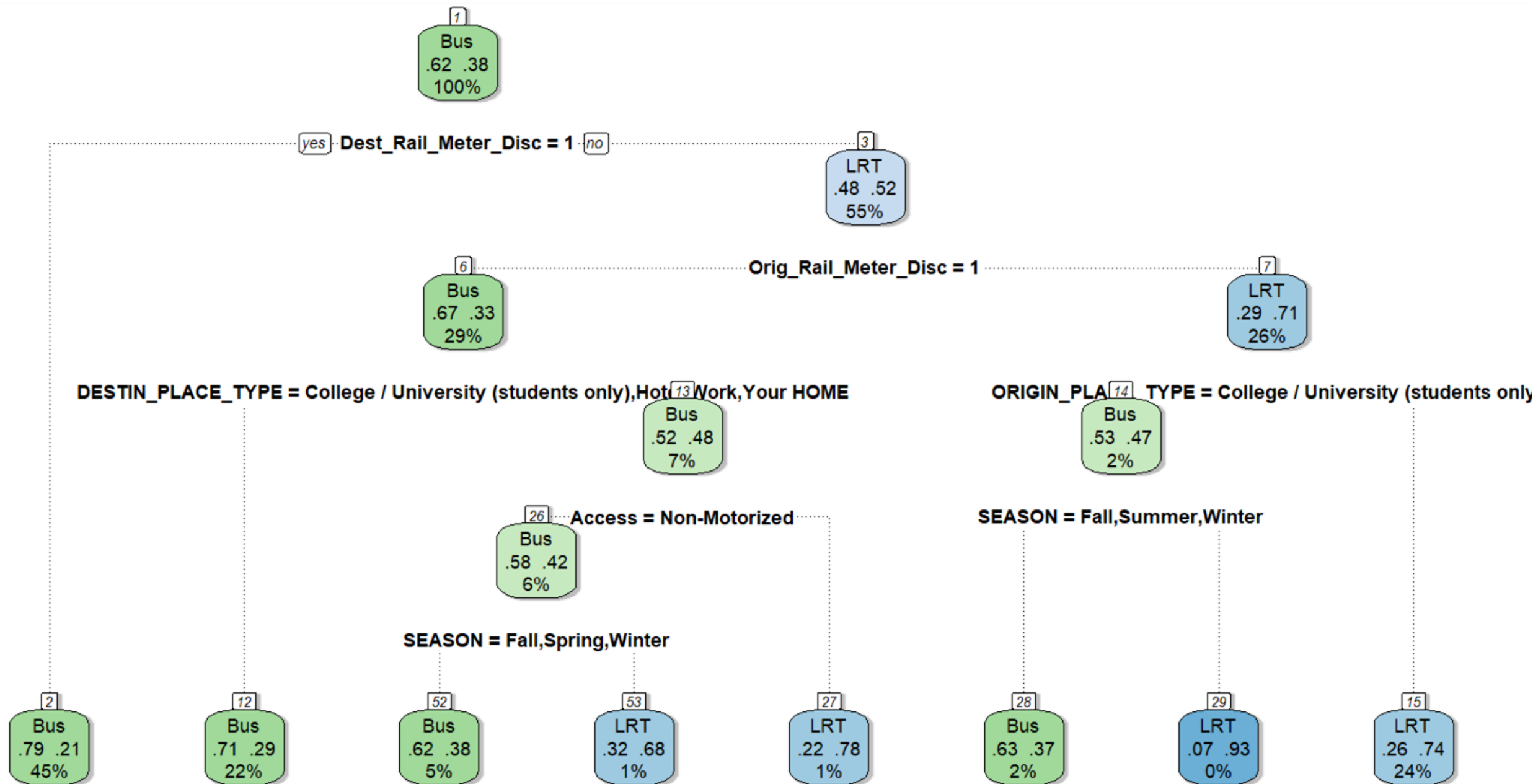
■ Pruning Result

1. Pruning resulted in 9 predictors to be utilized
2. 10 fold CV was implemented for tuning parameters
3. **Test error rate: 0.240**

6,098 Tests	Pred Bus	Pred LRT
True Bus	3,430 (88%)	459 (12%)
True LRT	1,005 (45%)	1,204 (55%)

■ Relative Error





Most Pruned Tree

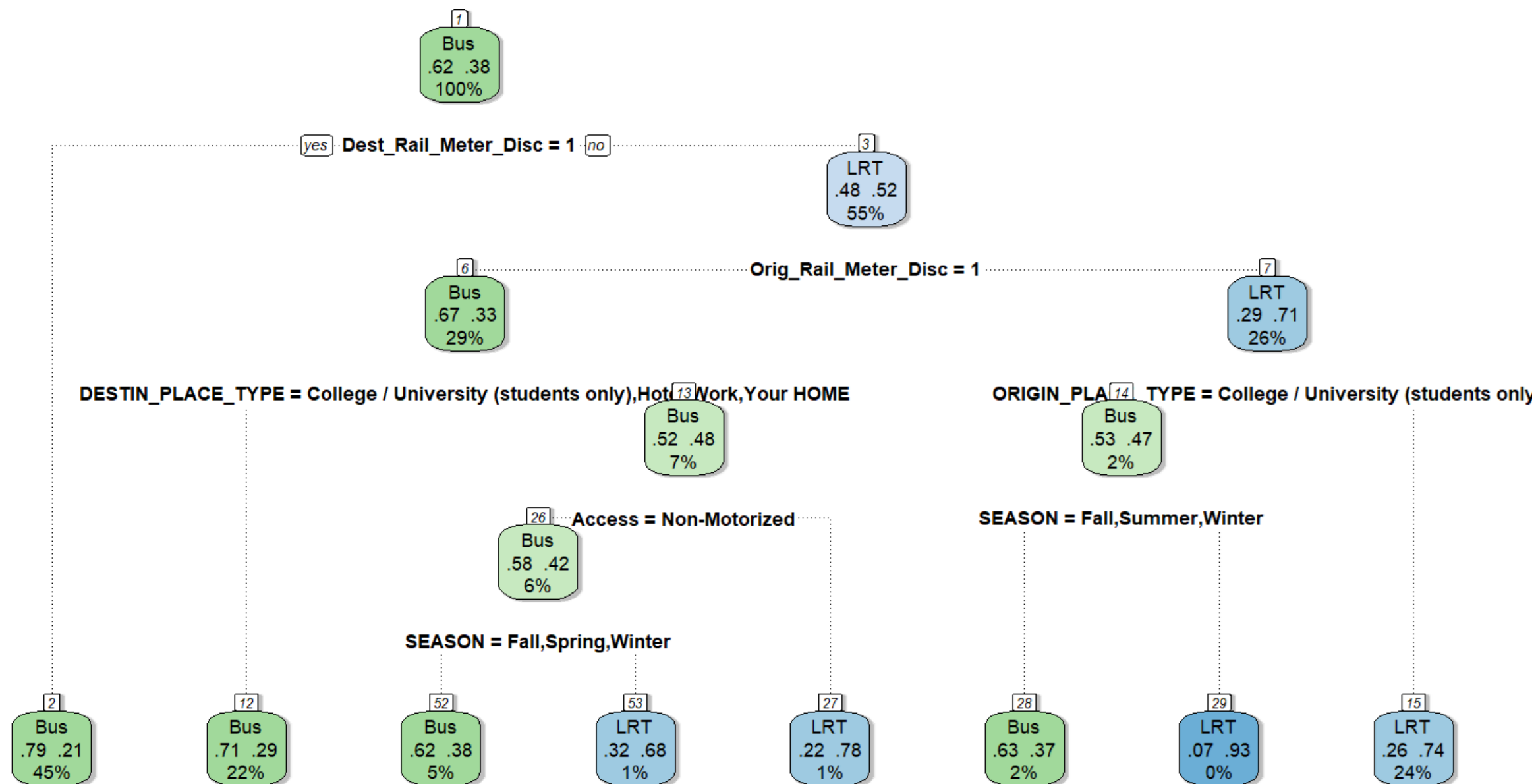
1. Pruning resulted in 6 predictors to be utilized

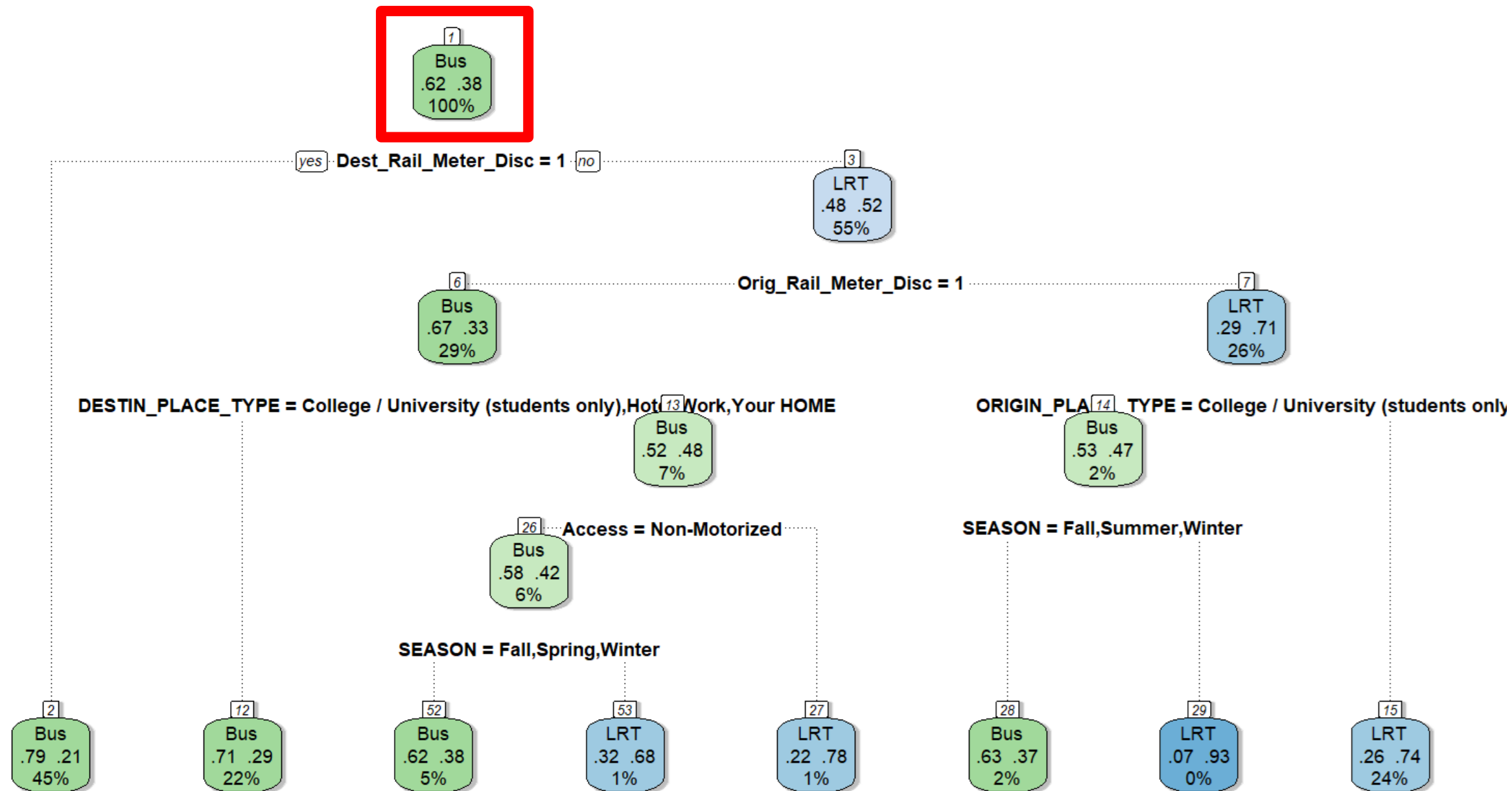
- Discretized Distance from Origin to nearest LRT station
- Discretized Distance from Destination to nearest LRT station
- Origin Place Type
- Destination Place Type
- Access Mode (Motorized or not)
- Season

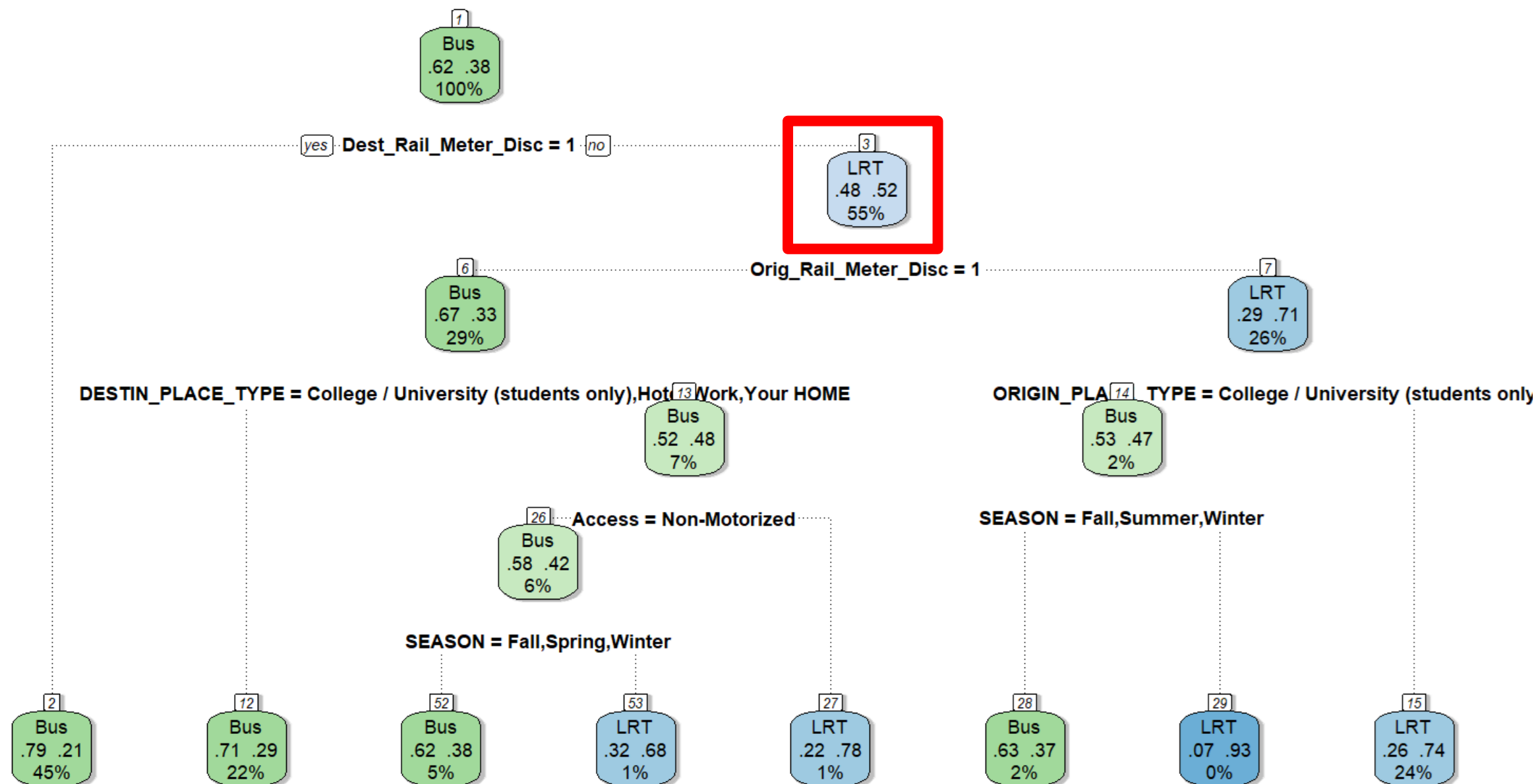
2. Most straightforward interpretation

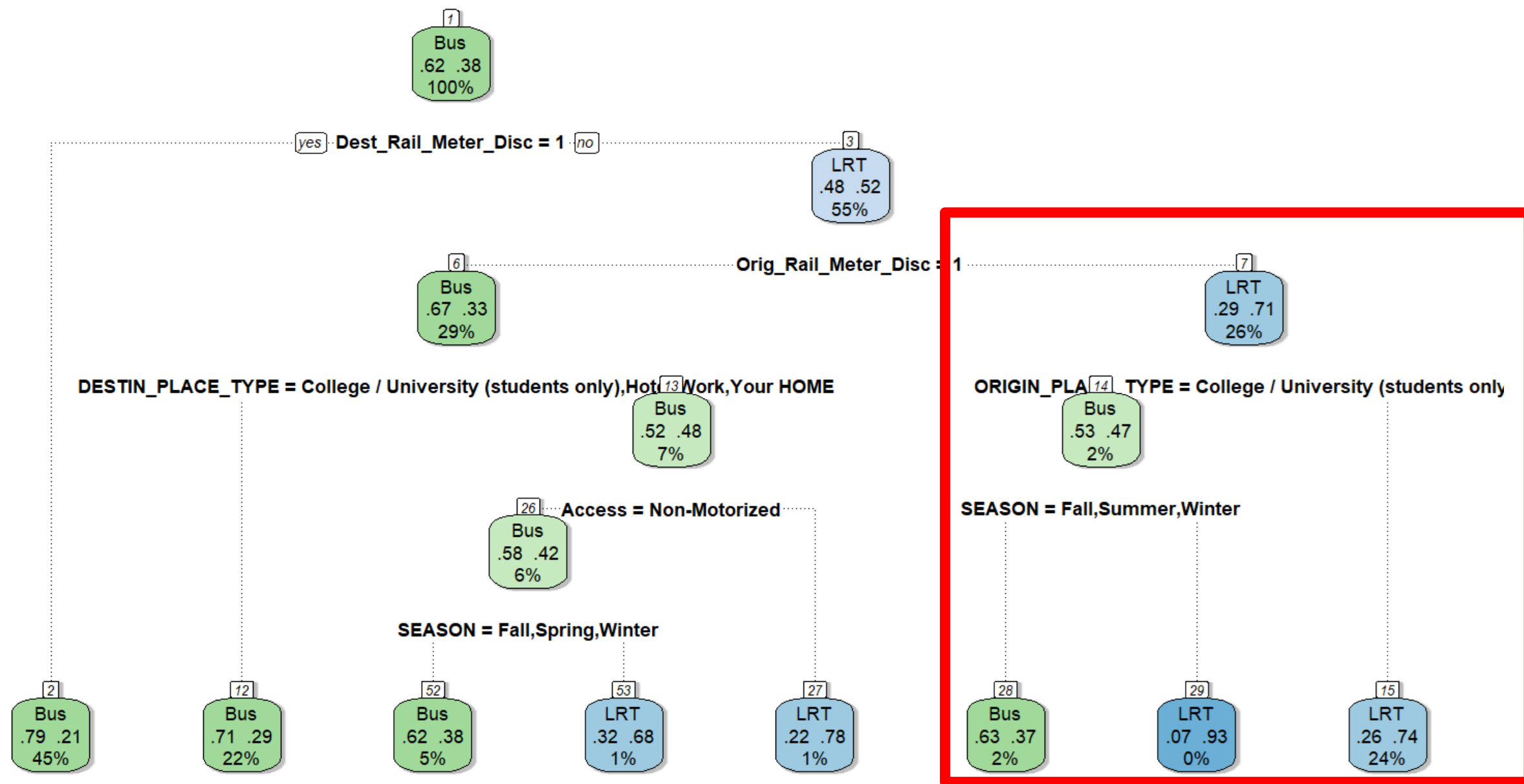
3. Test error rate: 0.249

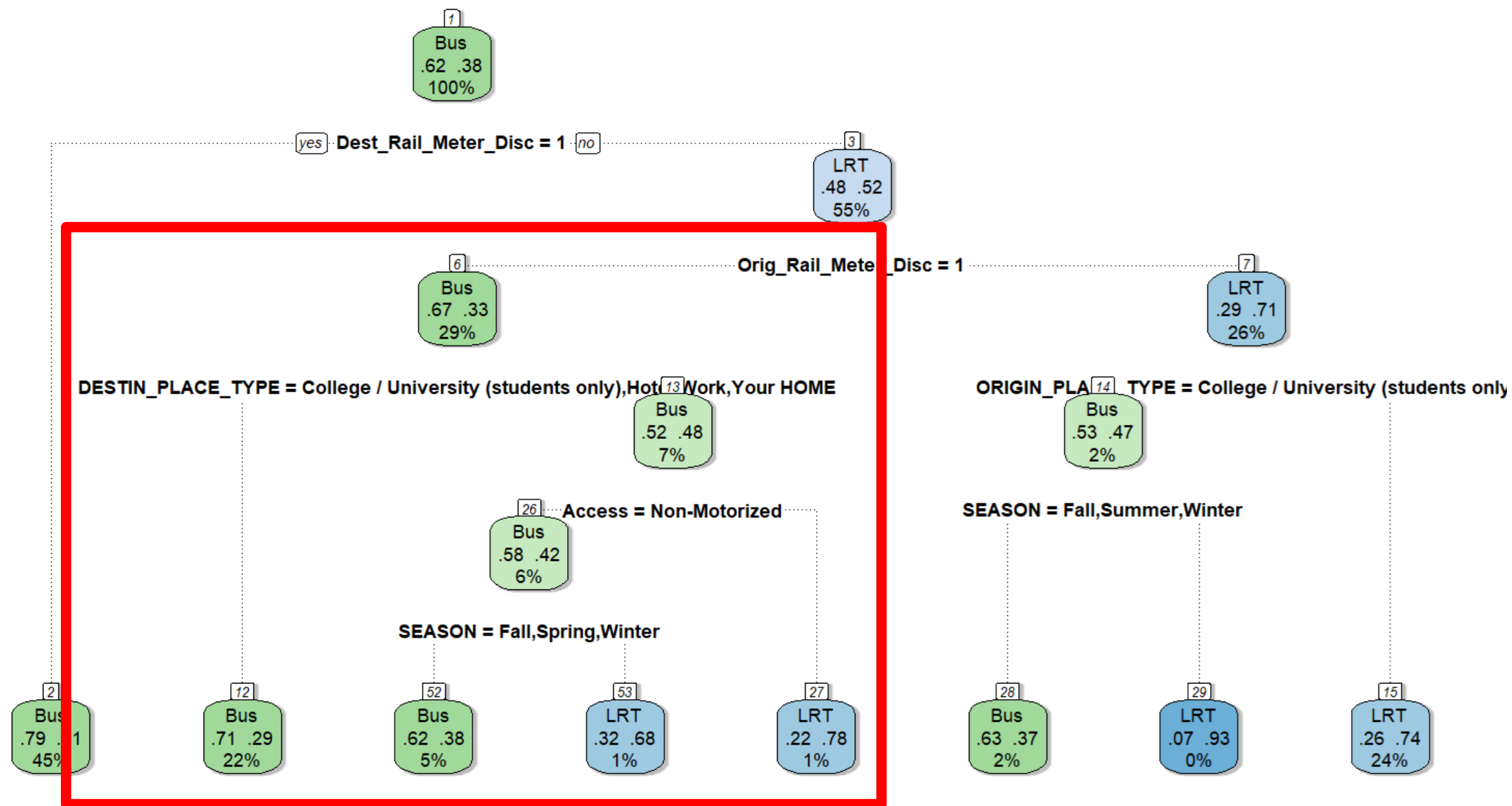
<i>6,098 Tests</i>	Pred Bus	Pred LRT
True Bus	3,462 (89%)	427 (11%)
True LRT	1,094 (50%)	1,115 (50%)











Most Pruned Tree

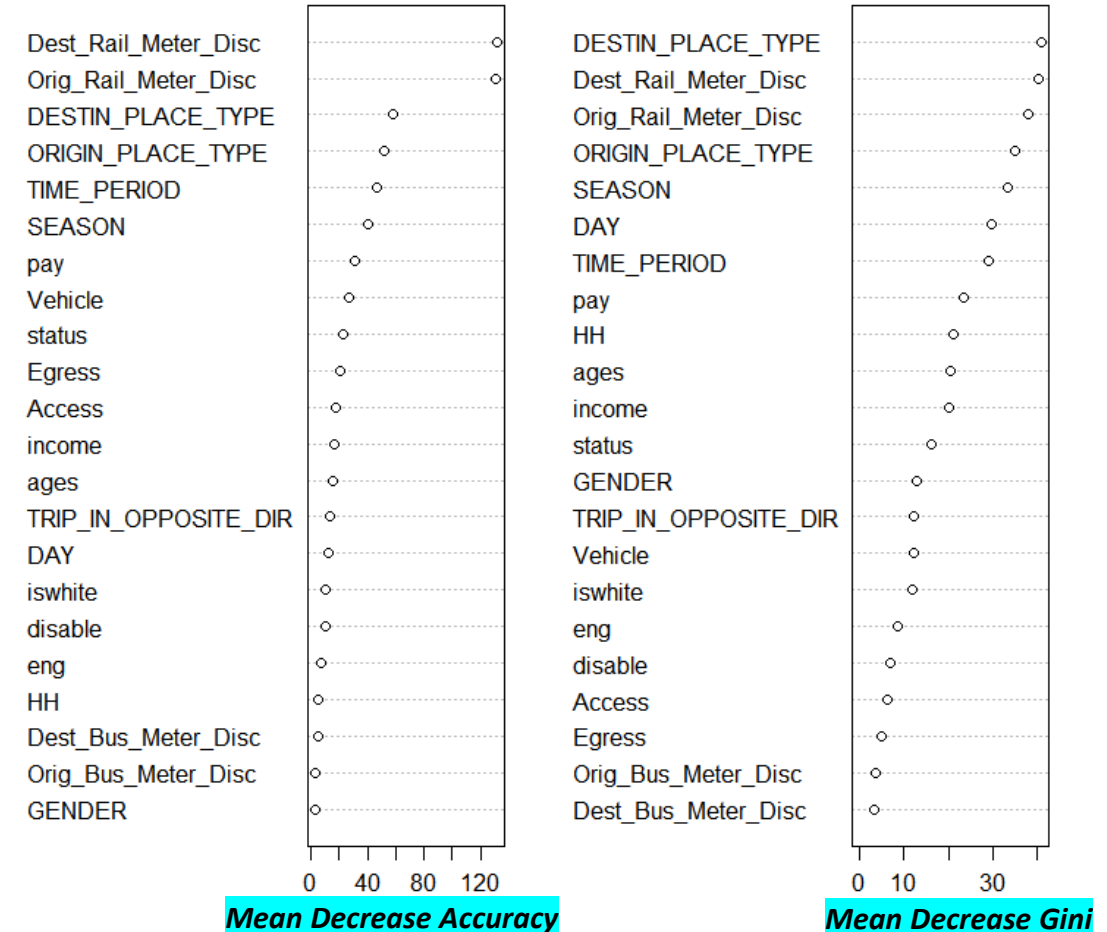
Random Forest

Model Fitting

1. All predictors were used
2. **Test error rate: 0.236**
3. Simple to implement, but higher time complexity
4. Slightly overfitting, but acceptable
5. Aligning features importance with DT
6. Does not doing so well with categorizing LRT

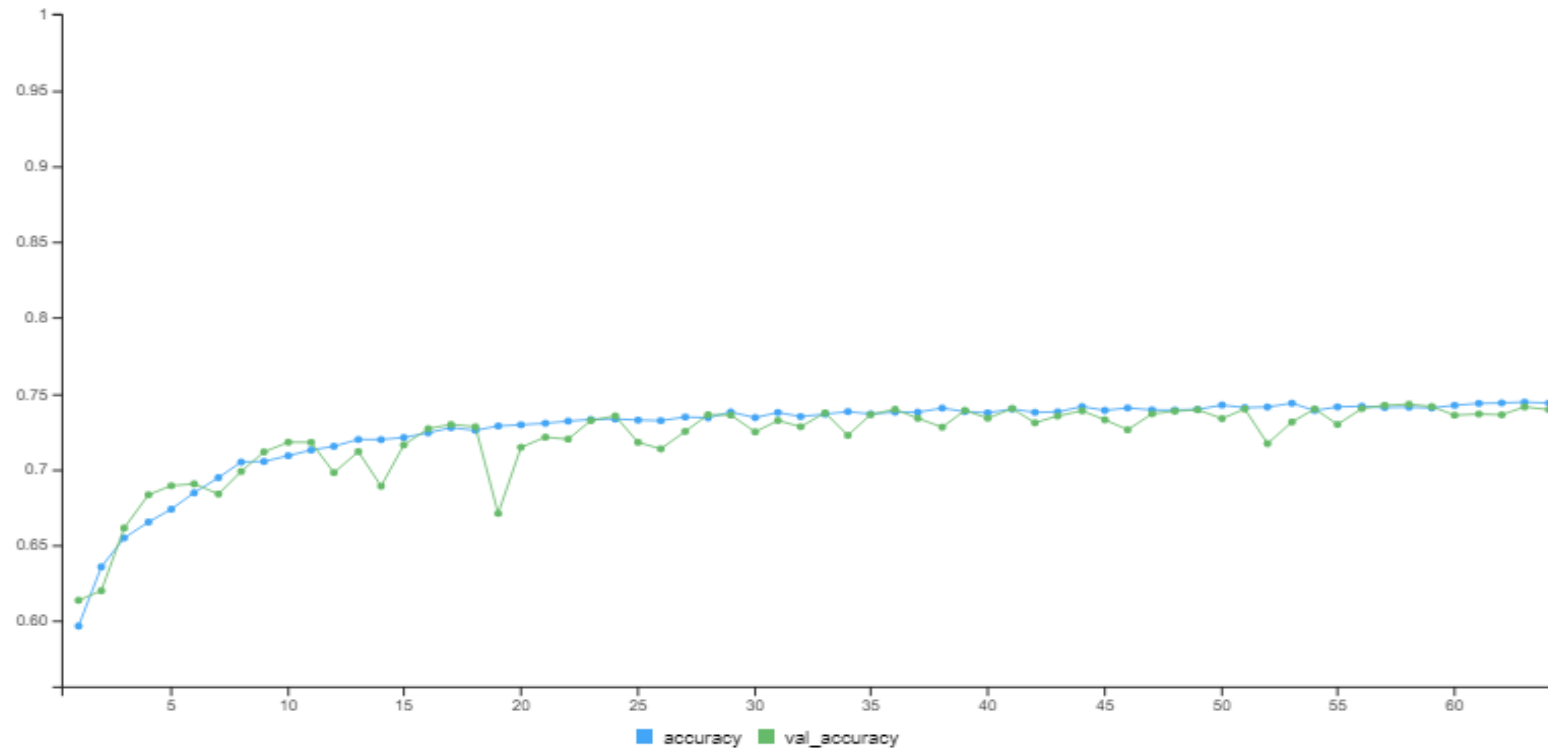
6,098 Tests	Pred Bus	Pred LRT
True Bus	3,480 (89%)	409 (11%)
True LRT	1,034 (47%)	1,175 (53%)

Importance Plots



Neural Network

1. FNN with between 1-5 layers and 32-256 nodes per layer and different activation functions
2. All results were similar, with the lowest test error of 26.0% with 2 layers of 32 nodes



6,098 Tests	Pred Bus	Pred LRT
True Bus	3,375 (87%)	514 (13%)
True LRT	1,050 (47%)	1,159 (53%)

Conclusion

Discussion

1. Smallest test error rate of **23.6% (Random Forest)**
2. **Behavioral interpretation** with decision trees and RF's predictor importance plot are available for better interpretability
3. The main purpose of the models are predicting & analyzing LRT for given demographic/geographic input -> **The accuracy for LRT** could be another metric (False LRT is costlier than False Bus)
4. **Baseline Decision Tree** had the most accurate LRT prediction
5. Incorporating **path attributes** (in-vehicle time, # of transfers...) predictors derivable from shortest path passenger assignment algorithms **would significantly increase** the model accuracy

Results Summary

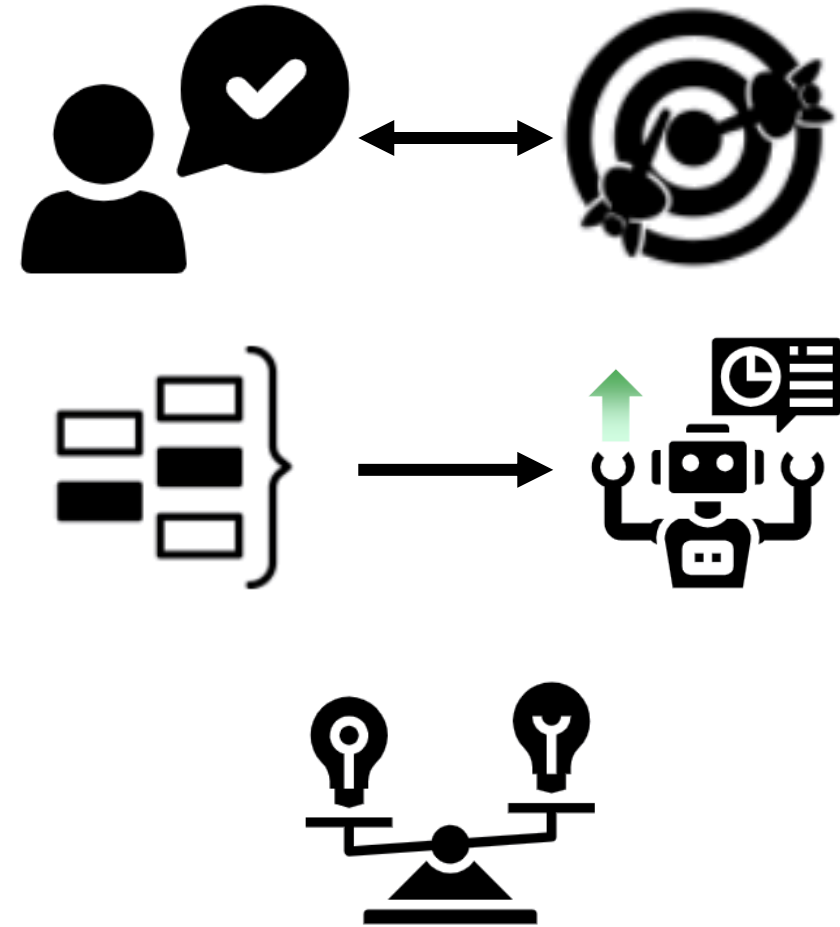
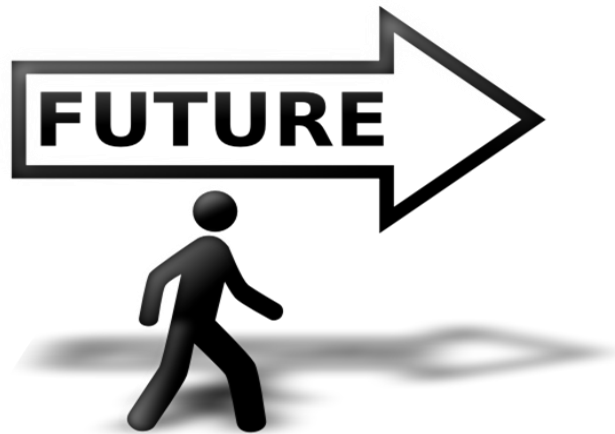
Method	Training Error Rate	Test Error Rate	Predicted LRT Given True LRT*
<i>Logistic Regression</i>	0.250	0.249	1,188 (-46%)
<i>Baseline Decision Tree</i>	0.183	0.256	1,350 (-39%)
<i>Pruned Decision Tree</i>	0.242	0.240	1,204 (-45%)
<i>Most Pruned Tree</i>	0.251	0.249	1,115 (-50%)
<i>Random Forest</i>	0.190	0.236	1,175 (-47%)
<i>Neural Network</i>	0.260	0.261	1,159 (-47%)

* Numbers in the parentheses: The difference proportion of the predicted LRT ridership compared to the actual ridership (2,209)

Future Work

■ Potential approaches

1. Implementing more models such as SGD Classifier or LDA
2. Applying Factors analysis since variables might have underlying relationships
3. Better feature engineering
4. Fine tuning Neural Network
5. Comparing model interpretability versus prediction accuracy



Thank you for your attention!

Questions and Answers

