

Hotels Dataset

MGSC 310 - 02

Arshia Sarma | Betsy Heredia | Kayla Cho | Nawal Alasmari

Motivation

- Determine which factors would affect a hotel business and advance their profit
- Focusing on predicting why customers would cancel their reservation
 - The more guests that cancel = less profit earned
- Useful for hotels to find which factors of their operation to improve on to keep customers from cancelling
 - Profits will increase as more guests stay committed



Exploratory Analysis & Summary Stats

- Dataset consists of information from hotels from Portugal
- Total dataset → 119,390 observations and 32 columns
- Pruned variables agent, company, country, reservation_status, and reservation_status_date
- 37% cancelled their reservation (for is_canceled)
- After pruning, we also upsampled our model with more repeated copies of our minority class

```
Rows: 119,390
Columns: 32
$ hotel
$ isCanceled
$ lead_time
$ arrival_date_year
$ arrival_date_month
$ arrival_date_week_number
$ arrival_date_day_of_month
$ stays_in_weekend_nights
$ stays_in_week_nights
$ adults
$ children
$ babies
$ meal
$ country
$ market_segment
$ distribution_channel
$ is_repeated_guest
$ previous_cancellations
$ previous_bookings_not_cancelled
$ reserved_room_type
$ assigned_room_type
$ booking_changes
$ deposit_type
$ agent
$ company
$ days_in_waiting_list
$ customer_type
$ adr
$ required_car_parking_spaces
$ total_of_special_requests
$ reservation_status
$ reservation_status_date

<fct> Resort Hotel, Resort Hotel, Resort Hotel, Resort Hotel, Resor...
<fct> 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<int> 342, 737, 7, 13, 14, 0, 9, 85, 75, 23, 35, 68, 18, 37, 68...
<int> 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2...
<fct> July, J...
<int> 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 27, 2...
<int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
<int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1...
<int> 0, 0, 1, 1, 2, 2, 2, 2, 3, 3, 4, 4, 4, 4, 4, 4, 4, 4, 1, 1, 4, 4...
<int> 2, 2, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 1...
<int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<fct> BB, BB, BB, BB, BB, FB, BB, HB, BB, BB, HB, BB, BB, BB, B...
<fct> PRT, PRT, GBR, GBR, GBR, PRT, PRT, PRT, PRT, PRT, U...
<fct> Direct, Direct, Direct, Corporate, Online TA, Online TA, Dire...
<fct> Direct, Direct, Direct, Corporate, TA/TO, TA/TO, Direct, Dire...
<fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<fct> C, C, A, A, A, C, C, A, D, E, D, G, E, D, E, A, G, F...
<fct> C, C, C, A, A, C, C, A, D, E, D, E, G, E, E, G, G, G, F...
<int> 3, 4, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<fct> No Deposit, No Deposit, No Deposit, No Deposit, No Deposit, N...
<fct> NULL, NULL, NULL, 304, 240, 240, NULL, 303, 240, 15, 240, 240...
<fct> NULL, NULL, NULL, NULL, NULL, NULL, NULL, NULL, NULL, N...
<int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<fct> Transient, Transient, Transient, Transient, Transient, Transi...
<dbl> 0.00, 0.00, 75.00, 75.00, 98.00, 98.00, 107.00, 103.00, 82.00...
<int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
<int> 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 3, 1, 0, 3, 0, 0, 0, 1, 1...
<fct> Check-Out, Check-Out, Check-Out, Check-Out, Check-Out, Check...
<chr> "7/1/15", "7/1/15", "7/2/15", "7/2/15", "7/3/15", "7/3/15", "...
```

Exploratory Analysis & Summary Stats

```
> nlevels(hotels$agent)
[1] 334
>
> nlevels(hotels$company)
[1] 353
>
> nlevels(hotels$country)
[1] 178
```

Total of 865 factors from agent,
company and country

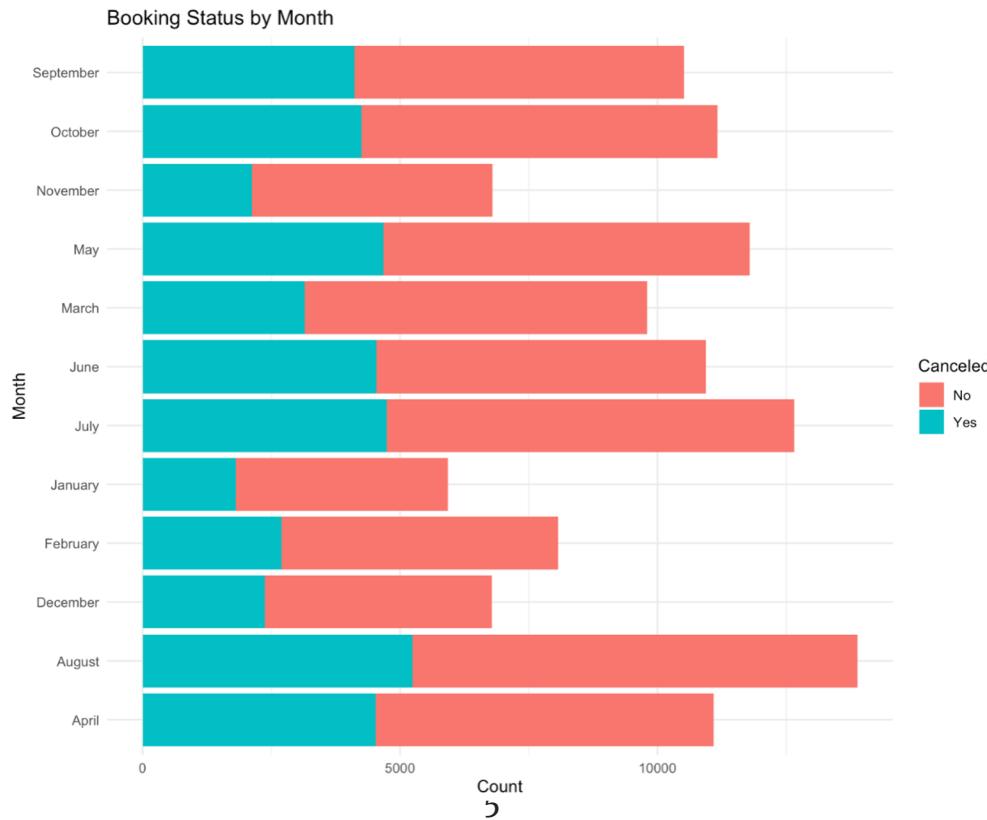
```
is_canceled
0:75166
1:44224
```

```
reservation_status
Canceled :43017
Check-Out:75166
No-Show  : 1207
```

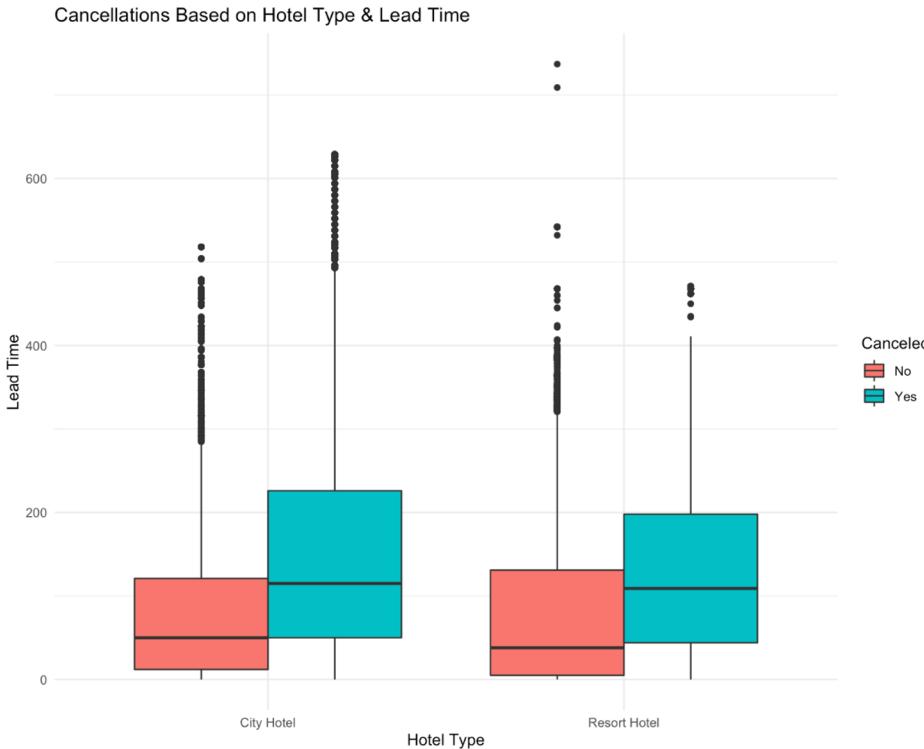
Multicollinearity

```
reservation_status_date
Min.    :2014-10-17
1st Qu.:2016-02-02
Median   :2016-08-07
Mean     :2016-07-30
3rd Qu.:2017-02-08
Max.     :2017-09-14
```

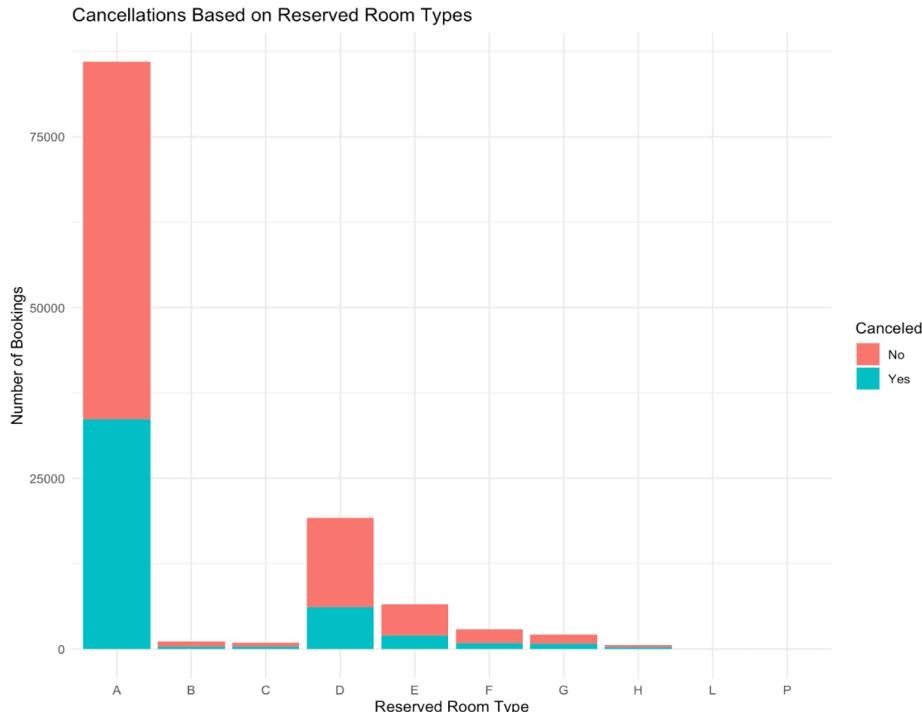
Data Visualizations: Booking Status v. Month



Data Visualizations: Hotel Type v. Lead Time



Data Visualizations: Cancellations v. Room Type



Logistic Model #1

- Predicting the categorical outcome of if a guest canceled their reservation
- Confusion matrix of training and test

Training

		Truth	
P R E D I C T I O N	0	0	1
	0	41,104	10,603
	1	19,073	49,486

Test

		Truth	
P R E D I C T I O N	0	0	1
	0	10,175	2,620
	1	4,813	12,457

Training

Accuracy: 75.32%
Sensitivity: 68.31%
Specificity:
82.35%

Test

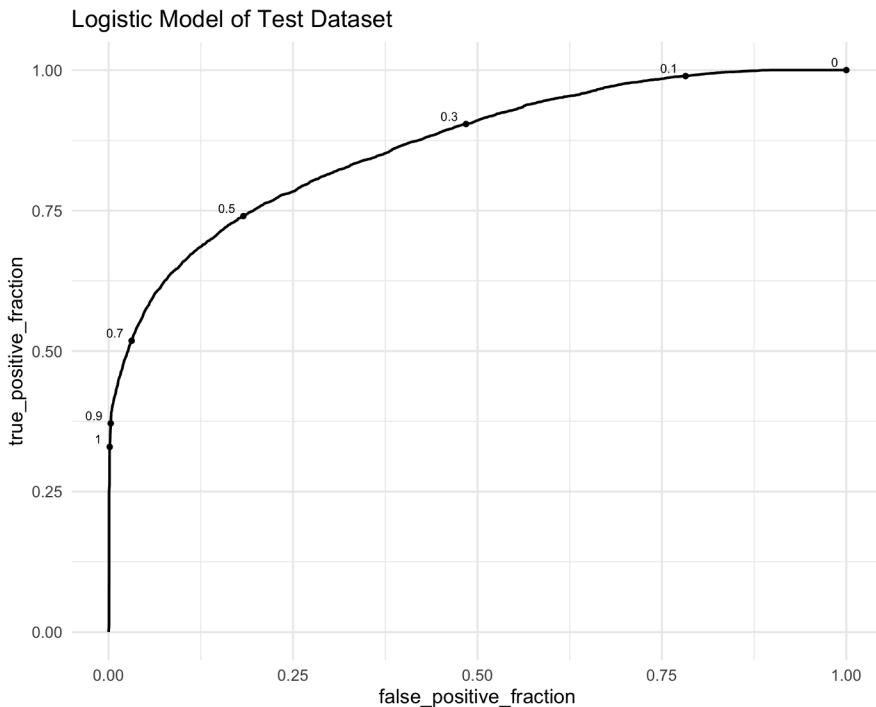
Accuracy: 75.27%
Sensitivity: 67.89%
Specificity: 82.62%

Scores and Takeaways

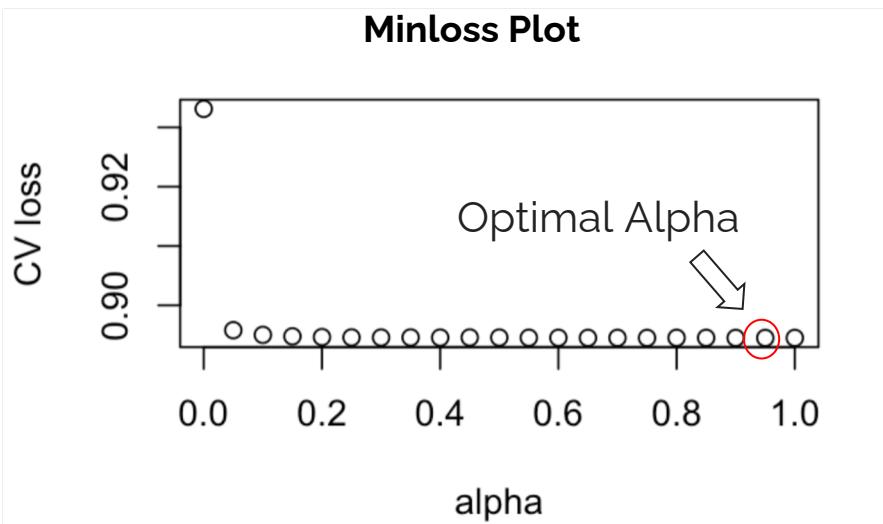
Conclusions

- Not the best model → ~0.26 R² and average accuracy score
- Accuracy, sensitivity, and specificity of the training and test are similar, as well as their R² and RMSE

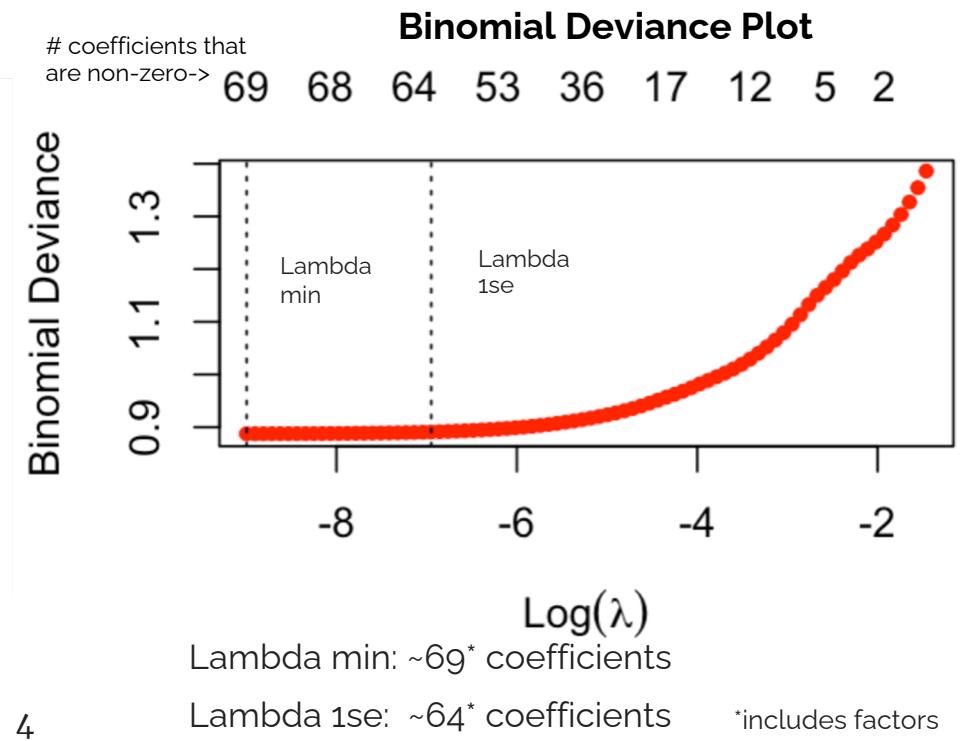
	AUC	R ²
Training	0.8644	0.2617
Test	0.8649	0.2609



Elastic Net Model



- Optimal Model Parameters:
- Alpha= 0.95

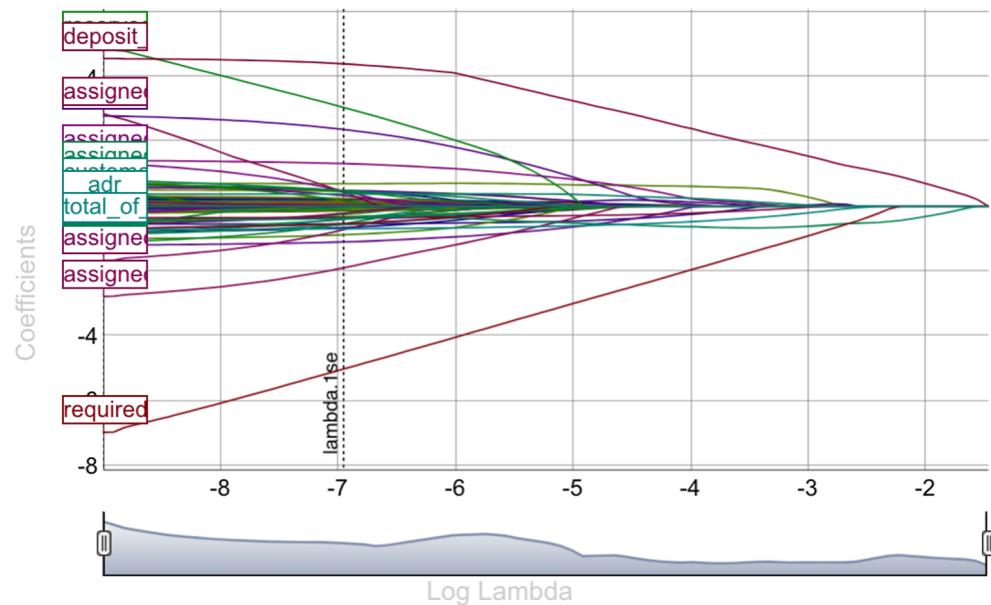


Elastic Net Coefficients

Noteable Coefficients

Required Car Parking Spaces	5.021
Deposit Type- Non Refund	4.394
Reserved Room Type- P	3.048
Previous Cancellations	2.368
Assigned Room Type - I	-1.9

Coefficients Plot



Elastic Net Coefficients Cont.

24 Variables

- | | |
|--------------------------------|-----------------------------|
| 1. Required car parking spaces | 16. Guests |
| 2. Deposit Type | 17. Hotel |
| 3. Reserved Room Type | 18. Babies |
| 4. Previous Cancellations | 19. Children |
| 5. Assigned Room Type | 20. Stays in Week Nights |
| 6. Distribution Channel | 21. Stays in Weekend Nights |
| 7. Market Segment | 22. Arrival Date Year |
| 8. Total of Special Requests | 23. ADR |
| 9. Meal | 24. Lead time |
| 10. Customer Type | |
| 11. Booking Changes | |
| 12. Previous Bookings | |
| 13. Is Repeated Guest | |
| 14. Adults | |
| 15. Arrival Date Month | |

**80 Coefficients ->
61 Coefficients**

Elastic Net Coefficients Cont.

Non-Important Variables

- **Days in waiting list**
- **Arrival date week number**
- **Arrival date day of month**

Combined Coefficients

- **Reserved_room_type_Other**
 - Reserved_room_type (C,D,H, L)
- **Market_Segement_Other**
 - Market_Segment (Aviation, Groups, Undefined)
- **Distribution_Channel_Other**
 - Distribution_Channel (Corporate, TA/TO, Undefined)

```
hotels_train_clean <- hotels_train %>% mutate(market_segment =  
fct_other(market_segment, drop= c("Aviation", "Direct", "Groups",  
"Undefined")),  
                                              reserved_room_type  
=fct_other(reserved_room_type, drop= c("C", "D", "H", "L")),  
distribution_channel=fct_other(distribution_channel,  
drop=c("Corporate", "TA/T0", "Undefined")))
```

```
hotels_test_clean <- hotels_test %>% mutate(market_segment =  
fct_other(market_segment, drop= c("Aviation", "Direct", "Groups",  
"Undefined")),  
                                              reserved_room_type  
=fct_other(reserved_room_type, drop= c("C", "D", "H",  
"L")),distribution_channel=fct_other(distribution_channel,  
drop=c("Corporate", "TA/T0", "Undefined")))
```

Elastic Net Model Results

<u>Train</u>	Truth		
P R E D I C T I O N		0	1
	0	40860	10505
	1	19317	49584

<u>Test</u>	Truth		
P R E D I C T I O N		0	1
	0	10135	2578
	1	4853	12499

Conclusions

- Similar performance as Logistic Regression (Test: 75.27%)
- Slightly better test performance
- Higher Specificity -> Better at predicting if a booking is not canceled

Accuracy

Train: 75.2033%
Test: 75.2836%

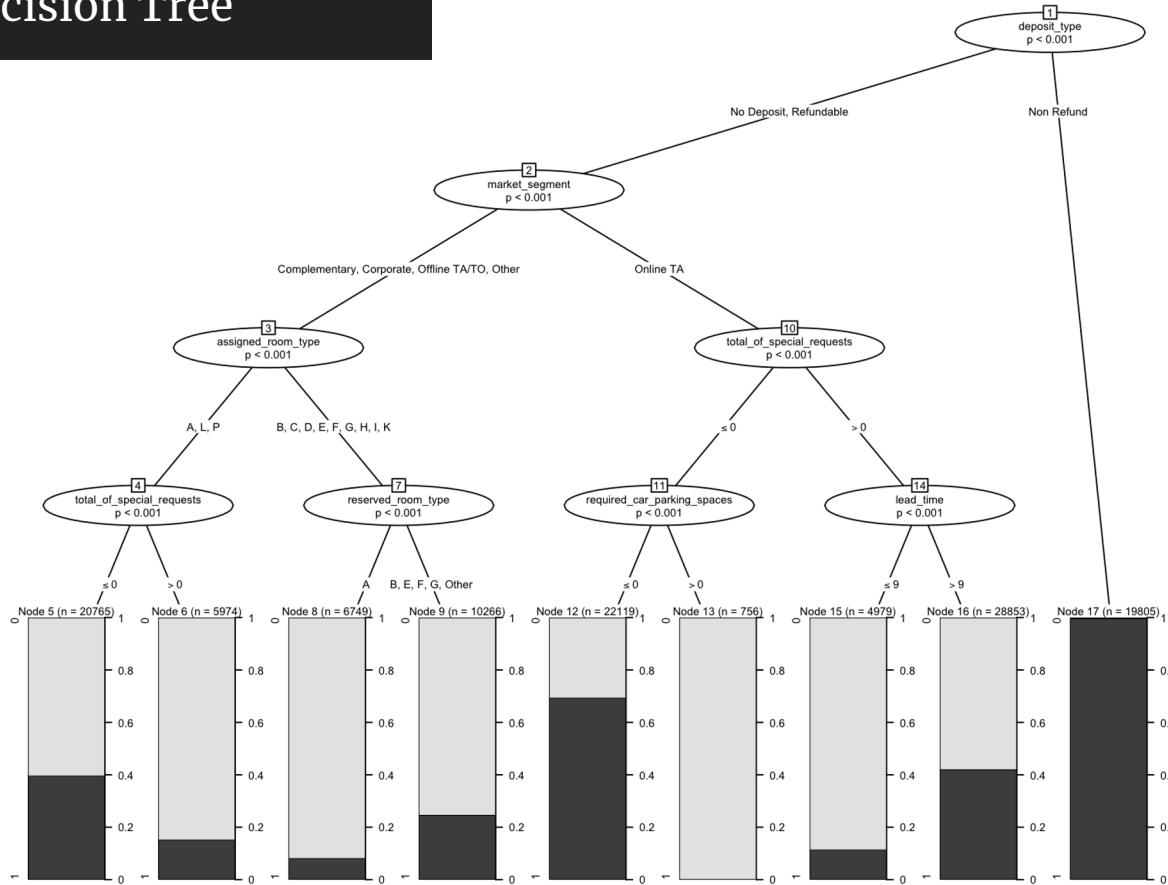
Sensitivity

Train: 67.8997%
Test: 67.6208%

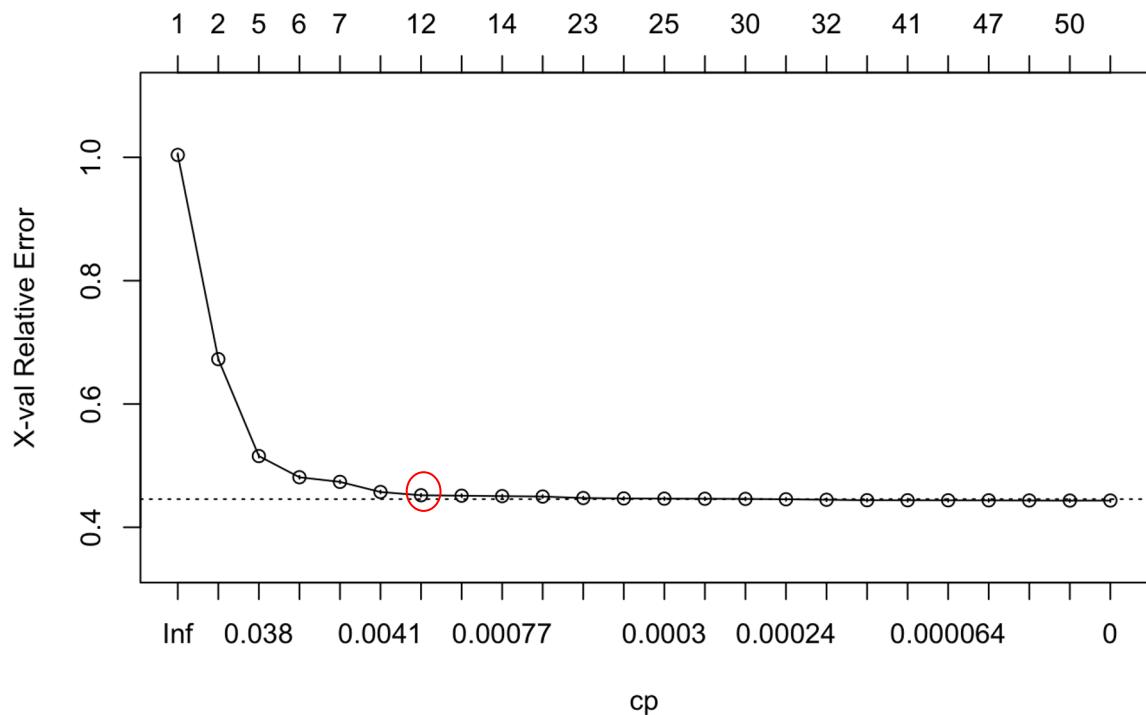
Specificity

Train: 82.5176%
Test: 82.9011%

Decision Tree



RPart



Decision Tree

Max Depth = 4

	RMSE	RSQ
Train	0.5141	0.2443
Test	0.5147	0.2438

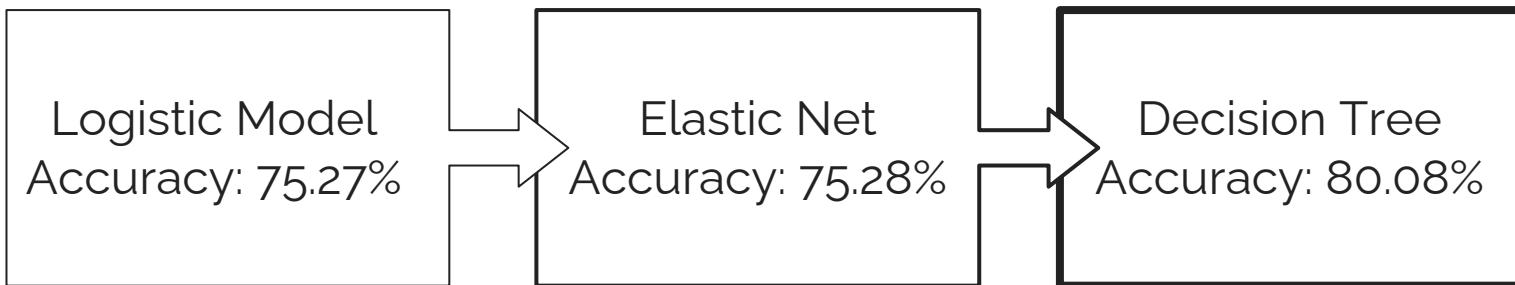
P R E D I C T I O N	Truth		
		0	1
	0	13,280	6,257
	1	1,708	8,820

Max Depth = 12

	RMSE	RSQ
Train	0.4427	0.3698
Test	0.4463	0.3620

P R E D I C T I O N	Truth		
		0	1
	0	11,932	2,932
	1	3,056	12,145

Comparison of Models



Conclusion and Final Takeaways

- Understand your dataset and what variables you want to predict
- Remove outliers, use the entire training sample for the elastic net model and model the outcome variable with a random forest.
- Logistic model with all variables performed better than with top 10 variables.
- Company should focus on reserve room type and people who cancelled their bookings in the past.
- Future directions:
 - Predict market segment
 - Predict distribution channel

Thanks!

ANY QUESTIONS?

Team Hotels

Credits

Special thanks to all the people who made and released these awesome resources for free:

- Presentation template by [SlidesCarnival](#)
- Photographs by [Unsplash](#)