

Question 1

(a) Model Construction

To begin, in Excel, we construct a binary logistic regression model to predict the probability of a wagon being sold at auction.

	<i>coeff b</i>	<i>s.e.</i>	<i>Wald</i>	<i>p-value</i>	<i>exp(b)</i>	<i>lower</i>	<i>upper</i>
Intercept	0.75866104	1.562302952	0.235812	0.6272477	2.13542		
X1 (Age)	-0.486756233	0.376927637	1.667658	0.1965725	0.61462	0.2936	1.28661
X2 (Condition)	3.210648662	1.447464045	4.920055	0.0265467	24.7952	1.45306	423.107

Figure 1.1 The constructed Model in Excel

We have the coefficient respectively:

$$\beta_0 \approx 0.7587$$

$$\beta_1 \approx -0.4868$$

$$\beta_2 \approx 3.2106$$

Given the outcome (Y), we want to predict whether or not a given wagon will be sold at the auction, with Y = 1 meaning the wagon is sold at the auction and Y = 0 meaning the wagon is not sold at the auction. The constructed binary logistic regression model can be interpreted as follows.

The intercept (β_0)

$$\beta_0 \approx 0.7587$$

The intercept represents the log-odds of a wagon being sold when all predictors (wagon age – X1 and condition – X2) are zero. Although, in reality, it is unlikely to see a wagon with no age (X1 = 0) and in an average condition (X2 = 0) at the same time, this provides a baseline from which the other coefficients adjust the log-odds.

Coefficient for Wagon Age (β_1)

$$\beta_1 \approx -0.4868$$

This coefficient represents the change in the log-odds of a wagon being sold for each one-year increase in its age. Since β_1 is negative, increasing the wagon's age decreases the log-odds of it being sold.

$$\text{Odds Ratio} = e^{\beta_1} \approx e^{-0.4868} \approx 0.6146$$

An odds ratio of 0.6146 means that for each additional year of age, the odds of the wagon being sold are multiplied by approximately 0.6146. In other words, for each additional year of age, the percentage decrease in odds of the wagon being sold is approximately 38.54% ($1 - 0.6146 = 0.3854$).

Coefficient for Wagon Condition (β_2)

$$\beta_2 \approx 3.2106$$

This coefficient represents the change in the log-odds of a wagon being sold when the condition changes from the baseline category (average, where $X_2 = 0$) to the other category (good, where $X_2 = 1$). Since β_2 is positive, having a good condition increases the odds of the wagon being sold.

$$\text{Odds Ratio} = e^{\beta_2} \approx e^{3.2106} \approx 24.7952$$

An odds ratio of 24.7952 means that a wagon in good condition has odds of being sold that are approximately 24.8 times higher than a wagon in average condition. This represents a significant increase in the odds of selling the wagon when its condition is good.

		coeff b	s.e.	Wald	p-value	exp(b)	lower	upper
50								
51	Intercept	0.75866104	1.562302952	0.2358115	0.6272477	2.13542		
52	X1 (Age)	-0.486756233	0.376927637	1.6676583	0.1965725	0.61462	0.2936	1.28661
53	X2 (Condition)	3.210648662	1.447464045	4.9200555	0.0265467	24.7952	1.45306	423.107
54								
55	Odds Ratio (X1) =			0.6146 = EXP(C52)				
56	Odds Ratio (X2) =			24.7952 =EXP(C53)				

Figure 1.2 Odds ratio calculations in Excel

Model Performance

1. Chi-square Goodness of Fit

To evaluate the significance of the model performance, we constructed the hypothesis as follows.

H_0 : The model does not fit the data better than a model with no predictors.

H_A : The model fits the data better than a model with no predictors.

Since the p-value (0.0076) is less than significant level α (0.05), we reject the null hypothesis (H_0). Thus, there is significant evidence to suggest that the binary logistic regression model fits the data better than a model with no predictors. In other words, this indicates that the predictors used in the model are significantly associated with the probability of a wagon being sold at auction.

2. Testing predictor significance

	coeff b	s.e.	Wald	p-value	exp(b)	lower	upper
Intercept	0.75866104	1.562302952	0.2358115	0.6272477	2.13542		
X1 (Age)	-0.486756233	0.376927637	1.6676583	0.1965725	0.61462	0.2936	1.28661
X2 (Condition)	3.210648662	1.447464045	4.9200555	0.0265467	24.7952	1.45306	423.107

Figure 2 P-value of the predictors

Coefficient for Wagon Age (β_1)

To consider the significance of wagon age in the model, we have the hypothesis that:

H_0 : Wagon age is not significant factor ($H_0: \beta_1 = 0$)

H_A : Wagon age is a significant factor ($H_A: \beta_1 \neq 0$)

As computed in Excel, we have p-value (β_1) $\approx 0.1966 > 0.05$. Since we have p-value (0.1966) larger than significant level α (at 0.05), we fail to reject H_0 . Therefore, conclude β_1 is zero and that wagon age is not significant at $\alpha = 0.05$.

Coefficient for Wagon Condition (β_2)

To consider the significance of wagon age in the model, we have the hypothesis that:

H_0 : Wagon age is not significant factor ($H_0: \beta_2 = 0$)

H_A : Wagon age is a significant factor ($H_A: \beta_2 \neq 0$)

As computed in Excel, we have p-value (β_2) $\approx 0.0265 < 0.05$. Since we have p-value (0.0265) less than significant level α (at 0.05), we reject H_0 . Therefore, conclude β_2 is not zero and that wagon condition is significant at $\alpha = 0.05$.

3. Classification Table

Classification Table			
	Suc-Obs	Fail-Obs	
Suc-Pred	6	1	7
Fail-Pred	3	10	13
	9	11	20
Accuracy	0.66667	0.90909	0.8
Cutoff	0.5		

Figure 3: Classification Table in Excel

At the cutoff probability of 0.5, the model has an accuracy of approximately 66.67% for successful auctions ($Y=1$, sold) and 90.91% ($Y=0$, not sold) for unsuccessful auctions. This implies that the model is **moderately accurate** in predicting the outcome for the wagon auction based on the wagon's age and condition.

However, while this might suggest that the model is moderately effective in distinguishing between wagons likely to be sold at auction and those that are not, it is likely weak in predicting successful cases (with an accuracy of 66.67%). Thus, this can be achieved through further optimization of the model's cutoff threshold and feature selection.

4. AUC – ROC

1. AUC (Area Under the Curve) ≈ 0.8384

It suggests a reasonably good ability of the constructed binary logistic regression model to distinguish between the two different outcomes of the auction based on the wagon's age and condition, given that an AUC value of 1 implies perfect discrimination, and 0.5 represents random chance.

2. The ROC Curve

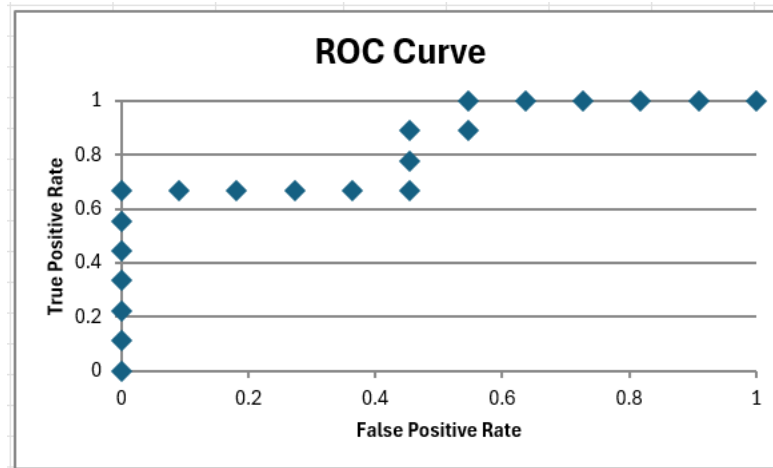


Figure 4: ROC Curve from the model

The curve shows the trade-off between the true positive rate (sensitivity) and the false positive rate (1 - specificity) at different cutoffs. As it stays close to the upper left corner, it indicates that the model performs moderately well.

- (b) To predict the probability that a 7-year-old wagon in average condition will sell at auction using the logistic regression model constructed above, we have the formula for the probability of selling at auction as follows.

$$P(Y=1, \text{ sold at auction}) = \frac{e^{\beta_0 + \beta_1 \times \text{Age} + \beta_2 \times \text{Condition}}}{1 + e^{\beta_0 + \beta_1 \times \text{Age} + \beta_2 \times \text{Condition}}}$$

As constructed in question 1a, we have:

$$\beta_0 \approx 0.7587$$

$$\beta_1 \approx -0.4868$$

$$\beta_2 \approx 3.2106$$

Given:

Age of the wagon = 7 (in years)

Condition = 0 (for average condition)

Thus, substitute to the equation above, we have:

$$P(Y=1) \approx \frac{e^{0.7587 - 0.4868 \times 7 + 3.2106 \times 0}}{1 + e^{0.7587 - 0.4868 \times 7 + 3.2106 \times 0}} \approx 0.0661$$

	A	B	C	D	E	F	G	H	I
50			<i>coeff b</i>	<i>s.e.</i>	<i>Wald</i>	<i>p-value</i>	<i>exp(b)</i>	<i>lower</i>	<i>upper</i>
51		Intercept	0.75866104	1.562302952	0.2358115	0.6272477	2.13542		
52		X1 (Age)	-0.486756233	0.376927637	1.6676583	0.1965725	0.61462	0.2936	1.28661
53		X2 (Condition)	3.210648662	1.447464045	4.9200555	0.0265467	24.7952	1.45306	423.107
54									
55		Odds Ratio (X1) =		0.6146	= EXP(C52)				
56		Odds Ratio (X2) =		24.7952	=EXP(C53)				
57									
58	(b)	P (Y=1) =		0.0661	=EXP(C51+C52*7+C53*0)/(1+EXP(C51+C52*7+C53*0))				
59				6.61%					

Figure 5: Compute the probability of selling at auction of the given wagon (in Excel)

The estimated probability that a 7-year-old wagon in average condition will sell at auction is approximately 6.61%. Using the cut-off at 0.5, we can infer that this wagon is unlikely to be sold at the auction. However, considering the model's 90.91% accuracy in predicting unsuccessful cases, this outcome is highly probable.

Question 2

Guest Speaker 1: Rijwa Abbas

- Current Job Title: Risk Analyst
- Organization: Ernst & Young (EY)
- [Rijwa Abbas's LinkedIn](#)

In Rijwa Abbas's role as a Risk Analyst at EY, she might employ **binary logistic regression analysis** to predict the likelihood of fraudulent activities within a financial institution.

In general, using historical transaction data, Rijwa could identify patterns and anomalies indicative of potential fraud. Logistic regression allows her to model the relationship between a binary outcome (fraudulent or non-fraudulent transaction) and multiple predictor variables such as transaction amount, frequency, location, and user behavior.

To begin, Rijwa would preprocess the data, cleaning and transforming it as necessary to ensure its quality and relevance. Next, she would then split the dataset into training and testing sets to develop and validate the logistic regression model. After fitting the model to the training data, Rijwa would evaluate its performance using metrics like accuracy, precision, recall, and F1-score on the test set. By analyzing the coefficients of the predictor variables, she can identify which factors contribute most significantly to the likelihood of fraud.

Ultimately, Rijwa would help her clients safeguard their financial assets and empower them to make informed decisions regarding resource allocation, fraud detection technology investments, and strategic planning to mitigate risks effectively in the ever-evolving landscape of financial security.

Guest Speaker 2: Laura Heo

- Current Job Title: Performance Education Lead
- Organization: Everperform
- [Laura Heo's LinkedIn](#)

In Laura Heo's role as the Performance Education Lead at Everperform, she could utilize **two-sample tests** to gain insights into the productivity trends and employee performance within her client's organisations. Since two-sample tests involve comparing the means of two independent groups to determine if there is a significant difference between them, Laura could apply this technique to evaluate the performance metrics of different teams, departments, or time periods.

For instance, Laura could collect data on sales metrics such as revenue, number of deals closed, and customer satisfaction scores from two different sales teams within a company. By conducting a two-sample test, she can determine if there is a statistically significant difference in the average *performance* between the two teams.

Additionally, Laura could use two-sample tests to evaluate the *productivity* levels of employees before and after the implementation of a new well-being program or productivity initiative. By analyzing the data using appropriate statistical tests, she can evaluate the effectiveness of these interventions in improving employee performance and well-being.

To summarize, by leveraging two-sample tests, Laura can assist her clients in making data-driven decisions to optimize sales strategies, and enhance overall performance and well-being within Everperform.

Guest Speaker 3: Richard Norrie

- Current Job Title: Cyber Strategy and Advisory Consultant
- Organization: Dell Technologies
- [Richard Norrie's LinkedIn](#)

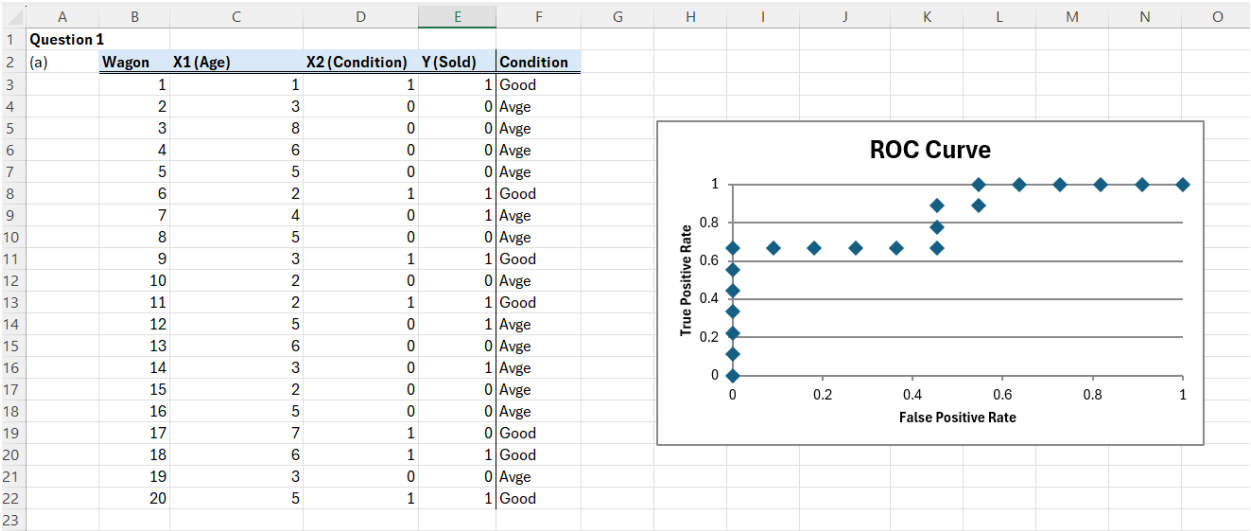
In his role as a Cyber Strategy and Advisory Consultant at Dell, Richard Norrie could employ **One-way Analysis of Variance (ANOVA)** to analyze and interpret various cybersecurity-related variables for his clients.

For example, Richard could collect data on the severity of cyber threats detected by a company's security systems, categorized into different types such as malware, phishing, and ransomware. By organizing this data and calculating the average severity within each category, Richard can then apply ANOVA to determine if there are statistically significant differences in the severity levels among different types of cyber threats. In the end, this technique would allow Richard to assess whether certain types of cyber threats pose a significantly greater risk to the organization's systems and operations compared to others. Consequently, it empowers him to provide advice to his clients on identifying priority areas for cybersecurity efforts and allocating resources effectively to mitigate the most critical threats.

Additionally, Richard could use ANOVA to compare the effectiveness of various cybersecurity measures in mitigating different types of cyber threats. By analyzing variance between groups, he can determine which strategies are most effective in minimizing the impact of cyber-attacks and recommend adjustments or enhancements to his client's organization accordingly.

In short, by leveraging ANOVA, Richard can make data-driven recommendations and provide tailored advisory services to help organizations strengthen their cybersecurity defenses, minimize vulnerabilities, and safeguard against potential threats effectively.

Appendix A: Data Table and ROC Curve



Appendix B: Binary Logistic Regression (Summary Table)

	A	B	C	D	E	F	G	H	I	J	K	L	M
24													
25		Logistic Regression											
26													
27		X1 (Age)	X2 (Condition)	Success	Failure	Total	p-Obs	p-Pred	Suc-Pred	Fail-Pred	LL	%Correct	
28		1	1	1	0	1	1	0.97019	0.97019	0.0298127	-0.03027	100	
29		3	0	0	1	1	0	0.33146	0.33146	0.6685441	-0.40265	100	
30		8	0	0	1	1	0	0.04167	0.04167	0.9583291	-0.04256	100	
31		6	0	0	1	1	0	0.10323	0.10323	0.8967734	-0.10895	100	
32		5	0	0	1	1	0	0.15774	0.15774	0.8422573	-0.17167	100	
33		2	1	1	0	1	1	0.95238	0.95238	0.0476161	-0.04879	100	
34		4	0	1	0	1	1	0.23355	0.23355	0.7664483	-1.45435	0	
35		5	0	0	1	1	0	0.15774	0.15774	0.8422573	-0.17167	100	
36		3	1	1	0	1	1	0.92477	0.92477	0.0752268	-0.07821	100	
37		2	0	0	1	1	0	0.44649	0.44649	0.5535072	-0.59148	100	
38		2	1	1	0	1	1	0.95238	0.95238	0.0476161	-0.04879	100	
39		5	0	1	0	1	1	0.15774	0.15774	0.8422573	-1.84679	0	
40		6	0	0	1	1	0	0.10323	0.10323	0.8967734	-0.10895	100	
41		3	0	1	0	1	1	0.33146	0.33146	0.6685441	-1.10426	0	
42		2	0	0	1	1	0	0.44649	0.44649	0.5535072	-0.59148	100	
43		5	0	0	1	1	0	0.15774	0.15774	0.8422573	-0.17167	100	
44		7	1	0	1	1	0	0.63692	0.63692	0.3630811	-1.01313	0	
45		6	1	1	0	1	1	0.74054	0.74054	0.2594609	-0.30038	100	
46		3	0	0	1	1	0	0.33146	0.33146	0.6685441	-0.40265	100	
47		5	1	1	0	1	1	0.82281	0.1771862	-0.19503	100	100	
48				9	11	20			9	11	-8.88373	80	

[illegible]

Appendix C: Model Performance Metrics

[illegible]

Appendix D: ROC Table

[illegible]

Appendix E: Question 1b Calculation

[illegible]