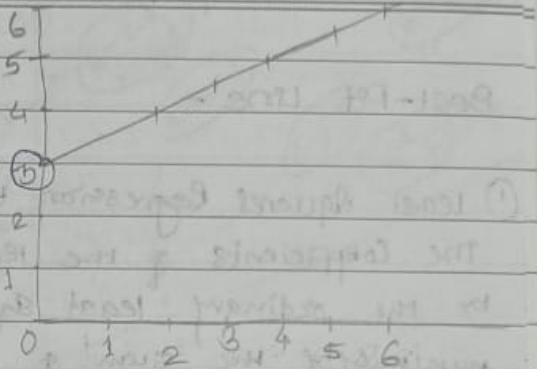


Regression

i | nautix.

A BNY MELLON COMPANY



① Intercept of the straight line -
what is the intercept of the given line? use the graph. Dependent given above to answer this question

-3, the value of y when $x=0$ in the x -Independent variable given straight line is 3, so -3 would be the intercept in this case

② Slope of a straight line -
what is the slope of the given line? use the graph given above to answer this question.

$-1/2$, the slope of any straight line can be calculated by $(y_2 - y_1) / (x_2 - x_1)$ where (x_1, y_1) & (x_2, y_2) are any two points through which the given line passes. This line passes $(0, 3)$ & $(2, 4)$, so the slope of this line would be

$$((4-3)/(2-0)) = 1/2$$

$$\frac{(y_2 - y_1)}{(x_2 - x_1)} \therefore \frac{(x_1, y_1), (x_2, y_2)}{0, 3 \quad (2, 4)} \therefore \frac{y_2 - y_1}{x_2 - x_1}$$

③ what - Equation of a straight line -
what would be the equation of the given line.

$$-y = \frac{x}{2} + 3 \text{ ; the std eqn of a straight line } y = mx + c$$

where m is the slope & c is the intercept. In this case

$$y = \frac{x}{2} + 3$$

Best-fit Line -

① Least Squares Regression Line -

The coefficients of the least squares regression line are determined by the ordinary least squares method — which basically means minimising the sum of the squares of the

y -coordinates of actual data — x -coordinates of predicted data

The ordinary least squares method has the criterion of the minimisation of the sum of square of residuals. Residuals are defined as the differences between the y -coordinates of actual data & the y -coordinates of predicted data.

② Best Fit Regression Line -

what is the main criterion used to determine the best fitting regression line?

— the line that minimises the sum of squares of distances of points from the regression line —

the criterion is given by the Ordinary Least Squares (OLS) method, which states that the sum of the squared of residuals should be minimum. This is explained by the option.

straight Strength of Simple Linear Regression

① Residual Sum of Squares (RSS)

Find the value of RSS for this regression line

- 6.25 ∵ the residual for all 5 points are -0.5, 1, 0, -2, 1. The sum of squares of all 5 residuals would be

$$0.25 + 1 + 4 + 1 = 6.25$$

② Total Sum of Squares (TSS)

Find the value of TSS for this regression line

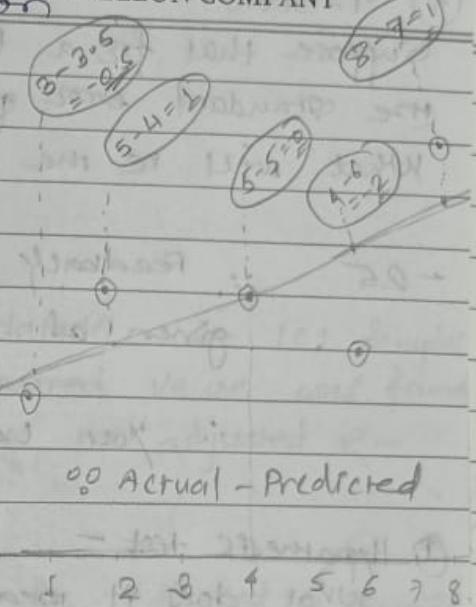
- 14 ∵ the average of y-value for all data points $(3+5+5+9+8)/5 = 25/5 = 5$. So $\bar{y}-\bar{y}$ term for each data point would be $-2, 0, 0, 1, 3$. So the squared sum of these terms would be $4+1+9=14$

- ③ R^2 - The RSS for this example comes out to be 6.25 & TSS comes out to be 14.

What would be the R^2 for this regression line

$$- 1 - (6.25 / 14) \quad \because R^2 \text{ value is given by } 1 - \frac{\text{RSS}}{\text{TSS}}$$

$$\text{So in this case } R^2 = 1 - \left(\frac{6.25}{14} \right) = 1 - 0.4464 = 0.5536$$



Hypothesis Testing for Linear Regression -

- ④ T-score -
- Suppose that for a linear Model, you got β_1 at 0.5. Also the standard error of β_1 was found out to be 0.02. What will be the value of t-score for β_1 ?

-25 ∵ Feedback Recall that the t-score for β_1 is given as

$$\frac{\beta_1}{SE(\beta_1)}$$

∴ You have : t-score = $\frac{0.5}{0.02} = 25$

① Hypothesis test -

What does it mean if you fail to reject the Null hypothesis in the case of simple linear Regression.

Given → Null Hypothesis (H_0) : $\beta_1 = 0$

Alternate Hypothesis (A_1) : $\beta_1 \neq 0$

β_1 & thus, the independent variable. It is associated with the insignificant in the prediction of the dependent variable.

feedback - Correct! the Null Hypothesis in simple linear regression is : $\beta_1 = 0$

thus, if we fail to reject the Null Hypothesis, it means that β_1 is indeed zero, & thus significant for the prediction of the independent variable.

Model Assessment & Comparison -

$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

① Calculating Adjusted R-Squared -

When a model was built from a dataset with 101 samples and 10 predictor variables, the R-squared value was found to be 0.7. What will the value of the adjusted R-squared be for the same model?

→ 0.67 : Feedback : The formula for Adjusted R-squared is given as :

$$1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

So, substituting the given values at the appropriate places gives us

$$1 - \frac{(1-0.7)(101-1)}{101-10-1} \approx 0.67$$

② Model Assessment

→ After performing inferences on a linear model built with several variables, you concluded that the variable 'r' was insignificant. This meant that the variable 'r'

→ Had a high P-value.

A high p-value means that the variable is not significant, and hence, doesn't help much in prediction.

R-squared values -

Suppose you built a model with some features. Now you add another variable to the model. Which of the following statements would be true.

- R-squared value will either increase or remain the same.

- The Adjusted R-squared value may increase, decrease

② Overfitting -

- Number of data points are less.
∴ Correct! Overfitting is the condition where in the model is so complex that it ends up memorising almost all the data points on the train set. Hence, this condition is more probable if the number of data points is less since the model passing through almost every point becomes easier.

③ VIF -

VIF is a measure of: $(1 - R^2)^{-1}$

- How well a predictor variable is correlated with all the other variables, excluding the target variable.

④ Calculating VIF -

Suppose you were predicting sales of a company using two variables 'Social Media Marketing' and 'TV Marketing'. You found out that the correlation between 'Social Media Marketing' & 'TV Marketing' is 0.9. What will be the approximate value of VIF for either of them?

- 5.26

Correct! The formula for VIF is given by:

$$VIF = \frac{1}{1 - R_i^2}$$

∴ Here, the R^2 -squared variable will be simply the correlation coefficient squared since we have only 2 variables. Hence, you have:

$$\therefore VIF = \frac{1}{1 - 0.9^2} \approx 5.26$$

$$= \frac{1}{1 - 0.81} = \frac{1}{0.19} \approx 5.26$$

① Dummy Variables -

Suppose you have 'n' Categorical variables, each with 'm' levels. How many dummy variables would you need to represent all the levels of all the categorical variables?

$-(m-1)*n$ ∵ Each of the dummy variable has 'm' levels. So to represent one categorical variable, you would require $(m-1)$ levels. Hence, to represent 'n' categorical variables, you would need $(m-1)*n$ dummy variables.

② Automated feature selection -

which of the following is/are an example of an automated approach for linear regression

- Recursive feature Elimination

- Stepwise Selection using AIC

- Regularisation

Redundant Variables -

After performing Inferences on a linear model built with several variables, you concluded that the variable 'r' was almost being described by other feature variable. This meant that the variable 'r' :

- Had a high VIF :: Correct! If the variable is being described well by the rest of the feature variable, it means it has a high VIF meaning it is redundant in the presence of the other variables.

① Comprehension Questions -

If $\beta_1 = \beta_2 = 0$ holds & $\beta_3 = 0$ fails to hold, then what can you conclude?

- There is no linear relationship between the outcome variable (Y) & x_3

Feedback - Since $\beta_3 = 0$ fails to hold, this means that x_3 is a significant variable in this linear regression model. Thus, we can say that there is a linear relationship between the outcome variable (Y) & x_3 .

② Comprehension Questions -

If $\beta_1 = \beta_2 = \beta_3 = 0$ holds true, then what can you conclude?

- There is no linear relationship between Y and any of the 3 independent variables

Feedback - If all the coefficients are found significant, then there cannot be a linear relationship between Y & any of the variables.

Python Coding Problem Mapping Variables

i | nautix.
A BNY MELLON COMPANY

'Yes' → 1 'No' → 0 'Maybe' → 0.5
'yes' → 1 'no' → 0 'maybe' → 0.5
'YES' → 1 'NO' → 0 'MAYBE' → 0.5

Import pandas as pd

```
df = pd.DataFrame({'Name': name, 'Response': response})
df['Response'] = df['Response'].str.lower()
df['Response'] = df['Response'].map({'yes': 1.0, 'maybe': 0.5,
                                      'no': 0})
```

print(df)

② Residuals -

In regression analysis, which of the statement is true.

- The mean of residuals is always equal to zero

Feedback - When a model gives you a "best fit" line, by

design it is made such that the mean of all residuals is
always zero

- The sum of residuals is always equal to zero

Feedback - When a model gives you a "best fit" line, by
design it is made such that the sum of all residuals is
always zero.

Sigmoid Curve -

This is the sigmoid curve equation:

$$y = P(\text{Diabetes}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} \quad \text{Here level say you}$$

take $\beta_0 = -15$ & $\beta_1 = 0.065$. Now, what will be the probability of diabetes for a patient with sugar level 220?

$$\begin{aligned} \therefore P(\text{Diabetes}) &= \frac{1}{1 + e^{(-15 + 0.065 \times 220)}} \\ &= \frac{1}{1 + e^{(0.7)}} \\ &= \frac{1}{1 + 2.01} \\ &= \frac{1}{3.01} \\ &\approx 0.33 \end{aligned}$$

(2) Sigmoid curve -

for the sigmoid curve ($\beta_0 = -15$ & $\beta_1 = 0.065$), what will be the probability of diabetes for a patient with sugar level 240?

$$\therefore P(\text{Diabetes}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$$\begin{aligned} &= \frac{1}{1 + e^{(-15 + 0.065 \times 240)}} \\ &= \frac{1}{1 + e^{(-15 + 15.6)}} \\ &= \frac{1}{1 + e^{-0.6}} \\ &= \frac{1}{1 + 0.5488} \end{aligned}$$

① Likelihood

Now, let's say that for the ten points in our example, the labels are as follows:

Point nos:	1	2	3	4	5	6	7	8	9	10
Diabetes:	No	No	No	Yes	No	Yes	Yes	Yes	Yes	Yes

In this case, the likelihood would be equal to:

$$-(1-P_1)(1-P_2)(1-P_3)(1-P_5)(P_4)(P_6)(P_7)(P_8)(P_9)(P_{10})$$

∴ Feedback -

Recall that likelihood is the product of $(1-P_i)$ for all non-diabetic patients & (P_i) for all diabetic patients. Hence, the likelihood is given by $(1-P_1)(1-P_2)(1-P_3)(1-P_5)$, (all non-diabetic patients) multiplied by $(P_4)(P_6)(P_7)(P_8)(P_9)(P_{10})$ (all diabetic patients).

A BNY MELLON COMPANY

① Log Odds -

So, let's say that the equation for the log odds is -

$$\ln\left(\frac{P}{1-P}\right) = -13.5 + 0.06x \quad \therefore \ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x$$

for $x=220$, the log odds are equal to

$$-13.5 + (0.06 * 220) = -0.3$$

For $x=231.5$, log odds are equal to:

$$-0.39 \quad \therefore -13.5 + (0.06 * 231.5) = 0.39$$

$$\frac{P}{1-P} = \text{odds}$$

$$\ln\left(\frac{P}{1-P}\right) = \text{Log Odds}$$

Shot on OnePlus

By pramodkhandare

① Standardising Variables -

a dataset with mean 50 & standard deviation 12, what will be the value of a variable with an initial value of 20 after you standardise it?

2.5

∴ the formula for standardising a value in a dataset given by $\frac{(X - \mu)}{\sigma}$

$$\therefore \frac{(20 - 50)}{12} = \frac{-30}{12} = -2.5$$

② Correlation (Table)

Command can be used to view the correlation table for the dataframe telecom? $\Rightarrow \text{telecom.corr()}$ will give you Correlation table for the dataframe telecom.

③ p-values -

After learning the coefficients of each variable, the model also produces a 'p-value' of each coefficient. Fill in the blanks so that the statement is correct:

The null hypothesis is that the coefficient is 0. If the p-value is small, you can say that the coefficient is significant & hence the null hypothesis is rejected.

- zero, can be rejected

∴ Feedback: Yes! Recall that the null hypothesis for any beta was $B_i = 0$. And if the p-value is small, you can say that the coefficient is significant, & hence you can reject the null hypothesis that $B_i = 0$.

③ Accuracy Calculation -

From the table you used for the last 2 questions, what will be the accuracy of the model?

Actual / Predicted	No	Yes
No	400	100
Yes	50	150

- 78.5% ∵ The accuracy of a model is given by:

$$\text{Accuracy} = \frac{\text{Correctly Predicted Labels}}{\text{Total Number of Labels}}$$

or

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FN + FP}$$

∴ The number of correctly predicted labels as you found out from the last question is equal to 550.

The total number of labels is $(400 + 100 + 50 + 150) = 700$.

Hence the accuracy becomes:

$$\text{Accuracy} = \frac{550}{700} = 0.785 \text{ or } 78.5\%$$

① Logistic Regression in Python -

Which of these methods is used for fitting a logistic regression model using statsmodels

- GLM() ∵ feedback: (Correct) the GLM() method is used to fit a logistic regression model using statsmodels

Log Odds :

Suppose you are working for a media service company like Netflix. They're launching a new show called 'Sacred Games' & you are building a logistic regression model which will predict whether a person will like it or not based on whether consumers have liked/disliked some previous shows. You have the data of five of the previous shows & you're just using the dummy variables for these five shows to build the model. If the variable is 1, it means that the consumer liked the show & if the variable is zero it means that the consumer didn't like the show. The following table shows the value of coefficients for these five shows that you got after building the logistic regression model.

Variable Name Coefficient Value

TrueDetective_Liked 0.47

ModernFamily_Liked -0.45

Mindhunter_Liked 0.39

Friends_Liked -0.23

Narcos_Liked 0.55

Now you have the data of three consumers Reetesh, Kshitij & Shanti for these 5 shows indicating whether or not they liked these shows. This is shown in the table below:

Consumer	True Detective_Liked	Modern Family_Liked	Mindhunter_Liked	Friends_Liked	Narcos_Liked
Reetesh	1	0	0	0	1
Kshitij	1	1	1	0	1
Shanti	0	1	0	1	1

Based on this data, which one of these three consumers is most likely to like to new show 'Sacred Games'?

Reetesh Feedback: Correct

To find the person who is most likely to like the show, you can use log odds. Recall the log odds is given by -

$$\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \dots + \beta_n x_n$$

Here, there are five variable for which the coefficients are given. Hence, the log odds become:

$$\ln\left(\frac{P}{1-P}\right) = 0.47x_1 - 0.45x_2 + 0.39x_3 - 0.23x_4 + 0.55x_5$$

As you can see, we have ignored the β_0 since it will be the same for all the three consumers. Now using the values of the 5 variables given, you get -

$$\begin{aligned} (\text{Log Odds})_{\text{Reetesh}} &= (0.47 \times 1) - (0.45 \times 0) + (0.39 \times 0) - (0.23 \times 0) + \\ &\quad (0.55 \times 1) \\ &= 1.02 \end{aligned}$$

$$\begin{aligned} (\text{Log Odds})_{\text{Kshitij}} &= (0.47 \times 1) - (0.45 \times 1) + (0.39 \times 1) - (0.23 \times 0) + \\ &\quad (0.55 \times 1) \\ &= 0.96 \end{aligned}$$

$$\begin{aligned} (\text{Log Odds})_{\text{Shreya}} &= (0.47 \times 0) - (0.45 \times 1) + (0.39 \times 0) - (0.23 \times 1) + \\ &\quad (0.55 \times 1) \\ &= -0.13 \end{aligned}$$

Shot on OnePlus

By pramodkhandare

You can clearly see, the log Odds of Reetesh is the highest, hence, the odds of Reetesh liking the show is the highest & hence, he is most likely to like the new show, sacred games.

Sensitivity & Specificity -

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

Actual / Predicted		Not Churn	Churn
Not Churn		TN	FP
Churn		FN	TP

$$\text{False Positive Rate} = \frac{FP}{TN + FP} \Leftrightarrow \text{False Positive Rate (FPR)}$$

$$\text{Positive Predictive Value} = \frac{TP}{TP + FP} \Leftrightarrow \text{True Positive Rate (TPR)}$$

$$\text{Negative Predictive Value} = \frac{TN}{TN + FN} \Leftrightarrow$$

$$\text{True Positive Rate (TPR)} = \frac{TP}{TP + FN}$$

$$\therefore TPR = \text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\therefore FPR = 1 - \text{Specificity} = 1 - \frac{TN}{TN + FP}$$

$$= \frac{FP}{TN + FP}$$

~~True Positive Rate -~~

You have the following table showcasing the actual 'churn' labels & the predicted probabilities for 5 customers.

Customer	Churn	Predicted Churn Probability
Thulasi	1	0.52
Aditi	0	0.56
Jaldeep	1	0.78
Ashok	0	0.45
Ananya	0	0.22

Calculate the True Positive rate & False Positive rate for the cutoffs of 0.4 & 0.5. Which of these cutoffs, will give you a better model?

Note - The good model is the one in which TPR is high & FPR is low.

- Cutoff of 0.5

Feedback:

Now, at the cutoff of 0.4, you get the following values of the predicted probabilities -

Customer	Churn	Pred. Churn Prob.	Pred. Churn Label
Thulasi	1	0.52	1
Aditi	0	0.56	1
Jaldeep	1	0.78	1
Ashok	0	0.45	1
Ananya	0	0.22	0

From the above table, you can easily calculate:

$$\text{True Positives} = 2$$

$$\text{False Positives} = 2$$

Also, from the original table, you have:

$$\text{Actual Positives} = 2$$

$$\text{Actual Negatives} = 3$$

Hence you get:

$$\text{TPR} = \frac{\text{True Positives}}{\text{Total Actual Positives}} = \frac{1}{2} = 100\%$$

$$\text{FPR} = \frac{\text{False Positives}}{\text{Total Actual Negatives}} = \frac{2}{3} \approx 67\%$$

Performing similar steps for a cutoff of 0.5 will give you

Customer	Churn	Pred. Churn Prob.	Pred. Churn Label
Thulasi	1	0.52	1
Aditi	0	0.56	0
Faideep	1	0.78	1
Ashok	0	0.45	0
Ananya	0	0.22	0

From the above table, you can easily calculate:

$$\text{True Positive} = 2$$

$$\text{False Positive} = 1$$

Also, from the original table, you have

$$\text{Actual Positives} = 2$$

$$\text{Actual Negatives} = 3$$

Hence you get -

$$TPR = \frac{\text{True Positives}}{\text{Total Actual Positives}} = \frac{2}{2} = 100\%$$

$$FPR = \frac{\text{False Positives}}{\text{Total Actual Negatives}} = \frac{1}{3} = 33\%$$

As you can see, with both the cutoffs, the TPR is 100% but for the cutoff of 0.5 you have a lower value of FPR so clearly, a cutoff of 0.5 gives you a better model.

Please note that 0.5 just gives the better model among 0.4 & 0.5. It might be possible that there is a cutoff point which gives an even better model.

② TPR and FPR →
Fill in the blanks:

When the value of TPR increases, the value of FPR

— Increases. ∵ feedback - Correct! This can be clearly seen from the ROC curve as well, when the value of TPR (on the Y-axis) is increasing, the value of FPR (on the X-axis) also increases

③ Area Under the curve —

You have the following five AUC (Area Under the curve) for ROC's plotted for five different models, which of these model is the best?

Model	A	B	C	D	E	Aus - <u>B(0.82)</u>
AUC	0.69	<u>0.82</u>	0.79	0.66	0.56	

Choosing the optimal Cut-off -

Suppose you created a dataframe to find out the optimal cut-off point for a model you built. The dataframe looks like the following:

Threshold	Probability	Accuracy	Sensitivity	Specificity
0.0	0.0	0.21	1.00	0.00
0.1	0.1	0.39	0.96	0.22
0.2	0.2	0.56	0.88	0.49
0.3	0.3	0.59	0.81	0.53
0.4	0.4	0.62	0.78	0.63
0.5	0.5	0.74	0.73	0.74
0.6	0.6	0.81	0.64	0.79
0.7	0.7	0.78	0.42	0.83
0.8	0.8	0.63	0.21	0.92
0.9	0.9	0.56	0.03	0.98

Based on the table above, what will the approximate value of the optimal cut-off be?

-0.5

Correct! the optimal cut-off point exists where the values of accuracy, sensitivity, and specificity are fairly decent if almost equal. At the cut-off of 0.5, the metric value are 0.74, 0.73 & 0.74 respectively. This is the optimal value of threshold that you can have.

Shot on OnePlus

By pramodkhandare

Precision & Recall -

i | nautix.
A BNY MELLON COMPANY

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$(\text{Recall} = \text{Sensitivity} = TPR)$$

① Calculating Precision →

Calculate the precision value for the following model:

Actual / Predicted	Not Churn	Churn
Not Churn	400 (TN)	100 (FP)
Churn	50 (FN)	150 (TP)

- 60%

∴ Feedback : Correct!

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{150}{100 + 150} = \frac{150}{250} = \frac{3}{5}$$

$$= 0.6 \approx 60\%$$

② F1-Score -

There is a measure known as F1-score which essentially combines both precision & recall. It is basically the harmonic mean of precision & recall & this formula is given by:

$$F = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The F1-Score is useful when you want to look at the performance of precision & recall together. Calculate the F1-score for the model below:

Shot on OnePlus
By pramodkhandare

By pramodkhandare

Actual / Predicted	Not Churn	Churn
Not Churn	400	100
Churn	50	150

67 % ::: Feedback: Correct!

$$F1\text{-Score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{Precision} + \text{Recall}}$$

$$= 2 \times \frac{\left(\frac{150}{100+150}\right) \left(\frac{150}{150+50}\right)}{\left(\frac{150}{100+150}\right) + \left(\frac{150}{150+50}\right)}$$

$$(91) \quad \text{Ans} = 2 \times \frac{(15\phi/25\phi) \times (15\phi/20\phi)}{(15\phi/25\phi) + (15\phi/20\phi)}$$

$$= 2 \times \frac{(3/5) \times (3/4)}{(3/5) + (3/4)}$$

$$= 2 \times \frac{(0.6 \times 0.75)}{0.6 + 0.75}$$

$$2 \times \frac{0.45}{1.35}$$

Ex 0.333

$$= 2 \times 0.333$$

$$= 0.6666 \approx 66.67\%$$

≈ 67%

359	1130
382	250
	1540
559	1050
223	1450
582	1020
1050	
(1450)	
1050	
350	
1400	0.77
	0.75
	8.02

⑤ Evaluation Metrics

Consider the same model given in the last question.

Patient ID	Heart Disease	Pred. Prob for Heart D	Pred. Lab
1001	0	0.84	0
1002	0	0.58	1
1003	0	0.73	1
1004	0	0.68	1
1005	0	0.21	0
1006	0	0.04	0
1007	1	0.48	0
1008	0	0.64	0
1009	0	0.61	1
1010	1	0.86	1

Calculate the values of Accuracy, Sensitivity, Specificity & Precision.
Which of these four metrics is the highest for the model.

→ Sensitivity

$$\therefore \text{Feedback: Correct!}$$

TP = 4	Actual	Predicted
FN = 2	1 ↪ 0 ↪ FN	1 ↪ TP
FP = 2	0 ↪ 1 ↪ FP	0 ↪ 0 ↪ TN
TN = 3		

$$\text{Precision} = \frac{4}{6} = 0.66 \approx 67\%$$

$$\text{Sensitivity} = \frac{4}{6} = 0.8 \approx 80\%$$

$$\text{Accuracy} = \frac{7}{10} = 0.7 \approx 70\%$$

$$\text{Specificity} = \frac{3}{5} = 0.6 \approx 60\%$$

Naive BayesBayes' Theorem & Its Building Blocks

Graded Questions -

	Courses	DS	ML	DL	BD	AI	Total
Male	80	60	40	50	30	260	
Female	70	40	50	70	10	240	
Total	150	100	90	120	40	500	

① Given this Contingency table, what is the probability that a randomly selected person joined Data Science DS?

$$\therefore P(\text{Person who Joined DS}) = \frac{150}{500} = \underline{\underline{.3 \approx 30\%}}$$

② Given this Contingency table, what is the probability that randomly selected female joined DS? In other words, what is the probability of a person joining DS given that she is female

$$\therefore P(\text{DS} | \text{Female}) = \frac{70}{240} = \underline{\underline{.29166 \approx 29.16\%}}$$

③ Consider a set containing all DL students OR all male students what is the Probability that a randomly selected person will belong to this set?
Hint: Use the formula: $(A \text{ or } B) = P(A) + P(B) - P(A \cap B)$

$$\therefore \frac{310}{500}$$

: Feedback :

This question deals with a probability concept called 'OR'. There is a formula for OR, which is -

$$P(A \text{ OR } B) = P(A) + P(B) - P(A \text{ AND } B)$$

In this example, you're looking at two things: DL & Male. So the question asked is - $P(\text{DL OR Male}) =$

$$P(\text{DL}) + P(\text{Male}) - P(\text{DL AND Male})$$

Using table in the question description, you see that -

$$P(\text{DL}) = 90/500$$

$$P(\text{Male}) = 260/500$$

$$P(\text{DL and Male}) = 40/500$$

therefore $P(\text{DL OR Male}) = P(\text{DL}) + P(\text{Male}) - P(\text{DL and Male})$

$$= 90/500 + 260/500 - 40/500$$

$$= (90+260-40)/500$$

$$= 310/500$$

$$= .62$$

$$\approx 62\%$$

Now, why do you subtract the probability of (Male & DL)? The answer is that when you count all the males and then count all the people who joined DL, there is an overlap because some males joined DL. This means you counted them twice, and so you have to subtract the extra count.

Shot on OnePlus

By pramodkhandare

$$P(A \text{ AND } B) = P(A) * P(B)$$

Given Bayes - With one feature

$$P(C = \text{Edible} | x = \text{Convex}) = \frac{P(x = \text{Convex} | C = \text{Edible}) \times P(C = \text{Edible})}{P(x)}$$

$$P(C_i | x) = \frac{P(x | C_i) P(C_i)}{P(x)}$$

Probability Calculation \rightarrow The probability of a convex mushroom being edible,

$P(C = \text{edible} | x = \text{CONVEX})$ given $P(x = \text{CONVEX})$

$- P(x = \text{CONVEX} | C = \text{edible}) \cdot P(C = \text{edible}) / P(x = \text{CONVEX})$

Mushroom Dataset

Sr.No.	Type of Mushroom	Cap. Shape	Cap. Surface
1.	Poisonous	Convex	Scaly
2.	Edible	Convex	Scaly
3.	Poisonous	Convex	Smooth
4.	Edible	Convex	Smooth
5.	Edible	Convex	Fibrous
6.	Poisonous	Convex	Scaly
7.	Edible	Bell	Scaly
8.	Edible	Bell	Scaly
9.	Edible	Convex	Scaly
10.	Poisonous	Convex	Scaly
11.	Edible	Flat	Scaly
12.	Edible	Bell	Smooth

(1) Probability Calculation -

Now let's say you picked a new mushroom whose cap-shape is CONVEX. What are the chances of this happening, i.e. What is the value of $P(X = \text{CONVEX})$?

- $8/12$

$$\because P(X = \text{CONVEX}) = 8/12$$

(2) $P(X = \text{CONVEX} | C = \text{edible})$?

- $4/8$

$\because P(X = \text{CONVEX} | C = \text{edible})$ means out of all the edible mushrooms, how many are CONVEX. Out of total 8 edible mushrooms, 4 are CONVEX. Thus, it is $4/8$.

(3) Probability Calculation -

In the previous questions, you have calculated that $P(C = \text{edible})$ is $8/12$, $P(X = \text{convex})$ is $8/12$ & $P(X = \text{CONVEX} | C = \text{edible})$ is $4/8$.

What is the probability that the convex mushroom is edible, $P(C = \text{edible} | X = \text{CONVEX})$?

- $4/8$

\therefore Feedback -

$$P(X = \text{CONVEX} | C = \text{edible}) \cdot P(C = \text{edible}) / P(X = \text{CONVEX})$$

$$= \left(\frac{4}{8}\right) \times \left(\frac{8}{12}\right) \times \frac{8}{12}$$

$$= \frac{4}{12} \times \frac{12}{6}$$

$$= \underline{\underline{4/8}}$$

Probability Calculation -

In the previous question, you find the probability of the convex mushroom being edible. What is the probability of the CONVEX mushroom being poisonous, $P(C = \text{poisonous} | X = \text{CONVEX})$?

4/8 ∵ Feedback: Since a mushroom can either be edible or poisonous, $P(C = \text{poisonous} | X = \text{convex})$ is $1 - P(C = \text{edible} | X = \text{CONVEX})$

$$= 1 - \frac{4}{8} = \frac{4}{8}$$

⑤ $P(C = \text{poisonous}) \Rightarrow 4/12$

⑥ Probability calculation -

What are the chances of a mushroom being CONVEX given it is poisonous. i.e. $(P(X = \text{CONVEX} | C = \text{poisonous}))$?

Feedback - $P(X = \text{CONVEX} | C = \text{poisonous})$ means out of all the poisonous mushroom, how many are CONVEX.

Out of the total 4 poisonous mushroom, all the 4 are CONVEX. Thus it is $\frac{4}{4} = 1$.

$$P(C = \text{edible} | X = \text{CONVEX}) =$$

$$\frac{P(X = \text{CONVEX} | C = \text{edible})}{P(C = \text{edible}) / P(X = \text{CONVEX})}$$

Refer Mushroom dataset :-

① Calculating conditional probability -

Say you take a new mushroom which is (CONVEX, SMOOTH). What is the numerator of $P(C = \text{edible} | X = \text{CONVEX, SMOOTH})$

$$\rightarrow P(\text{edible}) \times P(\text{CONVEX} | \text{edible}) \times P(\text{SMOOTH} | \text{edible})$$

② what is $P(\text{CONVEX} | \text{edible}) \Rightarrow 4/8$

③ what is $P(\text{SMOOTH} | \text{edible}) \Rightarrow 2/8$

④ What is $P(\text{CONVEX} | \text{Poisonous}) \Rightarrow 1 \quad (\because 4/4 = 1)$

⑤ What is $P(\text{SMOOTH} | \text{poisonous}) \Rightarrow 1/4$

⑥ $P(\text{CONVEX} | \text{edible}) = 4/8$
 $P(\text{SMOOTH} | \text{edible}) = 2/8$
 $P(\text{CONVEX} | \text{poisonous}) = 1$
 $P(\text{SMOOTH} | \text{poisonous}) = 1/4$

If all mushroom above 50% probability of being edible are classified as edible, is the CONVEX, SMOOTH mushroom edible.

\Rightarrow Cannot be decided, it is a tie

$$P(\text{edible} | \text{CONVEX, SMOOTH}) = \frac{P(\text{edible}) \cdot P(\text{CONVEX} | \text{edible}) \cdot P(\text{SMOOTH} | \text{edible})}{\text{denominator}}$$

$$= (8/12) (4/8) (2/8) / d = 1/12d$$

$$P(\text{poisonous} | \text{CONVEX, SMOOTH}) = \frac{P(\text{poisonous}) \cdot P(\text{CONVEX} | \text{poisonous}) \cdot P(\text{SMOOTH} | \text{poisonous})}{\text{denominator}}$$

$$= \frac{1}{12} \cdot \frac{1}{2} \cdot \frac{1}{4} / d = \frac{1}{12d}$$

The both numerators are equal to $1/12d$, this mushroom cannot be classified with a 50% threshold. Although if you would take a higher threshold, like 60% (which is reasonable since you don't want to take responsibility of people eating poisonous mushrooms), then it will be classified as poisonous. why? Because, when you set the threshold as 60%, you want the probability of edible | CONVEX, SMOOTH to at least 60%.

Prior, Posterior and Likelihood -

- Prior Probability - $P(\text{class} = \text{edible})$ or $P(\text{class} = \text{poisonous})$
- Likelihood - $P(X|\text{class})$
- $P(\text{class} = \text{edible} | X) = \text{Posterior Probability}$

Refer Mushroom Dataset -

② Likelihood Calculation -

Say you consider a (CONVEX, SCALY) mushroom. The likelihood is higher for it being -

- Poisonous :: Feedback -

$$\text{Likelihood} = P(X = (\text{CONVEX}, \text{SCALY}) | \text{class}) \cdot \text{class} = \text{edible}$$

$$= P(\text{CONVEX} | \text{Edible}) \cdot P(\text{SCALY} | \text{Edible}) = \frac{4}{8} \cdot \frac{5}{8} = \frac{20}{64} = \underline{\underline{31.25\%}}$$

$$\text{class} = \text{poisonous} : P(\text{CONVEX} | \text{poisonous}) \cdot P(\text{SCALY} | \text{poisonous})$$

$$= \frac{4}{4} \cdot \frac{3}{4} = \underline{\underline{75\%}}$$

(3) Likelihood Calculation -

The value of $P(X|Class) \cdot P(Class)$ where $X = (\text{CONVEX}, \text{SCALY})$ for both classes (edible and poisonous) are respectively:

$$\text{Edible} = 20.8\% ; \text{Poisonous} = 25.0\%$$

$$\therefore \text{Feedback} - P(\text{CONVEX}|\text{Edible}) \cdot P(\text{SCALY}|\text{EDIBLE}) \cdot P(\text{Edible}) \\ = (4/8)(5/8)(8/12) = 20.8\% ;$$

$$\text{Poisonous: } P(\text{CONVEX}|\text{poisonous}) \cdot P(\text{SCALY}|\text{poisonous}) \cdot P(\text{poisonous}) \\ = (4/4)(3/4)(4/12) = 25\%$$

(4) For the (**CONVEX, SCALY**) mushroom:

- The prior is in favour of edible, posterior in favour of poisonous.

$$\therefore \text{prior} = 8/12 \text{ & } 4/12 \text{ for edible & poisonous respectively;} \\ \text{posterior is } 20.8\% \text{ & } 25\%.$$

(1) Prior Probability - (Refer Dataset on next page)
What is the prior probability of a mail being spam,
 $P(\text{class} = \text{spam})$

- $7/15$ - \therefore Feedback: There are 7 spam mails in the data set.

(2) Naive Bayes Assumption -

What does Naive Bayes assume while classifying spam or ham mails?

That frequency of keywords like hurry, free, offer etc are conditionally independent of each other.

③ Likelihood Calculation -

Consider an email with the vector of features.

$X = (\text{free}, \text{data}, \text{weekend}, \text{click})$: what is the likelihood, $P(X|\text{spam})$?

$$\begin{aligned}
 & \text{Given feedback} \\
 P(X|\text{spam}) &= P(\text{free}|\text{spam}) \cdot P(\text{data}|\text{spam}) \cdot P(\text{weekend}|\text{spam}) \cdot P(\text{click}|\text{spam}) \\
 &= (2/7) (1/7) (1/7) (2/7) \\
 &= \frac{4}{2401}
 \end{aligned}$$

④ Likelihood Calculation + slides of moving on to next self features $X = (\text{free}, \text{data}, \text{weekend}, \text{click})$. what is the likelihood, $P(X|\text{ham})$?

$$\frac{2}{4096} \quad \because \text{Feedback}$$

$$\begin{aligned}
 P(X|\text{ham}) &= P(\text{free}|\text{ham}) \cdot P(\text{data}|\text{ham}) \cdot P(\text{weekend}|\text{ham}) \cdot P(\text{click}|\text{ham}) \\
 &= (1/8) (2/8) (1/8) (1/8) \\
 &= \frac{2}{4096}
 \end{aligned}$$

Dataset -

Sl. NO.	Class	Freq 1	Freq 2	Freq 3	Freq 4
1	Spam	free	buy	limited	hurry
2	Ham	reply	data	report	presentation
3	Ham	report	preach	file	end of day
4	Spam	limited	file	buy	click
5	Ham	meeting	timelines	limited	documents
6	Spam	hurry	data	buy	block
7	Spam	limited	sex	click	Viagra
8	Ham	presentn	end of day	data	report
9	Ham	reply	data	presum	click
10	Spam	free	reply	weekend	click
11	Spam	limited	click	free	hurry
12	Ham	meeting	end of day	weekend	data
13	Spam	hurry	weekend	block	offer
14	Ham	report	presentn	file	end of day
15	Ham	free	timelines	reply	offer

⑤ Calculate conditional probability
the value of $P(x|class) \cdot P(class)$ for class = spam for
 $x = (\text{free}, \text{data}, \text{weekend}, \text{click})$?

$\rightarrow (4/2401) (7/15)$

∴ feedback - $P(\text{class} = \text{spam} | x) = P(\text{class} = \text{spam}) \cdot P(x | \text{class} = \text{spam}) = (7/15) (4/2401)$

⑥ Posterior Probability -
What is the posterior for class = Ham (i.e. without division by denominator) for the feature vector $x = (\text{free}, \text{data}, \text{weekend}, \text{click})$?

$$(2/4096) (8/15) \quad \because \text{feedback}$$

$$\begin{aligned} P(\text{class} = \text{ham} | X) &= P(\text{class} = \text{ham}) \cdot P(X | \text{class} = \text{ham}) \\ &= (8/15)(2/4096) \end{aligned}$$

Learning from Model Selection:

- * Central issue in Machine learning can be said to be the study of - How to extrapolate learnings from a finite amount of data to explain or predict all possible inputs of the same type
- * Occam's Razor - A model should be simple as possible but robust
- * Regression - $Y = X + \beta w$
- * Model & Learning Algorithm - The learning algorithm is asked what needs to be done ; it figures out how it needs to be done & returns a model.

Finally, you learnt about the following four unique points regarding the usage of a simpler model whenever possible:-

- ① A simple model is usually more generic than a complex model. This becomes important because generic models are bound to perform better on unseen data sets.
- ② A simpler model requires fewer training data points. This becomes extremely important because in many cases, one has to work with limited data points.
- ③ A simple model is more robust and does not change significantly if the training data points undergo small changes.

(4) A simple model makes more errors in the training phase but it often may perform complex models when it views new data. This happens because of overfitting.

* Simplicity and Complexity -

Why are simpler models considered to be better than complex models? (Note: More than one option may be correct.)

- ⇒ Simpler models are generic i.e., they apply to a wider range of data
- ⇒ Complex models make assumptions about the data, which are likely to be wrong
- ⇒ Simpler models require less training data compared with complex
- ⇒ Simpler models are more robust

* Overfitting - Is an extreme case of overfitting?

⇒ The first person has mugged up all the possible questions from numerous textbooks & preparation material.

* Disadvantage of complexity -

- ⇒ They will need more training data to 'learn'
- ⇒ Despite the training, they may not learn and perform poorly in the real world

* Overfitting - Possibility of overfitting exist primarily?

⇒ Models are trained on a set of training data, but their efficacy is determined by their ability to perform well on unseen (test) data.

* Overfitting in Linear Regression - clear sign of overfitting in Linear Regression

The R-squared value is 0.90 and 0.30 on train and test data, respectively.

Bias & Variance

High Variance - Unstable & sensitive to changes in the training data

High Bias - Quantifies how accurate the model is likely to be on future (test) data. Extremely simple models are likely to fail in predicting complex real-world phenomena. Simplicity has its own disadvantage

* In practice, however, we often cannot have a model with a low bias and a low variance. As the model complexity increases, the bias reduces, whereas the variance increases & hence, the trade-off

① Bias & Variance

→ Student 1 - high Variance, Student 2 - high bias

→ student 1 - low bias, student 2 - low variance

② Model Variance - Regression

→ Model 2 - ∵ since it can change its coefficients (& the constant term) to fit the new training data

③ → straight line > Degree - 15 > Polynomial

⇒ ∵ feedback - The bias is high when the model is highly simple.

④ Variance — Why is the variance in the higher degree polynomial said to be higher than the other two models?
→ The model will change drastically from its current state if the current training data is altered.

Feedback: variance refers to changes in the model as a whole when trained on a different data set. Since the polynomial is trying to overfit the data, it will change drastically with respect to it.

⑤ Regularization — When is regularization typically performed?

→ While the learning algorithm uses the training data to produce a model.

⑥ Cross Validation — Why is it often not possible to use the validation set approach?

→ the data available for training & testing is limited.

* In K-Fold Cross Validation, you divide the training data into K-Groups of samples. If $K=4$ (say), you use $K-1$ folds to build the model and test the model on the K^{th} fold.

* Hyperparameters are used to "fine-tune" or regularize the model to keep it optimally complex.

* the learning algorithm is given the hyperparameters as the input, & it returns the model parameters as the output.

* Hyperparameters are not part of the final model output.

Shot on OnePlus

By pramodkhandare

① Low Variance -

weak learner \therefore weak learner creates simpler models that have a lower variance. They are not able to model complex relationships and, hence, create a more generic model.

② Bias - Variance -

- Linear regression will have a high bias and a low variance.
- A polynomial equation of degree 4 will have a low bias & a low variance

① Model variance -

Measure variance of a model

- By measuring how much does the estimates of the model change on the test data on changing the training data

② Regularization -

- It is a technique that is used to strike a balance between model complexity and model accuracy on training data

\because Regn doesn't improve accuracy; it improves the balance b/w accuracy and complexity.

③ Simple Models -

- Simpler Models will always have fewer test errors than a complex model

④ Simplicity of a Model -

- By the no. of features used in the model
- By the no. of nodes & depth of trees in case of a tree model.

⑤ K-FOLD Cross Validation -

- As K increases, the training time for k-fold cross validation increases
- With higher number of folds, the estimated error, on an average, is usually lower.
- You repeat the cross validation process 'K' times
- Each 'Kth' fold is used as the validation data once.
- A model trained with k-fold cross validation will overfit.