

Applied probability and statistics

Assignment - 2

Name: Nikhil S61
Reg: 241058024
Big Data Analytics

1) a) How often does she forecast rain?

$$\text{forecast rain} = P(\text{forecast rain} \cap \text{Actual rain}) + P(\text{forecast rain} \cap \text{no rain})$$

$$= 0.4 + 0.2 \Rightarrow \boxed{0.6}$$

no of times she predicted that forecast rain.

b) How often does she make a mistake?

no of times she predicted forecast and it is not a match with the actual truth.

$$= P(\text{forecast rain} \cap \text{no rain}) + P(\text{forecast no rain} \cap \text{Actual rain})$$

$$= 0.2 + 0.15 \Rightarrow \boxed{0.35}$$

c) Given that she just forecast rain, what is the chance that it will actually rain? (conditional probability):

$$P(\text{Actually rain} / \text{forecast rain}) = \frac{P(\text{Actually rain} \cap \text{forecast rain})}{P(\text{forecast rain})}$$

$$= \frac{0.4}{0.6} = \boxed{0.67}$$

d) Given that it rains today, what is the probability that she had forecast rain in last night's broadcast?

$$P(\text{forecast rain} / \text{actual rain}) = \frac{P(\text{forecast rain} \cap \text{actually rain})}{P(\text{actually rain})}$$

$$P(\text{actually rain}) = P(\text{forecast rain} \cap \text{Actual rain}) + P(\text{forecast no rain} \cap \text{Actual rain})$$

$$P(\text{actually rain}) = 0.4 + 0.15$$

$$= 0.55$$

$$\Rightarrow \frac{0.4}{0.55} \Rightarrow \boxed{0.72}$$

2) F \Rightarrow 52% M \Rightarrow 48% CS \Rightarrow 5% finc \Rightarrow 2%

a) The student is female given that the student is majoring in computer science

$$P(F) = 52\% \Rightarrow 0.52 \quad P(CS) = 5\% \Rightarrow 0.05 \quad P(F \cap CS) = 2\% = 0.02$$

$$P(F/CS) = \frac{P(CS/F) \cdot P(F)}{P(CS)} \Rightarrow P(CS/F) = \frac{P(CS \cap F)}{P(F)} = \frac{0.02}{0.52}$$

$$P(CS/F) = \boxed{0.038}$$

$$P(S/CS) = \frac{P(CS/S) \times P(S)}{P(CS)} \Rightarrow \frac{0.038 \times 0.52}{0.05} = 0.395 \approx 0.40$$

b) This student is majoring in computer science given that the student is female.

$P(CS/W) = 0.038\%$ calculated in the above problem.

$$3) P(BR) = 20\% \Rightarrow 0.20 \quad P(AR) = 40\% \Rightarrow 0.40 \quad P(GR) = 40\% \Rightarrow 0.40$$

$$P(mh/BR) = 25\% \Rightarrow 0.25 \quad P(mh/AR) = 20\% \Rightarrow 0.20 \quad P(mh/GR) = 10\% \Rightarrow 0.10$$

$$P(mh/Risk) = 1 - P(Mh/Risk) \Rightarrow P(mh/BR) + P(mh/AR) + P(mh/GR)$$

$$\Rightarrow \frac{P(mh \cap Risk)}{P(Risk)} = (0.25 \times 0.2) + \frac{P(BR)}{P(Risk)} (0.4 \times 0.2) + \frac{P(AR)}{P(Risk)} (0.4 \times 0.1)$$

$$\Rightarrow 1 - 0.17 = \boxed{0.83}$$

$$4) P(S) = 0.002 \quad P(e) = 0.002 \quad P(CS) = 0.01 \quad P(he) = 0.001$$

$$P(CA/S) = 0.25 \quad P(W/e) = 0.30 \quad P(W/CS) = 0.90 \quad P(W/he) = 0.10$$

$$\cancel{P(CA/S)} \times P(CS/CA) = \frac{P(W/CS) \times P(CS)}{P(W)} \Rightarrow \frac{0.90 \times 0.01}{P(W)}$$

law of total probability

$$P(W) = (P(S) \times P(W/S)) + (P(e) \times P(W/e)) + (P(CS) \times P(W/CS)) + (P(he) \times P(W/he))$$

$$(0.002 \times 0.25) + (0.002 \times 0.30) + (0.01 \times 0.90) + (0.001 \times 0.10)$$

$$0.0005 + 0.0006 + 0.009 + 0.0001$$

$$= P(W) = 0.0102 \Rightarrow \frac{0.90 \times 0.01}{0.0102} = \boxed{0.882}$$

$$5) P(L_1) = 0.80 \quad P(L_2) = 0.20$$

$$P(W/L_1) = 0.2 \quad P(W/L_2) = 0.9$$

$$i) P(L_1/W) = \frac{P(W/L_1) \times P(L_1)}{P(W)}$$

$$P(W) = P(W/L_1) \times P(L_1) + P(W/L_2) \times P(L_2)$$

$$= 0.8 \times 0.2 + 0.2 \times 0.9$$

$$P(W) = 0.16 + 0.18 = 0.34$$

$$P(L_1/W) = \frac{P(W/L_1) \times P(L_1)}{P(W)} = \frac{0.2 \times 0.8}{0.34} = 0.470.$$

$$P(L_2/W) = \frac{P(W/L_2) \times P(L_2)}{P(W)} = \frac{0.9 \times 0.2}{0.34} = 0.529$$

there is a 47% that the robot is in location L_1 and there is a 53% that the robot is in location L_2 .

6) a) Random donor and Random receiver.

Select the donor and check for all recipient that can take it:
for all the donors.

$P(\text{donor (All) and recipient (all who can take)})$

mutual exclusive events and donor and recipient are independent. so that add all the donors and recipients.

$$P(O^- \cap \text{who all can take}) + P(O^+ \cap \text{who all can take}) + P(A^- \cap \text{who all can take}) + P(A^+ \cap \text{who all can take}) \\ + P(B^- \cap \text{who all can take}) + P(B^+ \cap \text{who all can take}) + P(AB^- \cap \text{who all can take}) + P(AB^+ \cap \text{who all can take})$$

* $O^- \cap$ who all can take

$$0.066 (0.066 + 0.374 + 0.063 + 0.357 + 0.015 + 0.085 + 0.006 + 0.034) = 0.066$$

* $O^+ \cap$ who all can take.

$$0.374 (0.374 + 0.357 + 0.085 + 0.034) = 0.3179$$

* $A^- \cap$ who all can take.

$$0.063 (0.063 + 0.357 + 0.006 + 0.034) = 0.0289$$

* $A^+ \cap$ who all can take

$$0.357 (0.357 + 0.034) = 0.1395.$$

* $B^- \cap$ who all can take.

$$0.015 (0.015 + 0.085 + 0.006 + 0.034) = 0.0021$$

* $B^+ \cap$ who all can take

$$0.085 (0.085 + 0.034) = 0.0101$$

* $AB^- \cap$ who all can take

$$0.006 (0.006 + 0.034) = 0.0002$$

* $AB^+ \cap$ who all can take

$$0.034 (0.034) = 0.0011$$

$$P(\text{Donor (All) and (who all) can take}) = 0.0667 + 0.3179 + 0.0289 + 0.1395 + 0.0021 + 0.0101 + 0.0002 + 0.0011 \approx 0.56\%$$

b) O^- is the universal donor the blood drives should focus more on O^- blood group. the transfusion policies is prioritize in keeping stocks of O^- blood for emergencies. blood transfusion priority should be according to the. heavier blood types. like AB^- and look for other matching or surviving groups here if the blood is in stock of O^- if the recipient need O^- it should be prioritized first.

c) The given case where we have limited time to transfusion of blood. because of the limited time and only A^+ blood type is in availability. it is best to do the blood sampling for the wounded soldier. here we can get out any blood group. but if the blood type is A^+ or AB^+ of the wounded. we can immediate start transfusion if the blood group comes out to be of different blood type. we are helpless in this situations. and wait for any miracle to happen.

7) In a future society we assume that the law enforcement is strong and even if a crime is committed it can be easily or can be detected early after crime commitment. now in this case to ensure no innocent people are imprisoned. we should make the False positive near to zero so that any Innocent people should be not test/classified as positive (a crime maker) by the machine, so the false positive is low. so does the ~~Recall, TPR~~ sensitivity. is high for the machine. precision

8) Recall is more relevant performance metric. here because the false negative should be less so that the Recall is high (If false negative is high we see that model is predicting a spam mail as not spam and our data may be lost) so $FN \downarrow$ Recall \uparrow (TPR, sensitivity) \uparrow .

increasing the classification threshold generally decreases FP

- c) when the classification threshold increases, precision definitely increases.
 - d) when the classification threshold increases, quantity of TP (True positive) decreases.
 - e) when the classification threshold is decreased, the quantities TN and FN both non uniformly decrease.
 - f) decreasing the classification threshold generally decreases FN.
 - g) when the classification threshold is decreased, recall definitely increases.
 - h) when the classification threshold is decreased the quantities FP and FN both non uniformly increase.
- 9) a) An expensive and critical hydroelectric turbine operates 23 hours a day. An ML model evaluates vibration patterns and predicts when the turbine is operating without anomaly with an accuracy 99.99%. This ML model that its accuracy value high suggest that the model is doing good job, because ~~assuming~~ out of 1000 turbine the model is accurately predicting that 99.99 times that the turbine are working correctly it can not be less than 99.99% because if it less then this then there might be a damage in the system and may led to the crash and a huge property loss & life loss may happen so this should be high and here the model predicates it is a good model.
- 10) a) If model A has better precision and better recall than model B, then model A is probably better. as the F_1 score is high for model A so it is doing better.
- 11) An ROC curve is a plot of True positive rate vs False positive rate for different thresholds.
- 12) lowering the threshold classifies more items as positive, this increase in both True positive & False positive values

13) No change. AUC only comes about relative prediction probabilities.

14) It is a good model, because every time it predicts something the opposite happens. So we can think of the opposite decisions on prediction, we can help the model to go from zero to hero by telling the opposite or negation of the model as the model will mostly predict the exact wrong prediction.

15) Dashed black line represents random classification (True)

* ROC curve for any model can't fall below the dashed black line (False)

* The model represented by solid blue line is better than that represented by solid line (False)

16) TN \rightarrow True negative metric does not play a role in forming the precision-recall curve, because $\text{Recall} = \frac{TP}{TP+FN}$ and $\text{Precision} = \frac{TP}{TP+FP}$. So here we clearly see that the value of TN is not used to create this recall and precision metrics.

17) a) shows the overfit model.

b) shows the good fit model.

c) shows the underfit model.

18) a) metric = area under ROC curve because classifying the class as positive can buy, negative will not buy.

b) metric = precision-recall curves are more good with imbalanced datasets and focus mostly on positive classes.