

Assignment - 2

1. You have tracked the performance of the local meteorologist and compiled the following data.

$$P(\text{Forecast rain} \& \text{ actual rain}) = 0.4$$

$$P(\text{Forecast rain} \& \text{ no rain}) = 0.2$$

$$P(\text{Forecast no rain} \& \text{ actual rain}) = 0.15$$

$$P(\text{Forecast no rain} \& \text{ no rain}) = 0.25$$

- a) How often does she forecast rain?

$$\text{forecast rain} = P(\text{forecast rain} \& \text{ Actual rain}) +$$

$$P(\text{forecast rain} \& \text{ no rain})$$

$$= 0.4 + 0.2$$

$$= 0.6$$

- b) How often does she make a mistake?

mistake occurs when forecast doesn't match the actual weather

$$\text{mistake} = P(\text{for} \& \text{ NO rain}) + P(\text{g} \& \text{ no rain} \& \text{ Ar})$$

$$= 0.2 + 0.15$$

$$= 0.35$$

- c) Given that she just forecast rain, what is the chance that it will actually rain?

$$P(A_r | f_r) = ?$$

$$\frac{P(A_r \cap f_r)}{P(f_r)} = \frac{P(\text{actual rain} \& \text{ Forecast rain})}{P(\text{Forecast rain})}$$

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

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

$$P(F_r / A_r) = ? \quad P(F_r \cap A_r) = \frac{0.4}{0.55} = 0.72$$

$$\begin{aligned} P(\text{actual rain}) &= P(F_r + A_r) + P(\text{no rain} \cap A_r) \\ &= 0.4 + 0.15 \end{aligned}$$

$$P(A_r) = 0.55$$

2. Fifty-two percent of the students at a certain college are females. Five percent of the students in this college are majoring in Computer Science. Two percent of the students are women majoring in computer science. If a student is selected at random, find the conditional probability that:

a) the student is female given that the student is majoring in Computer Science

$$P(F) = 52\% = 0.52$$

$$P(CS) = 0.05$$

$$P(CS \cap F) = 2\% = 0.02$$

$$P(F / CS) = ?$$

Baye's theorem

$$P(CS \cap F) = \frac{P(CS \cap F)}{P(F)} = \frac{0.02}{0.52} = 0.038$$

$$\frac{P(F \cap CS)}{P(CS)} = \frac{0.038 \times 0.52}{0.05} = 0.3952 = 0.40$$

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

$$P(\text{cs} | w) = 0.038\% = \frac{P(\text{cs} \cap w)}{P(w)} = \frac{0.02}{0.52} = 0.038$$

3) Suppose that an insurance company classifies people who buy medical insurance from them into one of three classes: good, average & bad risks. As a data scientist for the company, you have access to the following customer data for the calendar year 2023-24.

Class	% of customers	% who had a major health problem
Bad risk	20%	25%
Average risk	40%	20%
Good risk	40%	10%

What is the probability that a new customer will not meet with a major health problem during 2024-25?

$$P(\text{Bad}) = 0.20$$

$$P(\text{mh} | \text{Bad}) = 0.25$$

$$P(\text{Avg}) = 0.40$$

$$P(\text{mh} | \text{Avg}) = 0.20$$

$$P(\text{Good}) = 0.40$$

$$P(\text{mh} | \text{Good}) = 0.10$$

$$= P(\text{Bad}) \times P(\text{mh} | \text{bad}) + P(\text{Avg}) \times P(\text{mh} | \text{Avg}) + P(\text{Good}) \times P(\text{mh} | \text{Good})$$

$$= 0.20 \times 0.25 + 0.40 \times 0.20 + 0.40 \times 0.10$$

$$= 0.05 + 0.08 + 0.04$$

$$= 0.17$$

$$= 1 - 0.17$$

$$P(\text{new customer not meet}) = 0.83$$

4. Suppose there are four failure modes are respectively 0.002, 0.002, 0.01 and 0.001. Given that it is a structural failure there is a 25% chance the plane will crash. the crash probabilities given the other three failure modes are 30%, 90% & 10% respectively for engine control system & human error failure mode. If a plane has crashed, what is the probability that it was due to a control system failure.

4 failure model for Single Engine plane

$$\text{Structural } P(s) = 0.002$$

$$P(c\sigma|s) = 25\% = 0.25$$

$$\text{engine } P(e) = 0.002$$

$$P(c\sigma|e) = 30\% = 0.30$$

$$\text{control System } P(cs) = 0.01$$

$$P(c\sigma|cs) = 90\% = 0.90$$

$$\text{human Error } P(he) = 0.001$$

$$P(c\sigma|he) = 10\% = 0.10$$

$$P(cs|c\sigma) = ? \quad \frac{P(c\sigma|cs) \cdot P(cs)}{P(c\sigma)} = \frac{0.90 \times 0.01}{0.0102} = 0.8823$$

Baye's theorem

$$\begin{aligned} P(c\sigma) &= P(s) \times P(c\sigma|s) + P(e) \times P(c\sigma|e) + P(cs) \times P(c\sigma|cs) + \\ &\quad P(he) \times P(c\sigma|he) \\ &= 0.002 \times 0.25 + 0.002 \times 0.30 + 0.01 \times 0.90 + \\ &\quad 0.001 \times 0.10 \\ &= 0.0005 + 0.0006 + 0.009 + 0.0001 \end{aligned}$$

$$\boxed{P(c\sigma) = 0.0102}$$

5. A robot, which only has a camera as a sensor, can either be in one of two locations: L_1 or L_2 . The robot doesn't know exactly where it is but based on all past observations, the robot thinks that there is an 80% chance that it is in L_1 and a 20% chance that it is in L_2 . Location L_2 is the only one that has a window. The robot's vision algorithm detects a window but its image recognition algorithm is not perfect; the probability of observing a window given there is no window at the location is 0.2 & the probability of observing a window given there is a window is 0.9. After incorporating the observation of a window, what is the robot's updated probability that it is in (1) L_1 , (2) L_2 ?

$$P(L_1) = 0.80 \quad P(w|L_1) = 0.2 \quad P(L_1|w) = ?$$

$$P(L_2) = 0.20 \quad P(w|L_2) = 0.9 \quad P(L_2|w) = ?$$

$$\begin{aligned} P(w) &= P(L_1) \times P(w|L_1) + P(L_2) \times P(w|L_2) \\ &= 0.8 \times 0.2 + 0.2 \times 0.9 \end{aligned}$$

$$= 0.16 + 0.18$$

$$P(w) = 0.34$$

$$P(L_1|w) = \frac{P(w|L_1) P(L_1)}{P(w)} = \frac{0.2 \times 0.8}{0.34} = 0.470$$

Baye's theorem

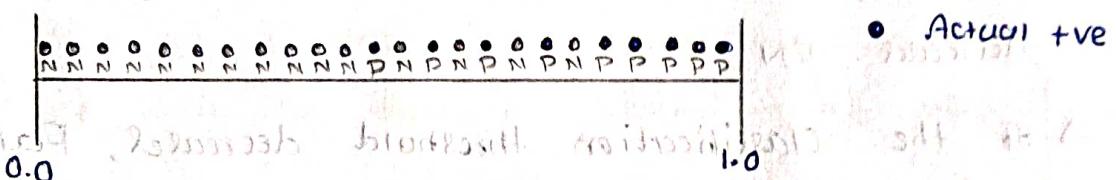
$$P(L_2|w) = \frac{P(w|L_2) P(L_2)}{P(w)} = \frac{0.9 \times 0.2}{0.34} = 0.529$$

Given that 52% chance will be there is a window in location L_2

7. In a future society, a machine is used to predict a crime before it occurs. If you were responsible for tuning this machine, what evaluation metric would you want to maximize to ensure no innocent people (people not about to commit a crime) are imprisoned.
- Here to ensure no innocent people are imprisoned we should make the False positive, ~~to zero so that any~~ thus the False positive low ~~definitely~~ definitely precision high for the machine.
8. Consider a classification model that separates email into two categories: "spam" or "not spam". Answer the following questions regarding precision & recall (a.k.a. sensitivity or TPR) by playing around with the threshold slider on the demo website here:
- a) which is more relevant performance metric in this case: recall or precision? Explain briefly why?
- False Positive:- It's not a spam email but is detected as a spam email
- False Negative:- It's a spam email but it's detected as a not a spam Email
- b) Increasing the classification threshold generally increase/decrease FP.
- FP is decreased
- When the classification threshold increased, precision
- If you increase the threshold, false positive decrease then automatically precision is definitely increased

- d) keeping in mind that $TP + FP + TN + FN = n$, which is the number of samples, when the classification threshold is increased, what happens to the quantity TP ?
→ when the classification threshold is increased, True positive decreased
- e) when the classification threshold is increased, the quantities TN & FN both increase uniformly.
→ Classification threshold is increased, the True Negative & False Negative both will increase uniformly.
- f) Decreasing the classification threshold generally increases FN .
- If the classification threshold decreased, False negative will be decreased.
- g) when the classification threshold is decreased, recall increases.
→ If the classification threshold is decreased, False negative also decreases then recall definitely increases.
- h) when the classification threshold is decreased, the quantities TP & FP both increase.
→ If the classification threshold decrease, the True positive & False positive both will increase.
10. consider two models: A and B, that evaluate the same dataset, which one of the following statements is true?
a) If model A has better precision & better recall than model B, then model A is probably better. (✓)

11. An ROC curve is a plot of True positive Rate vs False positive Rate for different thresholds
12. Lowering the classification threshold classifies more items as true / negative, thus increasing both True Positive & False Positive
13. AUC (Area under the Roc curve), one way of interpreting AUC is as the probability that the model ranks a random +ve Example more highly than a random -ve, ex. which are arranged from left to right in ascending order of prediction probabilities.

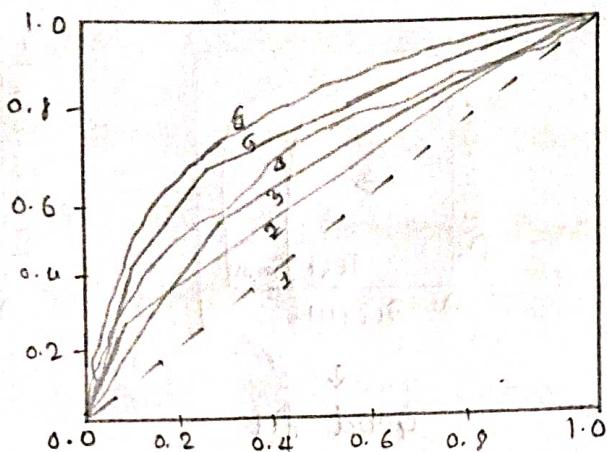


random positive (Blue) to right of a random -ve (Pencil) ex. AUC ranges from 0 to 1. If model predictions are 100% wrong, has an AUC of 0.0. If one whole predictions are 100% correct Auc of 1.0.

Suppose we multiplied all of the predictions from a given model by 0.5 (for ex, if the model predicts 0.4, we multiply by 0.5 to get a prediction of 0.2) how would it change the model's performance as measured by AUC?

→ No change. AUC only care about relative prediction probabilities.

15. The figure shows ROC curves for different models.



1. RF AUC = 0.684
2. DT AUC = 0.551
3. KNN AUC = 0.599
4. SVM AUC = 0.623
5. LR AUC = 0.718

• Dashed black line represents random classification \rightarrow True positive rate = False positive rate

• ROC curve for any model can't fall below the dashed black line \rightarrow False (Area must be high)

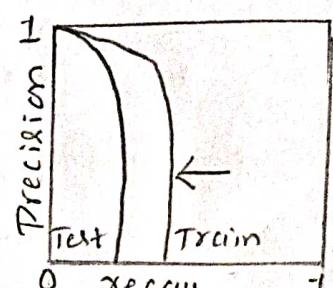
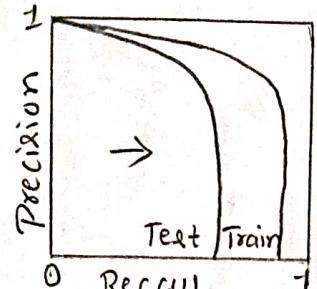
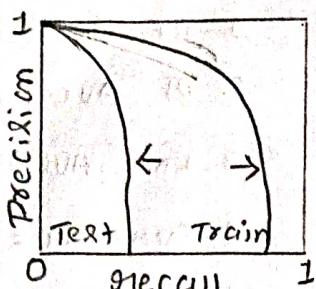
• The model represented by solid blue line is better than that represented by solid red line \rightarrow False \rightarrow Area must be high

16. Which one among TP, TN, FP, FN does not play a role in forming the precision-recall curve? What does the conclusion mean intuitively?

\rightarrow True Negative, for precision = $\frac{TP}{TP+FP}$, recall = $\frac{TP}{TP+FN}$

So if you plot a graph for precision & recall need only TP, FP, & FN there is no need of TN so that's why TN does not play a role in forming the precision-recall curve.

17. Identify which model overfit, underfit & which one is good fit.



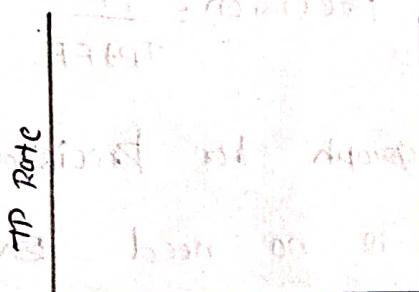
18. Explain which one among area under ROC and area under

Precision-recall curve would you use for the following scenario.

a) Identifying whether a customer will buy a product or discount or not when a customer is equally likely to do so \rightarrow ROC

b) Identifying a spam email when generally spam emails constitute 1% of the total emails \rightarrow Precision, Recall

14. Your friend shows you his model's ROC curve as follows



Is your friend's model any good, why? How can you help your friend's model go from zero to hero?

Yes it's good we can take decisions in the opposite way of the model predictions, this way friend model became zero to hero.

6. Suppose a person needs a blood transfusion. we show the compatibility chart x means compatibility
- a) what is the probability that a random person will be able to donate to another random person given no information about blood types of either the giver or the receiver?

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

mutual exclusive events and Donor and recipients independent so that add all the Donors and recipients

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

$$P(A^- \cap (A^-, A^+, AB^-, AB^+)) + P(A^+ \cap (A^+, AB^+)) + P(B^- \cap (B^-, B^+, AB^-, AB^+)) +$$

$$P(B^+ \cap (B^+, AB^+)) + P(AB^- \cap (AB^-, AB^+)) +$$

$$P(AB^+ \cap AB^+)$$

$$P(O^- \cap (\text{who all can take})) = 0.066 [0.066 + 0.374 + 0.063 + 0.357 + 0.015 + 0.085 + 0.006 + 0.034] = \underline{\underline{0.066}}$$

$$P(O^+ \cap (\text{who can take})) = 0.374 [0.374 + 0.357 + 0.085 + 0.034] = \underline{\underline{0.3179}}$$

$$P(A^- \cap (\text{who can take})) = 0.063 [0.063 + 0.357 + 0.006 + 0.034] = \underline{\underline{0.0289}}$$

$$P(A^+ \cap (\text{who can take})) = 0.357 [0.357 + 0.034] = \underline{\underline{0.1395}}$$

$$P(B^- \cap (\text{who can take})) = 0.015 [0.015 + 0.085 + 0.006 + 0.034] = \underline{\underline{0.0021}}$$

$$P(B^+ \cap (\text{who can take})) = 0.085 [0.085 + 0.034] = \underline{\underline{0.0101}}$$

$$P(AB^- \cap (\text{who can take})) = 0.006 [0.006 + 0.034] = 0.0002$$

$$P(AB^+ \cap (\text{who can take})) = 0.034 [0.034] = 0.0011$$

$$P(\text{Donor (A11)} \text{ and } (\text{who all can take})) =$$

$$0.066 + 0.3179 + 0.0284 + 0.1395 + 0.0021 + 0.0101 + 0.0002 \\ + 0.0011 = 0.56$$

$$P(\text{Donor (A11) & who can take}) = 56\%$$

b) Given the above data, what can you say about blood transfusion policy in a hospital regarding blood drives target & blood transfusion priority?

O negative is the universal donor the blood drives should focus more on O negative blood group the transfusion policies is prioritizing keeping stocks of O⁻ blood for emergencies. blood transfusion priority should be according to the rarest blood types like AB⁻ and look for other matching or receiving groups here if the blood is in stock of O negative if the recipient need of O negative should be prioritized first.

- c) In a battle field hospital a soldier is brought in for immediate blood transfusion. only blood type A+ is available in the supply for immediate use. we do not know the wounded soldier's blood type. there are two other soldiers present who are willing to donate their blood. we have time to do one blood typing before time becomes critical. what should we do?
- The given case where we have limited time to transfusion of blood. because of the limited time and only A+ve blood type is available in blood bank. so it is best to do the blood sampling for the wounded soldier here we ~~can~~ can get out of the blood group but if the blood type is A+ve or AB+ve of the wounded soldier we can immediate start transfusion, so if we got like this we have a two Probabilities. if the blood group comes out to be of different blood type we are helpless.
- d) In which of the following scenarios would a high accuracy value suggest that the ML model is doing a good job? Explain your answer briefly.
- An expensive and critical hydro-electric turbine operates 23 hours a day. An ML model evaluates vibration patterns and predicts when the turbine is operating without ~~any~~ anomaly with an accuracy 99.99%