**Section A)**

Q1)

(a)

A1)(a) More diversity means the ensemble trees would give different results & the a strong classifier could be made by combining all of them. otherwise we ~~wouldn't~~ would have ~~the~~ a single tree instead of multiple trees that produce the same results. However more diversity means lesser correlation b/w the trees which would make the model suffer from high variance. This is the tradeoff b/w correlation & diversity in RFs. The trees need to be correlated upto a certain extent which is achieved by selecting a subset of features for each model to train on & by creating bootstrap samples.

(b)

(b) When the no of features is relatively large compared to the no of data points the curse of dimensionality becomes a problem for Naive Bayes. We can reduce their dimensionality by doing feature selection or feature extraction through PCA or SVD. We can use TSNE for the same. Also, we can use K-fold cross-validation to tune the hyperparameters to to reduce overfitting of the sparse data present in higher dimensions or we can use ensemble methods too

(c)

(c) If some attributes are missing then it maybe that $P(Y=y|X=x_i)$ may be zero for x attribute. This will lead to its whole probability being zero (which shouldn't be the case). To mitigate this we can use Laplace smoothing (adding a constant value to the ~~numerator~~ & deno count of each attribute in training data). Example - Let's suppose we have the data of a football team that hasn't lost a match for the past 5 years. So, it doesn't mean that for the next year its probability of losing would be zero. However without Laplace smoothing the NB model would predict a probability zero.
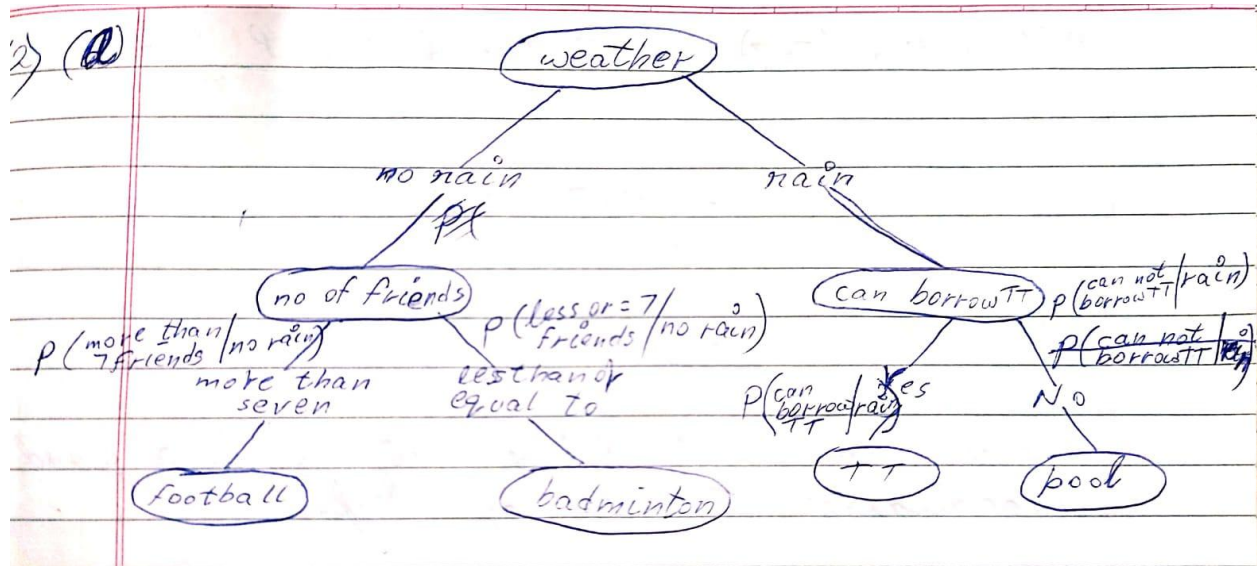
(d)

(d) Yes. If an attribute has higher cardinality then the information gain associated with that attribute would be higher as $P(X=x)$ would be higher. An alternative way could be to use the gini index (or ~~gain~~) (or gain Ratio). ~~gain~~ Example - Predicting whether a customer will ~~predict~~ tip the waiter based on 'meal type' & 'customer satisfaction' with 'meal type' being (Breakfast, Lunch, Dinner) & the other being ratings from 1-5. If we use IG then there may be bias towards 'customer satisfaction' since the cardinality is higher over here. However if we use Gain Ratio then it'll penalize the 'customer satisfaction' attribute thus providing a more balanced assessment

A2)
(a) TT here refers to Table Tennis.
Assumption -
P(no rain) = 0.5 and P(rain) = 0.5

## 2) (a)

weather
- no rain → no of friends
  - $P\left(\text{more than 7 friends} \mid \text{no rain}\right)$ more than seven → football
  - $P\left(\text{less or = 7 friends} \mid \text{no rain}\right)$ less than or equal to → badminton
- rain → Can borrow TT
  - $P\left(\text{can borrow TT} \mid \text{rain}\right)$
  - $P\left(\text{can not borrow TT} \mid \text{rain}\right)$
  - $P\left(\text{can borrow TT} \mid \text{rain}\right)$ yes → TT
  - $P\left(\text{can not borrow TT} \mid \text{rain}\right)$ No → pool

(b) Assumption: Rainy and Clear Weather are mutually exclusive. Thus 'not rainy' can be modeled as 'clear'.
So, P(~Prediction of rainy) = P(Prediction of Clear) and P(~Rainy) = P(Clear)

(b)    Let $P$ = prediction of the app that its rainy

     $R$ = It is Rainy

Thus,    $P(P/R) = 0.8$    — ①

      $P(\sim P/\sim R) = 0.9$ — ②   } given

      $P(P) = 0.3$

      $P(\sim P) = 0.7$

Hence, $P(R/P) = P(P/R) \times \dfrac{P(R)}{P(P)}$

~~$P(R) = P(R/R) \times P($~~

~~$P(R) = P(R,P) + P(R,\sim P)$~~

~~$= P(R/P) \times P(P) + P(R/\sim P) P(\sim P)$~~

Using the first two eq's –

     $P(P/R) = 0.8$

$\Rightarrow P(P,R)/P(R) = 0.8$

$\Rightarrow P(P,R) = 0.8 P(R)$    — ③

And,   $P(\sim P/\sim R) = 0.9$

    $\therefore P(P/\sim R) = 1 - 0.9 = 0.1$

$\Rightarrow P(P, \sim R)/P(\sim R) = 0.1$

$\Rightarrow P(P, \sim R) = 0.1 P(\sim R) = 0.1 - 0.1 P(R)$ — ④

Adding ③ & ④,

   $P(P, R) + P(P, \sim R) = 0.7 P(R) + 0.1$

We know that $P(P) = P(P,R) + P(P, \sim R)$

Thus,   $P(P) = 0.7 P(R) + 0.1$

     $\Rightarrow P(R) = \dfrac{2}{7}$

Thus, $P(R|P) = P(P|R) \dfrac{P(R)}{P(P)}$

$\Rightarrow P(R|P) = 0.8 \times \dfrac{2}{7 \times 0.3}$

$\Rightarrow \boxed{P(R|P) = \dfrac{16}{21}}$

**Alternate Solution considering that the probability of rainy weather and clear weather were given instead of probability of prediction of rainy and clear weather:**

(b) Let $P_r$ = Prediction of Rain through app
   & $R$ = Rainy Weather

Thus, $P(P_r | R) = 0.8$
   $P(\sim P_r | \sim R) = 0.9$
   $P(R) = 0.3$
   $P(\sim R) = 0.7$

Hence, a/q —

$P(R|P_r) = P(P_r|R) \times \dfrac{P(R)}{P(P_r)}$

$P(P_r) = P(P_r, R) + P(P_r, \sim R)$
   $= P(P_r|R)P(R) + P(P_r|\sim R) P(\sim R)$
   $= 0.8 \times 0.3 + (1 - 0.9)0.7$
   $= 0.24 + 0.07$
   $= 0.31$

$P(R|P_r) = \dfrac{0.8 \times 0.3}{0.31}$

$\boxed{P(R|P_r) = \dfrac{24}{31} = 0.774}$

(c)
Assumption: P(Good Mood) and P(Bad Mood) are same
Therefore, P(Good Mood) = 0.5 and P(Bad Mood) = 0.5

(c) $P(Gym \mid Good\ Mood) = 0.8$ — ①
$P(No\ Gym \mid Good\ Mood) = 0.2$ — ②
$P(Gym \mid Bad\ Mood) = 0.4$ — ③
$P(No\ Gym \mid Bad\ Mood) = 0.6$ — ④
$P(Cardio \mid Gym) = 0.5$ — ⑤
$P(Weight \mid Gym) = 0.5$ — ⑥

```
                  Mood
              /          \
          Good            Bad
         ①/ \②          ③/ \④
       Gym   No Gym    Gym    No gym
      ⑤/ \⑥          ⑤/ \⑥
  Cardio  Weight   Cardio  Weight
```

$$P(Cardio \mid Good\ Mood) = P(Cardio \mid Gym) \times$$
$$P(Gym \mid Good\ Mood)$$

$$= 0.5 \times 0.8 = 0.4$$

$$P(Cardio \mid Bad\ Mood) = P(Cardio \mid Gym) \times P\left(\begin{array}{c}Gym \mid \\ Bad\ Mood\end{array}\right)$$

$$= 0.5 \times 0.4 = 0.2$$

$$P(Cardio) = P(Cardio \mid Good\ Mood)\,P(Good\ Mood) +$$
$$P(Cardio \mid Bad\ Mood)\,P(Bad\ Mood)$$

$$= 0.4 \times 0.5 + 0.2 \times 0.5 = 0.3$$

$$P(Weights \mid Good\ Mood) = P(Weights \mid Gym) \times P\left(\begin{array}{c}Gym \mid \\ Good\end{array}\right)$$

$$= 0.5 \times 0.8 = 0.4$$

$$P(Weights \mid Bad\ Mood) = P(Weights \mid Gym) \times P\left(\begin{array}{c}Gym \mid Bad \\ Mood\end{array}\right)$$

$$= 0.5 \times 0.4 = 0.2$$

$$P(Weights) = P(Weights \mid Good\ Mood)\,P(Good\ Mood) +$$
$$P(Weights \mid Bad\ Mood)\,P(Bad\ Mood)$$

$$= 0.4 \times 0.5 + 0.2 \times 0.5 = 0.3$$

$$P(Gym) = P(Gym \mid Good\ Mood)\,P(Good\ Mood) + P\left(\begin{array}{c}Gym \mid Bad \\ Mood\end{array}\right)$$
$$P(Bad\ Mood)$$

$$= 0.8 \times 0.5 + 0.4 \times 0.5 = 0.6$$
$$P(No\ Gym) = 1 - 0.6 = 0.4$$

(d)

(d) $P(Good\ mood) = 0.6$

$P(not\ good\ mood) = 0.4$ $\left.\right\}$ given

$P(F = 7\ hours | Good\ mood) = 0.7$

$P(F = 7\ hours | Bad\ mood) = 0.45$

$P(Bad\ mood | F=7hours) = \dfrac{P(F=7hours|Bad\ mood) \times P(Bad\ mo}{P(F=7,}$

$P(F=7) = P(F=7, Good\ Mood) + P(F=7, Bad\ Mood)$

$\quad\quad = P(F=7 | Good\ Mood)\ P(Good\ Mood) +$

$\quad\quad\quad P(F=7 | Bad\ Mood)\ P(Bad\ Mood)$

$\quad\quad = 0.7 \times 0.6 + 0.45 \times 0.4$

$\quad\quad = 0.42 + 0.18$

$\quad\quad = 0.6$

Hence, $P(Bad\ Mood | F=7) = 0.45 \times \dfrac{0.4}{0.6} = 0.3$

$P(Good\ Mood | F=7) = P(F=7 | Good\ Mood)\ \dfrac{P(Good\ Mood)}{P(F=7)}$

$\quad\quad\quad = 0.7 \times \dfrac{0.6}{0.6}$

$\quad\quad\quad = 0.7$

Hence, the most likely outcome is a good mood after 7 hours of sleep.

Elaborating further on his decisions-

$P(cardio(F=7)) = P(cardio \mid good\ mood) P(good(F=7))$
$+ P(cardio \mid bad\ mood) P(bad\ mood(F=7))$
$= 0.5 \times 0.8 \times 0.7 + 0.5 \times 0.4 \times 0.3$

$P(cardio) = 0.34$

$P(cardio \mid Good) = 0.5 \times 0.8 \times 0.7 = 0.28$
$P(cardio \mid Bad) = 0.5 \times 0.4 \times 0.3 = 0.06$

$P(no\ gym) = P(no\ gym \mid good) \times P(good) + P(no\ gym \mid bad) P(bad)$
$= 0.2 \times 0.7 + 0.6 \times 0.3$
$= 0.32$

$P(weights \mid good\ mood) = 0.5 \times 0.8 \times 0.7 = 0.28$
$P(weights \mid bad\ mood) = 0.06$
$P(weights) = 0.34$

Hence, Rahul will go to the gym for cardio or weights if he has 7 hours of sleep.

# Section B

Q3)

*Decision Trees:-*

The best results obtained was with the Entropy criterion

```
Average accuracy of Decision Tree Classifier with Entropy as the splitting criterion:  0.7336666666666666
Average accuracy of Decision Tree Classifier with Gini as the splitting criterion:  0.7105000000000001
```

Through GridSearchCV -

```
Best accuracy:  0.8568840579710144
Best parameters:  {'max_depth': 3, 'max_features': 7, 'min_samples_split': 2}
```

*Random Forests:-*

```
Average accuracy of Random Forest Classifier with Entropy as the splitting criterion:  0.8003333333333333
Average accuracy of Random Forest Classifier with Gini as the splitting criterion:  0.7939999999999999
```

Through GridSearchCV -

```
Best parameters:  {'max_depth': 2, 'min_samples_split': 3, 'n_estimators': 250}
```

Random Forest Classifier with best parameters obtained using GridSearchCV -

```
Random Forest Classifier with best parameters-
Test Accuracy:  0.8333333333333334
Classification Report:
              precision    recall  f1-score   support

           0       0.84      0.84      0.84        32
           1       0.82      0.82      0.82        28

    accuracy                           0.83        60
   macro avg       0.83      0.83      0.83        60
weighted avg       0.83      0.83      0.83        60
```

# Section C

Q4)

Preprocessing steps that were implemented for the given dataset -
1) Checking for missing values
2) Merging the classes of output label column to make it into a binary classification problem
3) Encoding the categorical features.

The MyDecisionTree was implemented and evaluated. The following accuracy was obtained on 70:30 train-test split -

```
Testing Accuracy of MyDecisionTree: 0.9785714285714285
```

On comparing this with the sklearn's Decision Tree Classifer, the accuracy of it was -

```
Testing accuracy of sklearn DT: 0.9809523809523809
```

We can see that our implemented Decision Tree Class performs really good on the given dataset.