Name: Liangfan Pang        CWID:10453333

1.

1）Bias is the difference between the average prediction of our model and the correct value which we are trying to predict. Variance is the variability of model prediction for a given data point or a value which tells us spread of our data. If our model is too simple and has very few parameters then it may have high bias and low variance. On the other hand if our model has large number of parameters then it's going to have high variance and low bias. So we need to find the right/good balance without overfitting and underfitting the data. This tradeoff in complexity is a the bias-variance trade-off.

  Practical techniques to reduce bias:
    Increase hypothesis space (i.e., model complexity)
  Practical techniques to reduce variance:
    resampling (e.g., bagging models like random forest)
    bias of each tree is the same as full model but higher variance
    averaging many trees decreases variance without increasing bias
    theoretically the more trees, the less variance (if no computation limit)
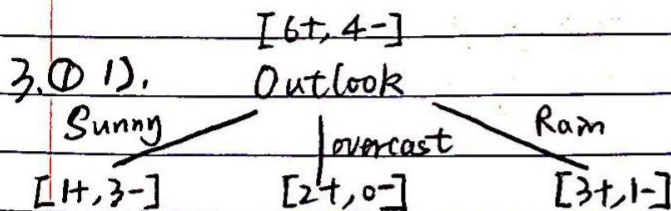
2）K-Fold cross-validation is where a given data set is split into a K number of sections/folds where each fold is used as a testing set at some point. Let's take the scenario of 5-Fold cross validation(K=5). Here, the data set is split into 5 folds. In the first iteration, the first fold is used to test the model and the rest are used to train the model. In the second iteration, 2nd fold is used as the testing set while the rest serve as the training set. This process is repeated until each fold of the 5 folds have been used as the testing set.

  As there is never enough data to train your model, removing a part of it for validation poses a problem of underfitting. By reducing the training data, we risk losing important patterns/ trends in data set, which in turn increases error induced by bias. So, what we require is a method that provides ample data for training the model and also leaves ample data for validation. K Fold cross validation does exactly that.

**2.** ① $P = \dfrac{TP}{TP+FP} = \dfrac{50}{50+40} = \dfrac{5}{9}$

② $r = \dfrac{TP}{TP+FN} = \dfrac{50}{50+30} = \dfrac{5}{8}$

③ $F_1\text{-score} = \dfrac{2pr}{p+r} = \dfrac{2}{\frac{1}{p}+\frac{1}{r}} = \dfrac{2}{\frac{9}{5}+\frac{8}{5}} = \dfrac{34}{5} = 6.8$

**3.** ① 1).

$$[6+, 4-]$$
$$\text{Outlook}$$

Sunny — | overcast — Rain

$[4+, 3-]$    $[2+, 0-]$    $[3+, 1-]$

$G(S, O) = E(S) - \dfrac{4}{10} \times E(Sunny) - \dfrac{2}{10} \times E(overcast) - \dfrac{4}{10} \times E(Rain)$
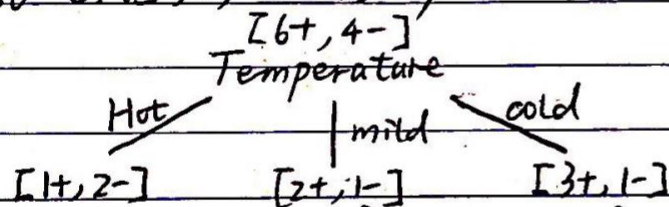
$E(S) = \dfrac{-6}{10} \times \log_2 \dfrac{6}{10} - \dfrac{4}{10} \times \log_2 \dfrac{4}{10}$

$E(Sunny) = \dfrac{-1}{4} \times \log_2 \dfrac{1}{4} - \dfrac{3}{4} \times \log_2 \dfrac{3}{4}$

$E(Overcast) = -1 \times \log_2 \dfrac{2}{2}$

$E(Rain) = \dfrac{-3}{4} \times \log_2 \dfrac{3}{4} - \dfrac{1}{4} \times \log_2 \dfrac{1}{4}$

So $G(S, O) = 0.3219$

$$[6+, 4-]$$
$$\text{Temperature}$$

Hot | mild   cold

$[1+, 2-]$    $[2+, 1-]$    $[3+, 1-]$
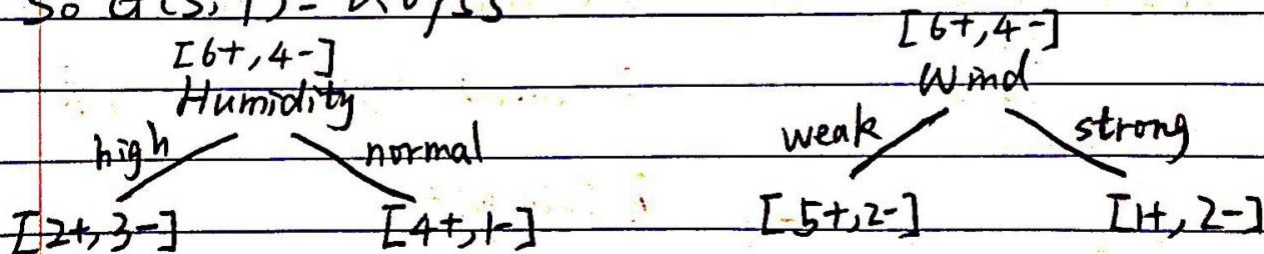
$G(S, T) = E(S) - \dfrac{3}{10} \times E(Hot) - \dfrac{3}{10} \times E(mild) - \dfrac{4}{10} \times E(cold)$

$E(hot) = \dfrac{-1}{3} \times \log_2 \dfrac{1}{3} - \dfrac{2}{3} \times \log_2 \dfrac{2}{3}$

$E(mild) = \dfrac{-2}{3} \times \log_2 \dfrac{2}{3} - \dfrac{1}{3} \times \log_2 \dfrac{1}{3}$

$E(cold) = \dfrac{-3}{4} \times \log_2 \dfrac{3}{4} - \dfrac{1}{4} \times \log_2 \dfrac{1}{4}$

So $G(S, T) = 0.0955$

$$[6+, 4-]$$
$$\text{Humidity}$$

high — | normal

$[2+, 3-]$    $[4+, 1-]$

$$[6+, 4-]$$
$$\text{Wind}$$

weak — strong

$[5+, 2-]$    $[1+, 2-]$

$G(S,H) = E(S) - \frac{1}{2} \times E(high) - \frac{1}{2} \times E(normal)$

$E(high) = \frac{-2}{5} \times \log_2 \frac{2}{5} - \frac{3}{5} \times \log_2 \frac{3}{5}$

$E(normal) = \frac{-4}{5} \times \log_2 \frac{4}{5} - \frac{1}{5} \times \log_2 \frac{1}{5}$

So $G(S,H) = 0.1245$

$G(S,W) = E(S) - \frac{7}{10} \times E(weak) - \frac{3}{10} \times E(strong)$

$E(weak) = \frac{-5}{7} \times \log_2 \frac{5}{7} - \frac{2}{7} \times \log_2 \frac{2}{7}$

$E(strong) = \frac{-1}{3} \times \log_2 \frac{1}{3} - \frac{2}{3} \times \log_2 \frac{2}{3}$

So $G(S,W) = 0.0913$

Since $G(S,O) > G(S,H) > G(S,T) > G(S,W)$

So the root node is Outlook.

2). $S_{Sunny} = \{D_1, D_2, D_8, D_9\}$   $S_{Rain} = \{D_4, D_5, D_6, D_{10}\}$   $S_{Overcast} = \{D_3, D_7\}$

$G(S_{Sunny}, Humidity) = E(Sunny) - \frac{3}{4} \times E(high) - \frac{1}{4} E(normal) = E(sunny) = 0.8113$

$\underset{\text{Humidity}}{[H, 3-]}$

high / normal

$[0+, 3-]$   $[1+, 0-]$

$E(high) = -0 \times \log_2 0 - 1 \times \log_2 \frac{3}{3} = 0$

$E(normal) = -1 \times \log_2 \frac{1}{1} - 0 \times \log_2 0 = 0$

$G(S_{Sunny}, Temperature) = E(Sunny) - \frac{2}{4} \times E(hot) - \frac{1}{4} \times E(mild) - \frac{1}{4} \times E(cold) = 0.8113$

$\underset{\text{Temperature}}{[1+, 3-]}$

hot / mild \ cold

$[0+, 2-]$   $[0+, 1-]$   $[1+, 0-]$

$E(hot) = 0$

$E(mild) = 0$

$E(cold) = 0$

$G(S_{Sunny}, Wind) = E(Sunny) - \frac{3}{4} \times E(weak) - \frac{1}{4} \times E(strong) = 0.1226$

$\underset{\text{Wind}}{[1+, 3-]}$

weak / strong

$[1+, 2-]$   $[0, 1-]$

$E(weak) = -\frac{1}{3} \times \log_2 \frac{1}{3} - \frac{2}{3} \times \log_2 \frac{2}{3}$

$E(strong) = 0$

∴ $G(S_{Sunny}, H) = G(S_{Sunny}, T) > G(S_{Sunny}, W)$

$G(S_{Rain}, Humidity) = E(Rain) - \frac{1}{4} \times E(high) - \frac{3}{4} \times E(normal) = 0.1226$

[3+, 1-]
Humidity

high $\diagup$ $\diagdown$ normal

[1+, 0-]　　[2+, 1-]

$E(high) = 0$

$E(normal) = \frac{2}{3} \times \log_2 \frac{2}{3} - \frac{1}{3} \times \log_2 \frac{1}{3}$

$G(S_{Rain}, Temperature) = E(Rain) - 0 \times E(hot) - \frac{2}{4} \times E(mild) - \frac{2}{4} \times E(cold) = 0.3113$
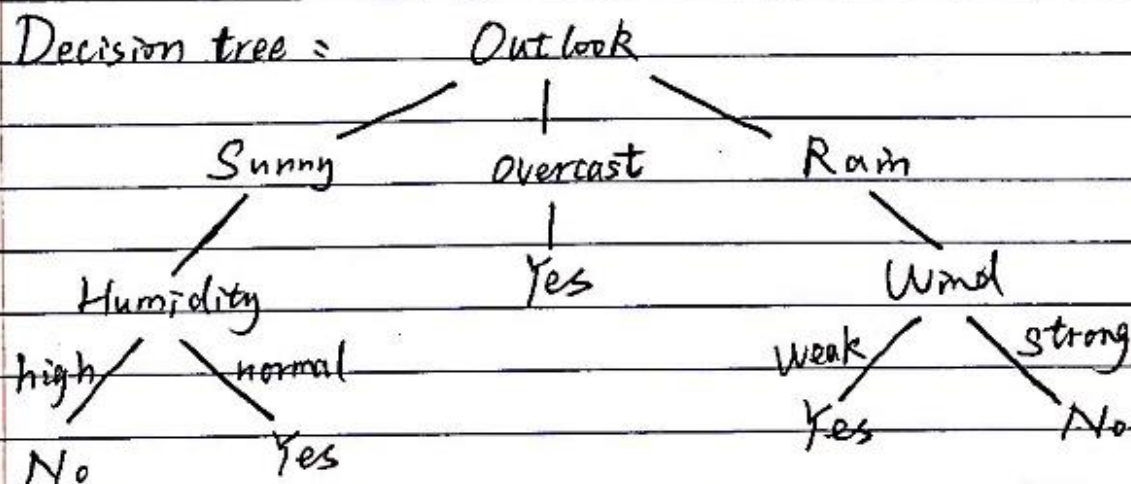
[3+, 1-]
Temperature

hot $\diagup$ |mild\ cold

[0+, 0-]　[2+, 0]　[1+, 1-]

$E(mild) = 0$

$E(hot) = 0$

$E(cold) = \frac{1}{2} \times \log_2 \frac{1}{2} - \frac{1}{2} \times \log_2 \frac{1}{2}$ #

$G(S_{Rain}, Wind) = E(Rain) - \frac{3}{4} \times E(weak) - \frac{1}{4} \times E(strong) = 0.8113$

[3+, 1-]
Wind

weak $\diagup$ $\diagdown$ strong

[3+, 0]　　[0, 1-]

$E(weak) = 0$

$E(strong) = 0$

So, $G(S_{Rain}, Wind) > G(S_{Rain}, Temperature) > G(S_{Rain}, Humidity)$.

Decision tree =

Outlook

Sunny　　overcast　　Rain

Humidity　　Yes　　Wind

high $\diagup$ \normal　　　weak $\diagup$ \strong

No　　Yes　　　　Yes　　No

② 　[6,4]

　$\diagup$ | $\diagdown$

[1,3] [2,0] [3,1]

$\hat{p}_1 = \frac{6}{10} \times 4 = 2.4$　$\hat{p}_2 = \frac{6}{10} \times 2 = 1.2$　$\hat{p}_3 = \frac{6}{10} \times 4 = 2.4$

$\hat{n}_1 = \frac{4}{10} \times 4 = 1.6$　$\hat{n}_2 = \frac{4}{10} \times 2 = 0.8$　$\hat{n}_3 = \frac{4}{10} \times 4 = 1.6$

$Q = \sum_{i=1}^{3} \frac{(p_i - \hat{p}_i)^2}{\hat{p}_i} + \frac{(n_i - \hat{n}_i)^2}{\hat{n}_i} = 3.75$

df = 2　Thus, $P(X^2 > CV) = 0.15$

4. The outputs of individual classifiers are 1 1 2

$\beta(w_1 | d_{1,1}(x)=1) = \frac{40}{70}$      $\hat{P}(w_2 | d_{1,1}(x)=1) = \frac{30}{70}$

$\beta(w_1 | d_{2,1}(x)=1) = \frac{20}{40}$      $\beta(w_2 | d_{2,1}(x)=1) = \frac{20}{40}$

$\beta(w_1 | d_{3,2}(x)=1) = \frac{0}{10}$      $\beta(w_2 | d_{3,2}(x)=1) = \frac{10}{10}$

Thus, class $1 = \frac{40}{70} \times \frac{20}{40} \times \frac{0}{10} = 0$

class $2 = \frac{30}{70} \times \frac{20}{40} \times \frac{10}{10} = 0.21$

So the final decision is class 2.