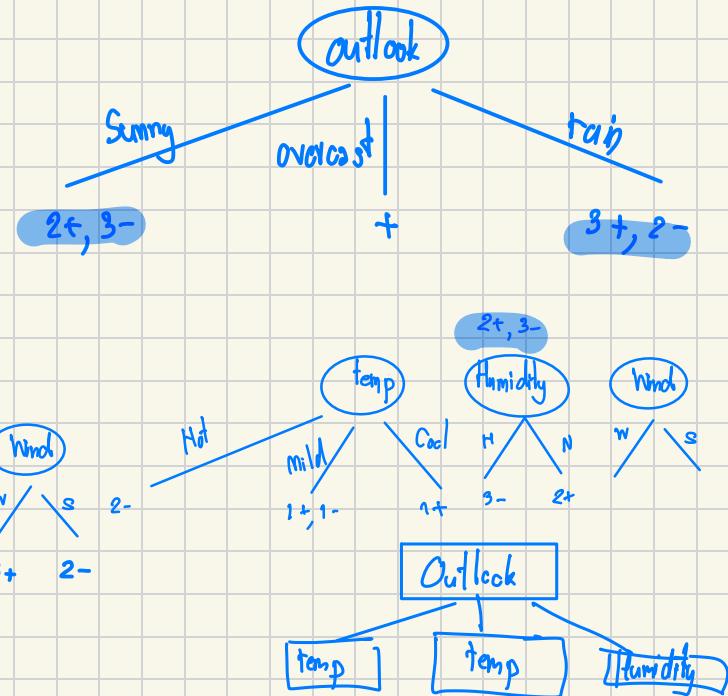


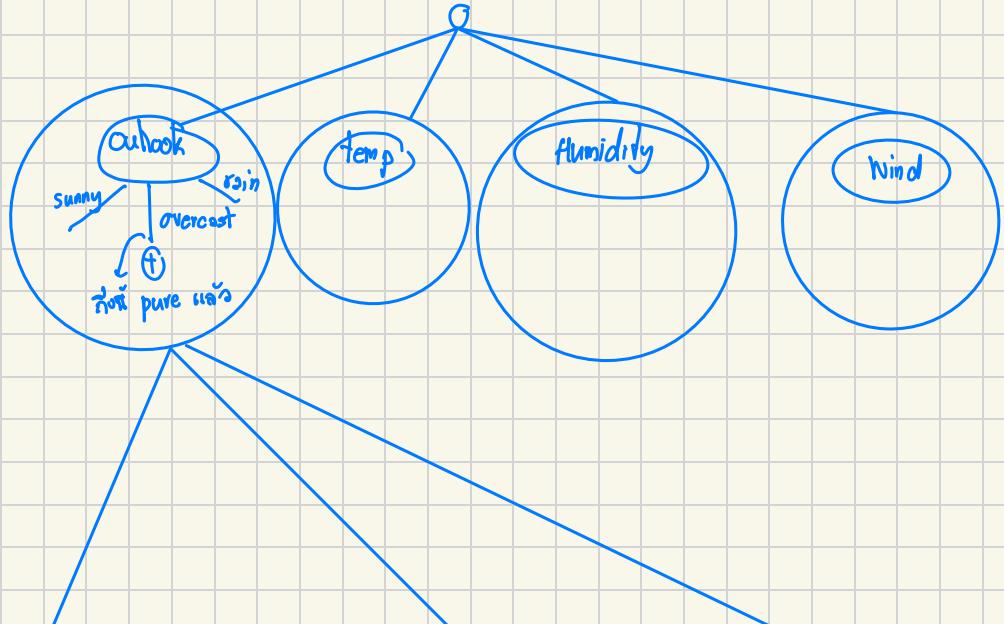
Decision Tree

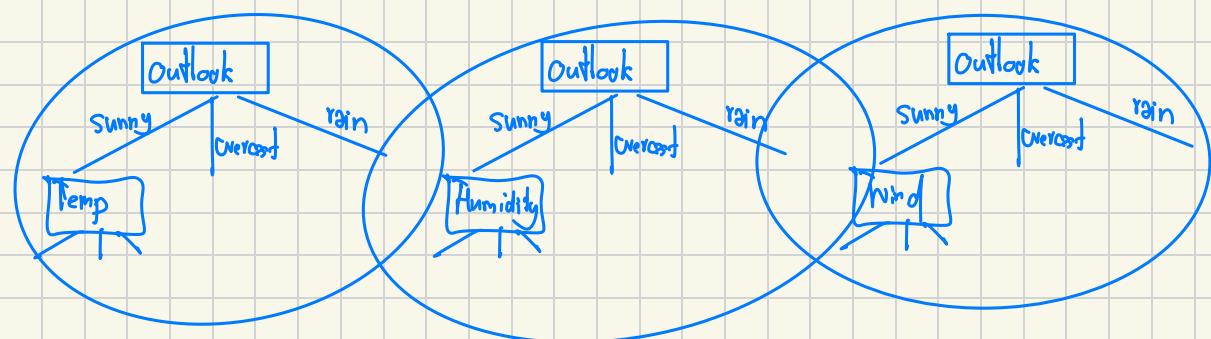
Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



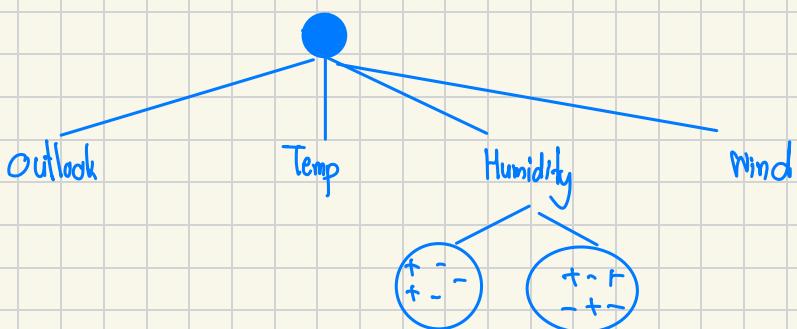
Prefers smaller trees

AJ : BFS





Greedy

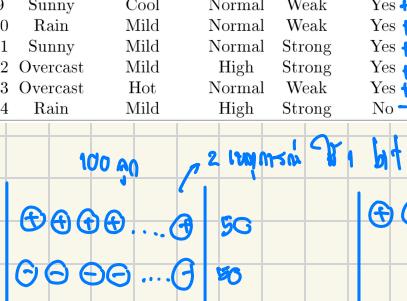


9+, 5 -

Entropy ສູນທະນາຄະນະ thermodynamic ດີວ່າດໍາລັງການສົບສຳຫຼວງ
(Information theory)



2 ຜົນຍອດທີ່ໄດ້ຮັບ



100 ກບ

2 ລະບຽບ ທີ່ 1 bit



2 ລະບຽບ ທີ່ 1 bit ມີ 1 bit ໄດ້ມີ



1 ລະບຽບ 0 bit

ກໍານົດເປົ້າ

ເສັ້ນຫຼືກົດ

$$\text{Entropy } (S) = \sum_v -p_v \log_2 p_v$$

$$\text{Entropy } \{+, +, +, +, -, -, -, -\} = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= -\frac{1}{2} \log_2 2^{-1} - \frac{1}{2} \log_2 2^{-1}$$

$$= \frac{1}{2} + \frac{1}{2} = 1 \text{ น้อยกว่า } 1 \text{ bit หมายความว่าไม่แน่นอน}$$

$$\text{Entropy } \{+, +, +, +, +, +, +\} = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -1 \cancel{\log_2 1} - 0 \cancel{\log_2 0}$$

$$= 0 \text{ bit}$$

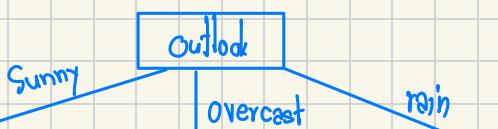
$$\text{Entropy } \{+, +, +, +, +, -\} = -0.99 \cancel{\log_2 0.99} - 0.01 \cancel{\log_2 0.01} \approx 0$$

$$\approx 0.052 \dots$$

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$$S_0 = \{+, +, +, +, +, +, +, -, -, -, -, -, -\}$$

$$\text{Entropy } (S_0) = -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) = 0.940$$



Entropy

Weighted
Average

$$-\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.971$$

$$4+0- = 0$$

$$3+2- = 0.971$$

$$= \frac{5}{14} (0.971) + \frac{4}{14} (0) + \frac{5}{14} (0.971)$$

หากเลือกเล่นเที่ยว 3 วัน

พื้นที่อยู่ต่อวัน ยกเว้นต้องใช้

ก่อนหน้า outlook - หลัง outlook

$$\text{Gain}(S_0, \text{Outlook}) = 0.940 - \text{Weighted Average}$$

$$= 0.940 - 0.694 = 0.246$$

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Humidity

High

3+, 4-

Normal

6+, (1-) ไม่เกิดข้างกันคือ

Entropy $= -\frac{3}{7} \log_2 \left(\frac{3}{7}\right) -$

$$\frac{4}{7} \log_2 \left(\frac{4}{7}\right)$$

$$= 0.985$$

$$-\frac{6}{7} \log_2 \left(\frac{6}{7}\right) - \frac{1}{7} \log_2 \left(\frac{1}{7}\right)$$

$$= 0.592$$

$$\text{Weighted Average} = \frac{7}{14} (0.985) + \frac{7}{14} (0.592) \\ = 0.789$$

$$\text{Gain}(S_0, \text{Humidity}) = 0.940 - 0.789 = 0.151$$

pure ฝ่ายใดฝ่ายหนึ่ง

Outlook

Sunny

2+, 3-

$$S_1 = \{2+, 3-\}$$

$$\text{Gain}(S_1, \text{Temp})$$

$$\text{Gain}(S_1, \text{Humidity})$$

$$\text{Gain}(S_1, \text{Wind})$$

Rain

3+, 2-

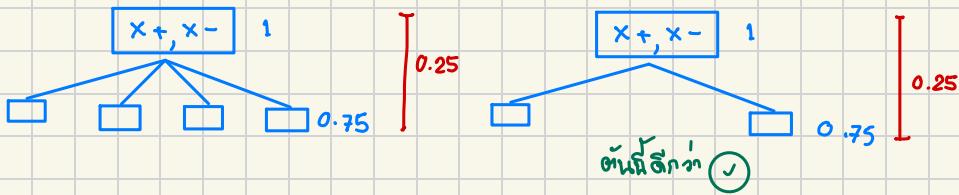
$$S_2 = \{3+, 2-\}$$

$$\text{Gain}(S_2, \text{Temp})$$

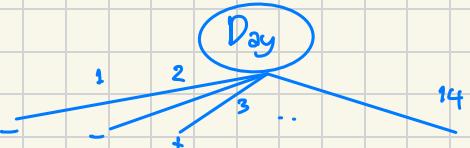
$$\text{Gain}(S_2, \text{Humidity})$$

$$\text{Gain}(S_2, \text{Wind})$$

4/2/2024 Gain → Information Gain នាមខ្លោយបង្កើតការពិនិត្យ



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

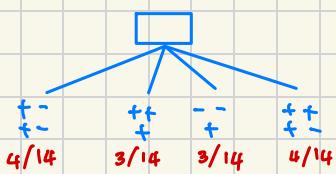


Gain Ratio

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}(A)}$$

$$\text{SplitInfo}(D) = - \sum_{j=1}^n \frac{|D_j|}{|D|} \log_2 \frac{|D_j|}{|D|}$$

$$= - \frac{4}{14} \log_2 \left(\frac{4}{14} \right) - \frac{3}{14} \log_2 \left(\frac{3}{14} \right) - \frac{3}{14} \log_2 \left(\frac{3}{14} \right) - \frac{4}{14} \log_2 \left(\frac{4}{14} \right)$$

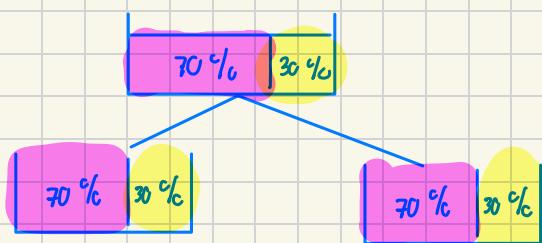
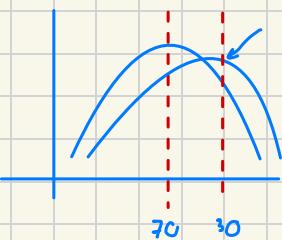


$$2 \text{ សំណើ សំលេងរូប} : -\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) = 1$$

$$4 \text{ សំណើ សំលេងរូប} = -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right)$$

$$= -1 \log_2 \left(\frac{1}{4} \right)$$

$$= -1 \log_2 2^{-2} = 2$$

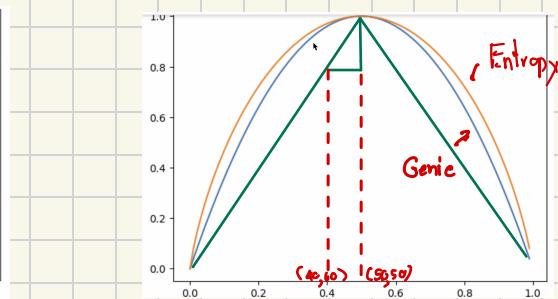
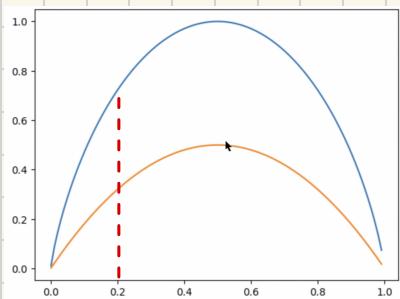
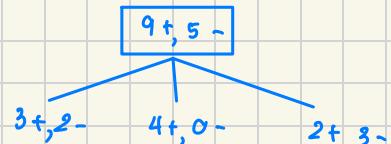


ໃນຕາມ ML ມີລາຍລະອຽດ ທີ່ມີ distributed data
ແລະ ຂອງມາຈິງໄປໄວ້ໄດ້ເທົ່າມາເປັນ 50 : 50

Gini Impurity

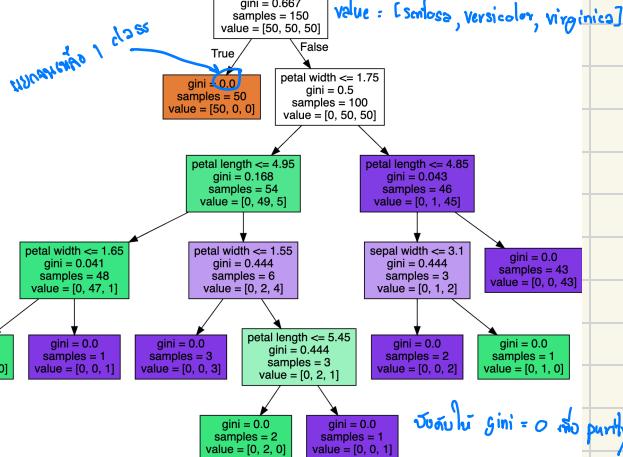
$$\text{Gini}(D) = 1 - \sum_v P_v^2$$

$$= \frac{5}{14} (\text{Gini}(3+, 2-)) + \frac{4}{14} (\text{Gini}(4+, 0-)) + \frac{5}{14} (\text{Gini}(2+, 3-))$$



ກົດຕິກຳ Genie ຕື່ອ ຂັບໄດ້ຮັງກວ່າ ກັບຫຼິໄວ້ສົດສະພາທີ sensitive ກວ່າ entropy

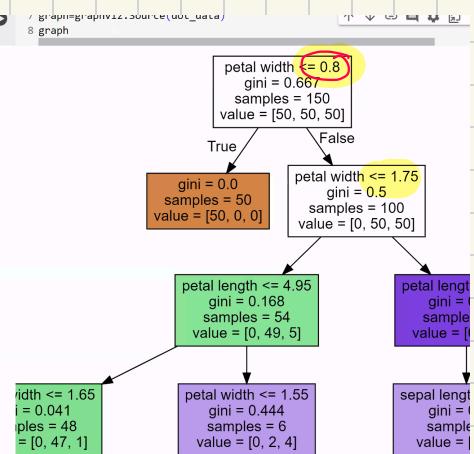
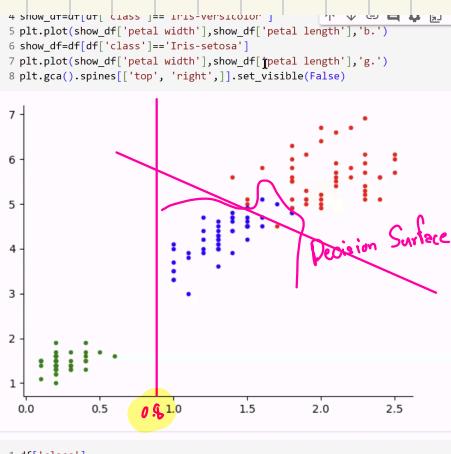
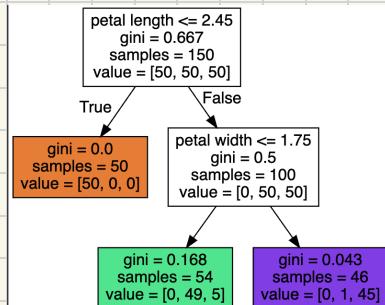
ກັບຫຼິໄວ້ກຳນົດຢູ່ນັກ ດັ່ງນັ້ນກັບຫຼິໄວ້ກຳນົດຢູ່ນັກ



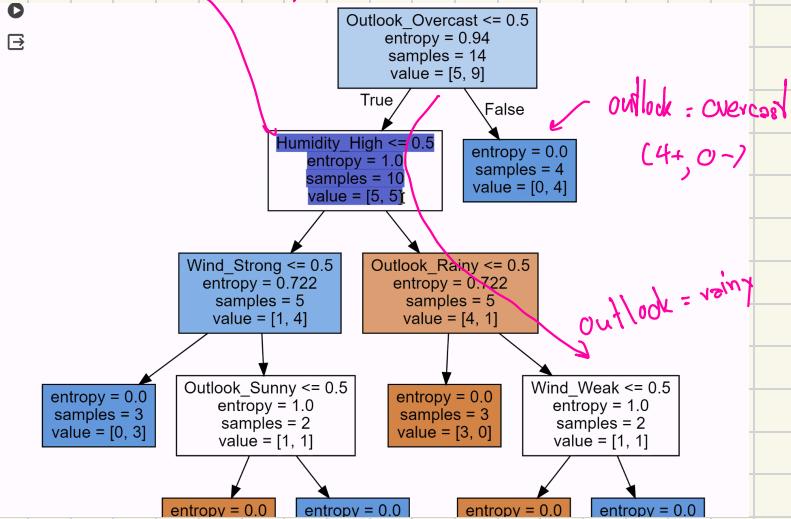
ரೂಪಕ್ಕೆ gini = 0 ಸುಧಾರಣೆ

min_sample_split = 4 ರೊಗೆ sample < 4 ಇಲ್ಲಿನ್ನೂ

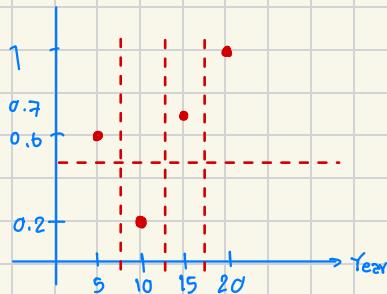
min_impurity_decrease = 0.05 ಈಂದೆ > 0.05 ರೊಗೆ split ಏಂ



in outlook sunny no rainy



sklearn սընթացակարգություն continuous միավոր 3 բարելու



$$R = [0.1, -0.3, 0.1, 0.2]$$

$$\text{Similarity Score} = \frac{(\sum R)^2}{\#R + \gamma}$$

$$\gamma = 0$$

minimum output R score = 0.5

Year < 7.5

0.1

$[-0.3, 0.1, 0.2]$

$$\text{Similarity Score} = \frac{(0.1 - 0.3 + 0.1 + 0.2)^2}{4 + 0}$$

$$= \frac{0.01}{4} = 0.0025$$

$$\frac{(0.1)^2}{1+0}$$

$$= 0.01$$

$$\frac{(-0.3 + 0.1 + 0.2)^2}{3+0}$$

$$= 0$$

$$\text{Gain} = S_{\text{left}} + S_{\text{right}} - S_{\text{root}}$$

$$= 0.01 + 0 - 0.0025$$

$$= 0.0075$$

Year < 12.5

0.1, -0.3

0.1, 0.2

$$\text{Gain} = 0.02 + 0.045 - 0.0025$$

$$= 0.065 - 0.0025$$

$$= 0.0625$$

$$\frac{(-0.2)^2}{2} = 0.02$$

$$\frac{(0.1 + 0.2)^2}{2+0} = 0.045$$

Year < 17.5

0.1, -0.3, 0.1

0.2

$$\text{Gain} = 0.0033 + 0.04 - 0.0025$$

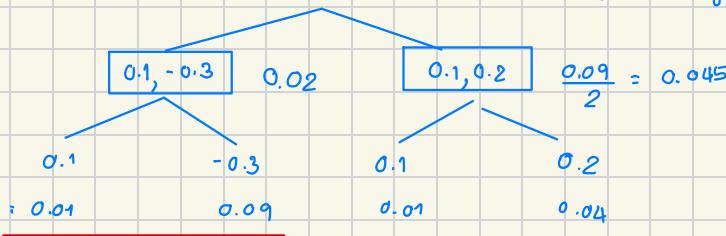
$$= 0.040$$

$$\frac{(-0.1)^2}{3} = 0.0033$$

$$\frac{(0.2)^2}{1} = 0.04$$

Year < 12.5

Gain > 8

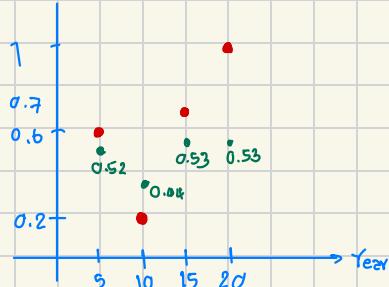


Year < 12.5

Year < 7.5

0.1 0.2

$$\text{Output} = \frac{\bar{Z}R}{\#R + \lambda} = 0.15$$



$$\text{new_output} > 0.5 + \text{learning rate} \times \text{output}$$

$$= 0.5 + 0.2 \times \text{output}$$

$$0.5 + (0.2)(0.1) = 0.52$$

$$0.5 + (0.2)(-0.3) = 0.44$$

$$0.5 + (0.2)(0.15) = 0.53$$

R	0.08	-0.24	0.07	0.17
---	------	-------	------	------

Year < 7.5

$$\text{Similarity score} = \frac{(0.08 - 0.24 + 0.07 + 0.17)^2}{4+0}$$

$$= \frac{(0.08)^2}{4}$$

$$\text{Similar} = \frac{(0.08)^2}{1} = \frac{(-0.24 + 0.07 - 0.17)^2}{3} = 0.0064$$

$$= 0.0016$$

$$= 0.0064 = 0$$

$$\text{Gain} > 0.0064 + 0 - 0.0016 = 0.0048$$

Year < 18.5

0.08, -0.24

0.07 0.17

$$\frac{0.16^2}{2}$$

$$\frac{0.24^2}{2}$$

$$Gain = 0.0128 + 0.0288 - 0.0016$$

$$= 0.0128$$

$$= 0.0288$$

$$+ 0.0284 \text{ Q}$$

```
[27] # clf.fit(X,y)
```

```
XGBClassifier(base_score=None, booster=None, callbacks=None,
    colsample_bylevel=None, colsample_bynode=None,
    colsample_bytree=None, device=None, early_stopping_rounds=None,
    enable_categorical=False, eval_metric=None, feature_types=None,
    gamma=None, learning_rate=None, max_delta=None, max_depth=None,
    max_cat_threshold=None, max_cat_to_prove=None,
    max_delta_step=None, max_depth=None, max_leaves=None,
    min_child_weight=None, missing='nan', monotone_constraints=None,
    multi_strategy=None, n_estimators=None, n_jobs=None,
    num_parallel_tree=None, random_state=None, ...)
```

```
[28] 1 clf.feature_importances_
```

தகவமளிக்குவது attribute

```
array([0.20695841, 0.03330886, 0.07800627, 0.04108072, 0.02753872,
       0.53931534, 0.          , 0.02886179, 0.01932171, 0.02560823],
      dtype=float32)
```

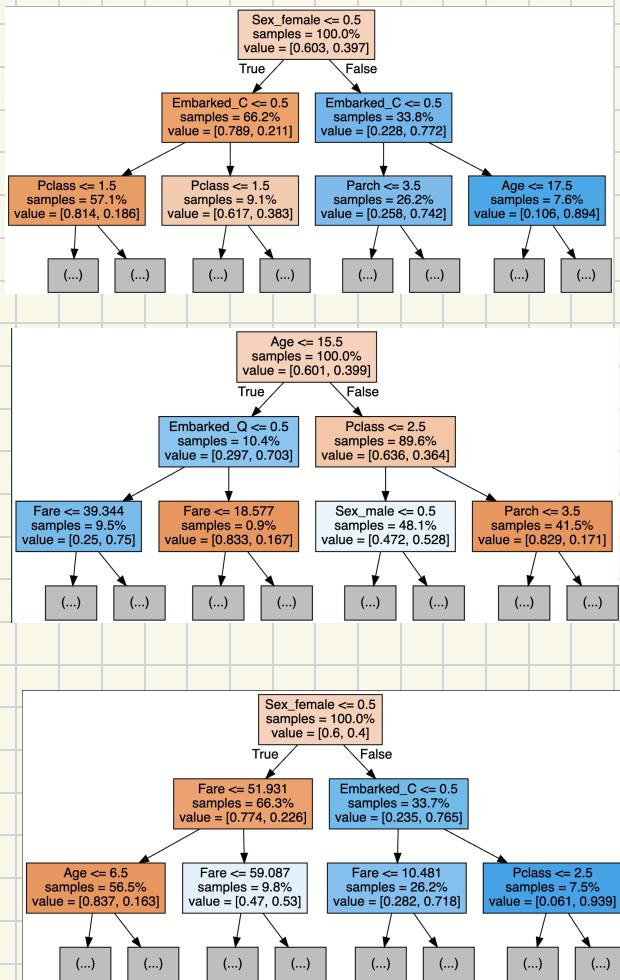
```
[29] 1 clf.feature_names_in_
```

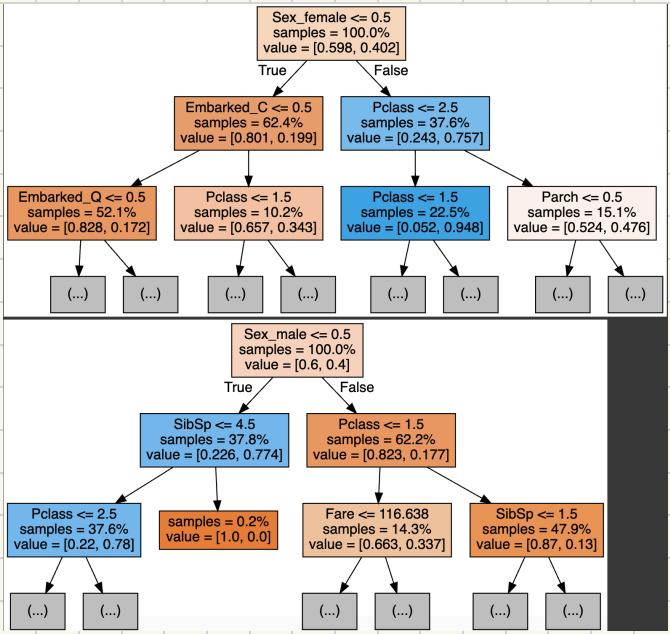
```
array(['Pclass', 'Age', 'SibSp', 'Parch', 'Fare', 'Sex_female',
       'Sex_male', 'Embarked_C', 'Embarked_Q', 'Embarked_S'], dtype='<U10')
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.metrics import accuracy_score
3 X_train,X_test,Y_train,Y_test=train_test_split(X,y,test_size=0.2)
4
5 rf=RandomForestClassifier()
6 xgb=XGBClassifier()
7 rf.fit(X_train,Y_train)
8 xgb.fit(X_train,Y_train)
9 y_pred=rf.predict(X_test)
10 print(f'RandomForest {accuracy_score(y_pred,Y_test)}')
11 y_pred=xgb.predict(X_test)
12 print(f'XGBoost {accuracy_score(y_pred,Y_test)}')
```

```
RandomForest 0.7482517482517482
XGBoost 0.7692307692307693
```

Random Forest display





https://colab.research.google.com/drive/1WWx_aR1LWnVA_wuFgVKgLuwr1Px_I5lc?usp=sharing