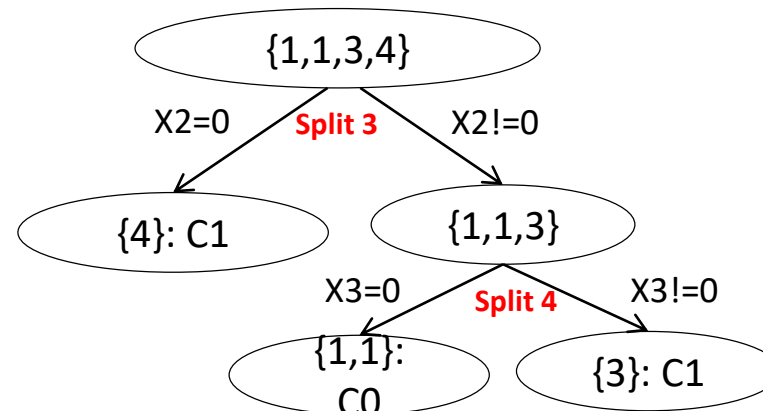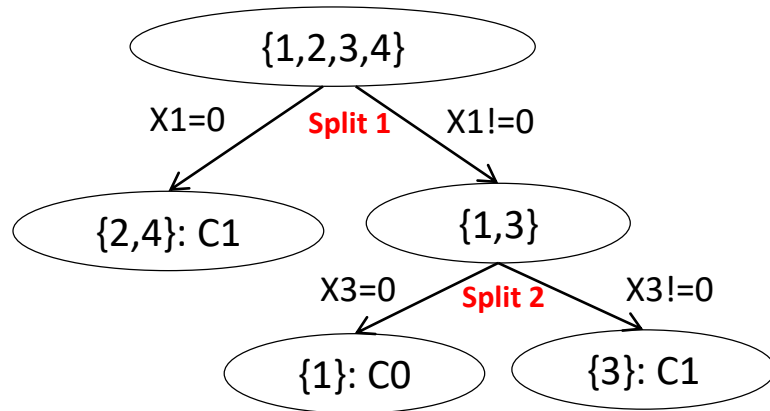# Lecture 12: Variable Importance in Tree Models

Instructor: Prof. Shuai Huang

Industrial and Systems Engineering

University of Washington

# Importance score in Random Forest (RF)

| ID | $X_1$ | $X_2$ | $X_3$ | Class |
|----|-------|-------|-------|-------|
| 1 | 1 | 1 | 0 | C0 |
| 2 | 0 | 0 | 0 | C1 |
| 3 | 1 | 1 | 1 | C1 |
| 4 | 0 | 0 | 1 | C1 |

```
        {1,2,3,4}
   X1=0   Split 1   X1!=0
  {2,4}: C1        {1,3}
              X3=0  Split 2  X3!=0
            {1}: C0        {3}: C1
```

```
        {1,1,3,4}
   X2=0   Split 3   X2!=0
  {4}: C1         {1,1,3}
              X3=0  Split 4  X3!=0
           {1,1}:            {3}: C1
            C0
```

# Importance score in Regularized RF (RRF)

| ID | $X_1$ | $X_2$ | $X_3$ | Class |
|----|-------|-------|-------|-------|
| 1  | 1     | 1     | 0     | C0    |
| 2  | 0     | 0     | 0     | C1    |
| 3  | 1     | 1     | 1     | C1    |
| 4  | 0     | 0     | 1     | C1    |

The regularized impurity gain of variable $X_i$ at a node is calculated as

$$Gain'(X_i) = \begin{cases} \lambda \cdot Gain(X_i) & X_i \notin F \\ Gain(X_i) & X_i \in F \end{cases}$$

# Importance score in Guided RRF (GRRF)

In GRRF, instead having one $\lambda$ for all variables, each variable $X_i$ can have its own $\lambda_i$:

$$Gain'(X_i) = \begin{cases} \lambda_i \cdot Gain(X_i) & X_i \notin F \\ Gain(X_i) & X_i \in F \end{cases},$$

where $\lambda_i$ is

$$\lambda_i = (1 - \gamma)\lambda_0 + \gamma * w_i,$$

where $\lambda_0$ controls the base regularization, $w_i \in [0,1]$ is a prior of importance of each variable $v_i$, and $\gamma \in [0,1]$ controls the weight from the prior.