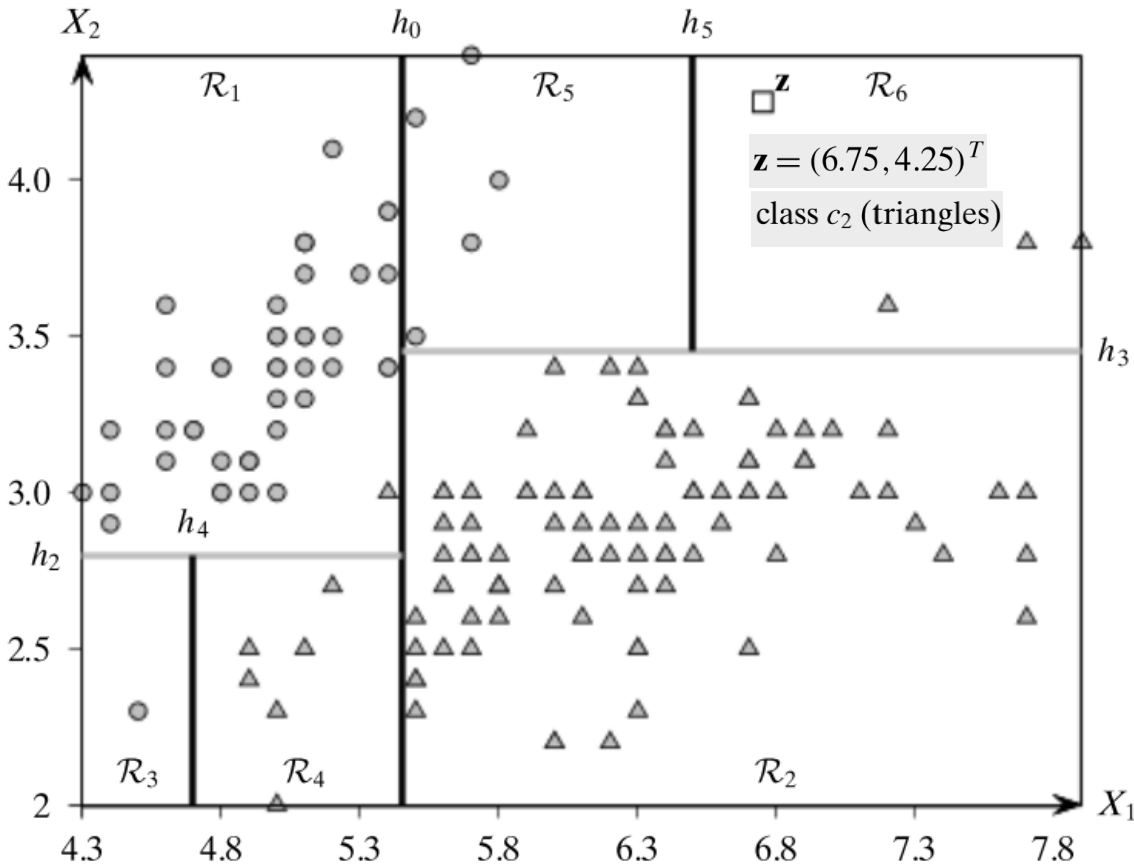


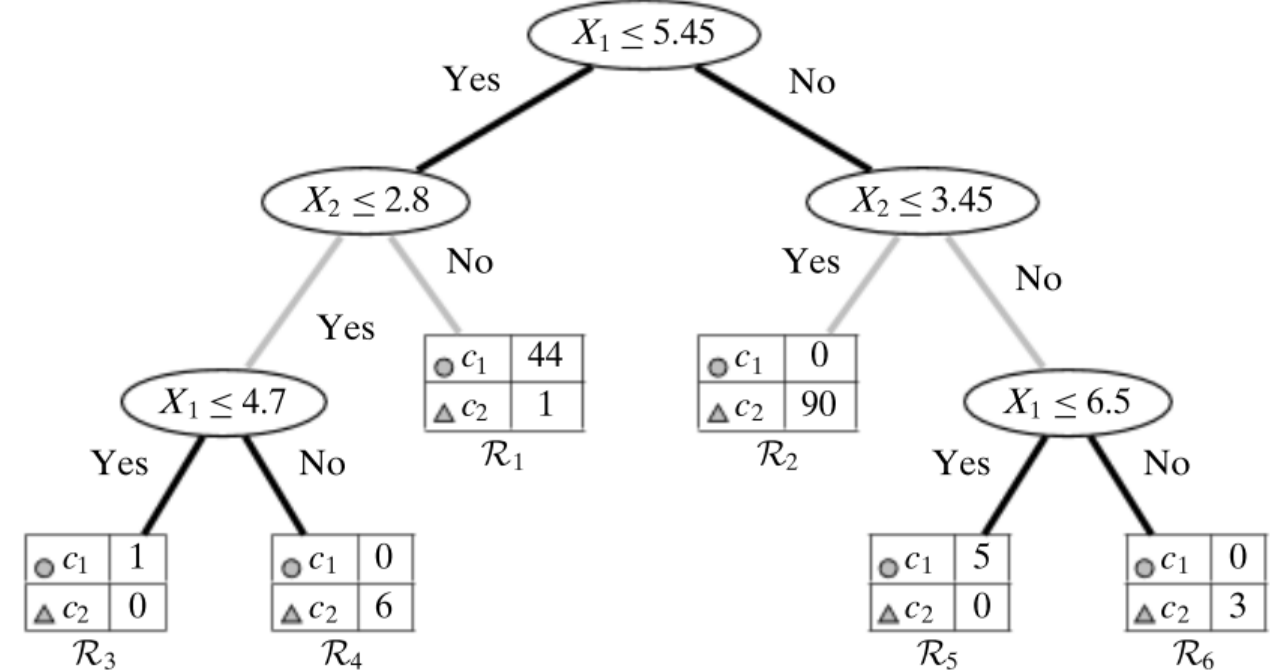
# دسته‌بندی گر درخت تصمیم

Decision Tree Classifier



(a) Recursive Splits

هر یک از مناطق  $\mathcal{R}_i$  با تقسیم مکرر ابرصفحه‌ها موازی محورها تا زمانی که نقاط درون یک پارتیشن از نظر برجسب دسته خود نسبتاً خالص باشند، یعنی بیشتر نقاط متعلق به یک دسته هستند.



(b) Decision Tree

Figure 19.1. Decision trees: recursive partitioning via axis-parallel hyperplanes.

برای نسبت دادن یک نقطه‌ی آزمون  $x$  به یک دسته، فقط باید مشخص کنیم به کدام منطقه تعلق دارد و در آن منطقه، دسته‌ی غالب کدام است.

ابرصفحه‌های موازی محور (Axis-Parallel Hyperplanes)

$$h(\mathbf{x}): \mathbf{w}^T \mathbf{x} + b = 0$$

$$h(x): \mathbf{e}_j^T \mathbf{x} + b = x_j + b = 0$$

ابرصفحه‌ی  $h(x) = 0$  که عمود بر مشخصه‌ی  $X_j$  می‌باشد و دامنه‌ی  $\mathcal{R}$  برای مشخصه‌ی  $X_j$  را به دو ناحیه  $\mathcal{R}_Y$  و  $\mathcal{R}_N$  تقسیم می‌کند. ناحیه‌ی  $\mathcal{R}_Y$  نقاطی هستند که در شرط  $h(x) \leq 0$  صدق می‌کنند و ناحیه‌ی  $\mathcal{R}_N$  نقاط در ناحیه‌ی مقابل هستند. **نقطه‌ی تقسیم (Split Point)** ناحیه‌ی  $\mathcal{R}$  با انتخاب  $v = -b$  مشخص می‌شود. با این تقسیم مجموعه داده‌های  $\mathbf{D}$  نیز به دو زیر مجموعه‌ی  $\mathbf{D}_Y$  و  $\mathbf{D}_N$  بنابر این که شرط  $X_j \leq v$  را ارضا نمایند و یا خیر، تقسیم می‌شوند.

$$\mathbf{D}_Y = \{\mathbf{x}^T \mid \mathbf{x} \in \mathbf{D}, x_j \leq v\}$$

$$\mathbf{D}_N = \{\mathbf{x}^T \mid \mathbf{x} \in \mathbf{D}, x_j > v\}$$

برای **مشخصه‌های دسته‌ای (Categorical Attributes)**، زیردامنه‌ای از  $X_j$  به شکل  $V \subset \text{dom}(X_j)$  برای تقسیم ناحیه  $\mathcal{R}$  انتخاب می‌شود. ناحیه‌ی  $\mathcal{R}_Y$  نقاطی از مجموعه داده‌های  $\mathbf{D}$  هستند که  $x_i \in V$

**خلوص (Purity)** برای ناحیه‌ی  $\mathcal{R}_j$  نسبت داده‌های دسته‌ای غالب در ناحیه به کل داده‌های موجود در آن ناحیه است.

$$\text{purity}(\mathbf{D}_j) = \max_i \left\{ \frac{n_{ji}}{n_j} \right\}$$

که در آن  $n_j = |\mathbf{D}_j|$  تعداد داده‌ها در ناحیه‌ی  $\mathcal{R}_j$  و  $n_{ji}$  تعداد نقاط در این ناحیه با برچسب دسته‌ی  $c_i$  می‌باشد.

---

**Algorithm 19.1:** Decision Tree Algorithm

---

**DECISIONTREE** ( $\mathbf{D}, \eta, \pi$ ):

- 1  $n \leftarrow |\mathbf{D}|$  // partition size
- 2  $n_i \leftarrow |\{\mathbf{x}_j | \mathbf{x}_j \in \mathbf{D}, y_j = c_i\}|$  // size of class  $c_i$
- 3  $\text{purity}(\mathbf{D}) \leftarrow \max_i \left\{ \frac{n_i}{n} \right\}$
- 4 **if**  $n \leq \eta$  **or**  $\text{purity}(\mathbf{D}) \geq \pi$  **then** // stopping condition
  - 5  $c^* \leftarrow \arg\max_{c_i} \left\{ \frac{n_i}{n} \right\}$  // majority class
  - 6 create leaf node, and label it with class  $c^*$
  - 7 **return**
- 8  $(\text{split point}^*, \text{score}^*) \leftarrow (\emptyset, 0)$  // initialize best split point
- 9 **foreach** (*attribute*  $X_j$ ) **do**
  - 10 **if** ( $X_j$  is numeric) **then**
    - 11  $(v, \text{score}) \leftarrow \text{EVALUATE-NUMERIC-ATTRIBUTE}(\mathbf{D}, X_j)$
    - 12 **if**  $\text{score} > \text{score}^*$  **then**  $(\text{split point}^*, \text{score}^*) \leftarrow (X_j \leq v, \text{score})$
  - 13 **else if** ( $X_j$  is categorical) **then**
    - 14  $(V, \text{score}) \leftarrow \text{EVALUATE-CATEGORICAL-ATTRIBUTE}(\mathbf{D}, X_j)$
    - 15 **if**  $\text{score} > \text{score}^*$  **then**  $(\text{split point}^*, \text{score}^*) \leftarrow (X_j \in V, \text{score})$
- // partition  $\mathbf{D}$  into  $\mathbf{D}_Y$  and  $\mathbf{D}_N$  using  $\text{split point}^*$ , and call recursively
- 16  $\mathbf{D}_Y \leftarrow \{\mathbf{x}^T | \mathbf{x} \in \mathbf{D} \text{ satisfies } \text{split point}^*\}$
- 17  $\mathbf{D}_N \leftarrow \{\mathbf{x}^T | \mathbf{x} \in \mathbf{D} \text{ does not satisfy } \text{split point}^*\}$
- 18 create internal node  $\text{split point}^*$ , with two child nodes,  $\mathbf{D}_Y$  and  $\mathbf{D}_N$
- 19 **DECISIONTREE**( $\mathbf{D}_Y$ ); **DECISIONTREE**( $\mathbf{D}_N$ )

---

شرط توقف رسیدن به حداکثر خلوص  $\pi$  و یا حداقل تعداد نقاط  $\eta$  می باشد.  
مقادیر  $\pi$  و  $\eta$  از پیش مشخص می شوند.

معیار انتخاب نقطه‌ی تقسیم (Split Point)

$$X_j \in V \text{ یا } X_j \leq v$$

آنترپی (Entropy) از رابطه زیر تعریف می‌شود که در آن  $P(c_i|\mathbf{D})$  احتمال دسته‌ی  $c_i$  در مجموعه داده‌ی  $\mathbf{D}$  و  $k$  تعداد کل دسته‌ها است. اگر ناحیه‌ای کاملاً خالص از دسته‌ی  $c_i$  باشد آنگاه  $P(c_i|\mathbf{D}) = 1$  و  $P(c_j|\mathbf{D})_{j \neq i} = 0$  و در نتیجه  $H(\mathbf{D}) = 0$  می‌شود و اگر ناحیه کاملاً مخلوط از هر  $k$  دسته باشد  $P(c_i|\mathbf{D}) = \frac{1}{k}$  و آنترپی می‌شود  $H(\mathbf{D}) = \log_2 k$  پس هر چه ناحیه خالص‌تر باشد آنترپی کمتر خواهد شد.

$$H(\mathbf{D}) = - \sum_{i=1}^k P(c_i|\mathbf{D}) \log_2 P(c_i|\mathbf{D})$$

$$H(\mathbf{D}_Y, \mathbf{D}_N) = \frac{n_Y}{n} H(\mathbf{D}_Y) + \frac{n_N}{n} H(\mathbf{D}_N)$$

نقطه‌ی تقسیمی که ناحیه را به نواحی خالص‌تر تقسیم کند، بهره‌ی اطلاعاتی (Information Gain) را افزایش می‌دهد.

$$Gain(\mathbf{D}, \mathbf{D}_Y, \mathbf{D}_N) = H(\mathbf{D}) - H(\mathbf{D}_Y, \mathbf{D}_N)$$

از هر طریق آنتروپی، بهره‌ی اطلاعاتی و یا عدد جینی بخواهیم نقطه‌ی تقسیم را محاسبه کنیم باید تابع جرمی احتمال (احتمال پسین)  $P(c_i|\mathbf{D})$  را برآورد نماییم.

**Algorithm 19.2:** Evaluate Numeric Attribute (Using Gain)

**EVALUATE-NUMERIC-ATTRIBUTE ( $\mathbf{D}, X$ ):**

- 1 sort  $\mathbf{D}$  on attribute  $X$ , so that  $x_j \leq x_{j+1}, \forall j = 1, \dots, n-1$
- 2  $\mathcal{M} \leftarrow \emptyset$  // set of midpoints
- 3 **for**  $i = 1, \dots, k$  **do**  $n_i \leftarrow 0$
- 4 **for**  $j = 1, \dots, n-1$  **do**
- 5     **if**  $y_j = c_i$  **then**  $n_i \leftarrow n_i + 1$   
       // running count for class  $c_i$
- 6     **if**  $x_{j+1} \neq x_j$  **then**
- 7          $v \leftarrow \frac{x_{j+1} + x_j}{2}; \mathcal{M} \leftarrow \mathcal{M} \cup \{v\}$  // midpoints
- 8         **for**  $i = 1, \dots, k$  **do**
- 9              $N_{vi} \leftarrow n_i$  // Number of points such that  $x_j \leq v$  and  $y_j = c_i$
- 10 **if**  $y_n = c_i$  **then**  $n_i \leftarrow n_i + 1$   
       // evaluate split points of the form  $X \leq v$
- 11  $v^* \leftarrow \emptyset; score^* \leftarrow 0$  // initialize best split point
- 12 **forall**  $v \in \mathcal{M}$  **do**
- 13     **for**  $i = 1, \dots, k$  **do**
- 14          $\hat{P}(c_i|\mathbf{D}_Y) \leftarrow \frac{N_{vi}}{\sum_{j=1}^k N_{vj}}$
- 15          $\hat{P}(c_i|\mathbf{D}_N) \leftarrow \frac{n_i - N_{vi}}{\sum_{j=1}^k n_j - N_{vj}}$
- 16      $score(X \leq v) \leftarrow Gain(\mathbf{D}, \mathbf{D}_Y, \mathbf{D}_N)$  // use Eq. (19.5)
- 17     **if**  $score(X \leq v) > score^*$  **then**
- 18          $v^* \leftarrow v; score^* \leftarrow score(X \leq v)$
- 19 **return**  $(v^*, score^*)$

$$\hat{P}(X \leq v | c_i) = \frac{\hat{P}(X \leq v \text{ and } c_i)}{\hat{P}(c_i)} = \left( \frac{1}{n} \sum_{j=1}^n I(x_j \leq v \text{ and } y_j = c_i) \right) / (n_i / n)$$

$$= \frac{N_{vi}}{n_i}$$

که در آن  $N_{vi}$  تعداد نقاطی با برچسب  $c_i$  هستند که شرط  $x_j \leq v$  را نیز ارضاع می‌کنند.

$$\hat{P}(c_i | \mathbf{D}_Y) = \hat{P}(c_i | X \leq v) = \frac{N_{vi}}{\sum_{j=1}^k N_{vj}}$$

$$\hat{P}(c_i | \mathbf{D}_N) = \hat{P}(c_i | X > v) = \frac{\hat{P}(X > v | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \hat{P}(X > v | c_j) \hat{P}(c_j)} = \frac{n_i - N_{vi}}{\sum_{j=1}^k (n_j - N_{vj})}$$

مشخصه‌های دسته‌ای (Categorical Attributes)

**Algorithm 19.3:** Evaluate Categorical Attribute (Using Gain)

**EVALUATE-CATEGORICAL-ATTRIBUTE (D, X, l):**

```

1 for  $i = 1, \dots, k$  do
2    $n_i \leftarrow 0$ 
3   forall  $v \in \text{dom}(X)$  do  $n_{vi} \leftarrow 0$ 
4 for  $j = 1, \dots, n$  do
5   if  $x_j = v$  and  $y_j = c_i$  then  $n_{vi} \leftarrow n_{vi} + 1$  // frequency statistics
   // evaluate split points of the form  $X \in V$ 
6  $V^* \leftarrow \emptyset$ ;  $\text{score}^* \leftarrow 0$  // initialize best split point
7 forall  $V \subset \text{dom}(X)$ , such that  $1 \leq |V| \leq l$  do
8   for  $i = 1, \dots, k$  do
9      $\hat{P}(c_i | \mathbf{D}_Y) \leftarrow \frac{\sum_{v \in V} n_{vi}}{\sum_{j=1}^k \sum_{v \in V} n_{vj}}$ 
10     $\hat{P}(c_i | \mathbf{D}_N) \leftarrow \frac{\sum_{v \notin V} n_{vi}}{\sum_{j=1}^k \sum_{v \notin V} n_{vj}}$ 
11     $\text{score}(X \in V) \leftarrow \text{Gain}(\mathbf{D}, \mathbf{D}_Y, \mathbf{D}_N)$  // use Eq. (19.5)
12    if  $\text{score}(X \in V) > \text{score}^*$  then
13       $V^* \leftarrow V$ ;  $\text{score}^* \leftarrow \text{score}(X \in V)$ 
14 return  $(V^*, \text{score}^*)$ 

```

$$\begin{aligned}\hat{P}(X = v | c_i) &= \frac{\hat{P}(X = v \text{ and } c_i)}{\hat{P}(c_i)} \\ &= \left( \frac{1}{n} \sum_{j=1}^n I(x_j = v \text{ and } y_j = c_i) \right) / (n_i / n) \\ &= \frac{n_{vi}}{n_i}\end{aligned}$$

$$\hat{P}(c_i | \mathbf{D}_Y) = \frac{\sum_{v \in V} \hat{P}(X = v | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \sum_{v \in V} \hat{P}(X = v | c_j) \hat{P}(c_j)} = \frac{\sum_{v \in V} n_{vi}}{\sum_{j=1}^k \sum_{v \in V} n_{vj}}$$

$$\hat{P}(c_i | \mathbf{D}_N) = \hat{P}(c_i | X \notin V) = \frac{\sum_{v \notin V} n_{vi}}{\sum_{j=1}^k \sum_{v \notin V} n_{vj}}$$

که در آن  $n_{vi}$  تعداد نقاطی هستند که برچسب  $c_i$  دارند و مشخصه‌ی  $X$  آن‌ها در دسته‌ی  $V$  قرار می‌گیرد.

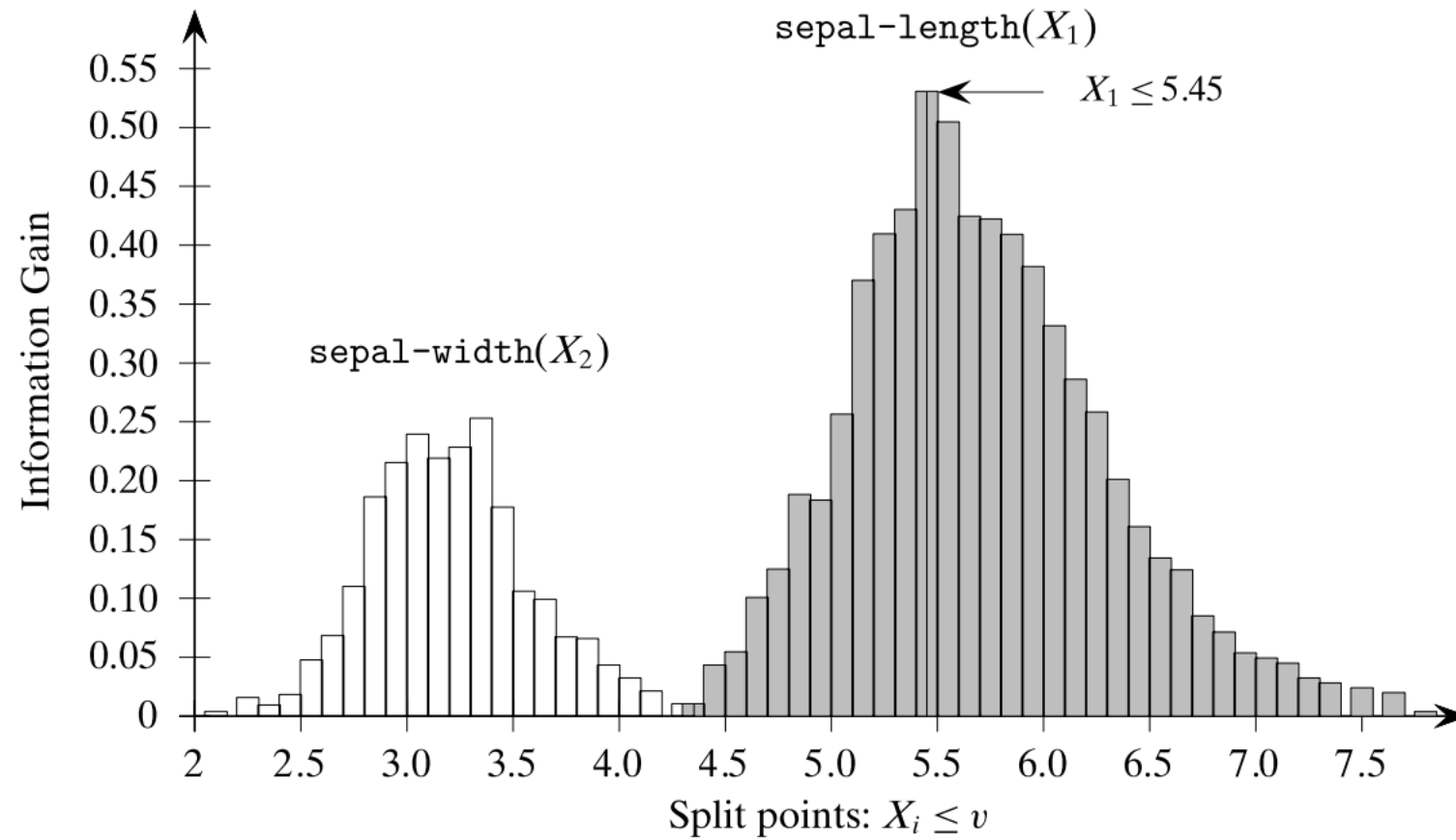


Figure 19.3. Iris: gain for different split points, for sepal length and sepal width .