

# CS3244 Machine Learning

AY2021/22 Semester 2

## 1. Concept Learning

### 1.1. Definitions

- **Concept:** A boolean-valued function over a set of input instances (each comprising input attributes).
- **Hypothesis Space ( $H$ ):** Each hypothesis  $h \in H$  ( $h: X \rightarrow \{0,1\}$ ) is represented by a conjunction of constraints on input attributes.
  - Each constraint can be a specific value, don't care (?), or no value allowed ( $\emptyset$ ).
  - An input instance  $x \in X$  satisfies a hypothesis  $h \in H$  if and only if  $h(x) = 1$ . In other words,  $h$  classifies  $x$  as a +ve example.
- **Concept Learning:** Given an unknown target concept  $c: X \rightarrow \{0,1\}$  and noise-free training examples  $D$ : +ve and -ve examples of the target concept, determine a hypothesis  $h \in H$  that is consistent with  $D$ .
  - A hypothesis  $h$  is consistent with a set of training examples  $D$  if and only if  $h(x) = c(x)$  for all  $(x, c(x)) \in D$ .
  - $h$  is consistent with  $D$  if and only if every +ve training instance satisfies  $h$  and every -ve training instance does not satisfy  $h$ .
- **Inductive Learning Assumption:** Any hypothesis found to approximate the target function well over a sufficiently large set of training examples will also approximate the target function well over other unobserved examples.
- **Version Space ( $VS_{H,D}$ ):** The version space ( $VS_{H,D}$ ) w.r.t. hypothesis space  $H$  and training examples  $D$ , is the subset of hypotheses from  $H$  consistent with  $D$ :
 
$$VS_{H,D} = \{h \in H \mid h \text{ is consistent with } D\}$$
  - If  $c \in H$ , then a large enough  $D$  can reduce  $VS_{H,D}$  to  $\{c\}$ .
  - If  $D$  is insufficient, then  $VS_{H,D}$  represents the uncertainty of what the target concept is.
  - The general boundary  $G$  of  $VS_{H,D}$  is the set of maximally general members of  $H$  consistent with  $D$ .
  - The specific boundary  $S$  of  $VS_{H,D}$  is the set of maximally specific members of  $H$  consistent with  $D$ .
  - **Version Space Representation Theorem (VSRT):**

$$VS_{H,D} = \{h \in H \mid \exists s \in S \exists g \in G, g \geq_g h \geq_g s\}$$
  - An input instance  $x$  satisfies every hypothesis in  $VS_{H,D}$  if and only if  $x$  satisfies every member of  $S$ .

- An input instance  $x$  satisfies none of the hypotheses in  $VS_{H,D}$  if and only if  $x$  satisfies none of the members of  $G$ .

### 1.2. Properties of Hypotheses

- $h_j$  is more general than or equal to  $h_k$  ( $h_j \geq_g h_k$ ) if and only if any input instance  $x$  that satisfies  $h_k$  also satisfies  $h_j$ :
  - $\forall x \in X, (h_k(x) = 1) \rightarrow (h_j(x) = 1)$
  - $\geq_g$  relation defines a partial order (reflexive, antisymmetric and transitive).
  - $h_j$  is more general than  $h_k$  if and only if  $h_j \geq_g h_k$  and  $h_k \not\geq_g h_j$ .
  - $h_j$  is more specific than  $h_k$  if and only if  $h_k$  is more general than  $h_j$ .
- Two hypotheses are syntactically distinct if any one of their attributes (specific values, ?,  $\emptyset$ ) are different.
- Two hypotheses are semantically distinct if any one of their attributes (specific values, ?) are different. Specially, all hypotheses containing  $\emptyset$  are semantically identical.

### 1.3. Find-S Algorithm

1. Initialise  $h$  to most specific hypothesis in  $H$
2. For each positive training instance  $x$ :
  - For each attribute constraint  $a_i$  in  $h$ :
    - If  $x$  satisfies  $a_i$  in  $h$ :
      - Do nothing
    - Else:
      - Replace  $a_i$  by the next more general constraint that is satisfied by  $x$
3. Output hypothesis  $h$

- **General idea:** Start with most specific hypothesis. Whenever it wrongly classifies a +ve training example as -ve, minimally generalize it to satisfy the input instance.
- **Correctness:** Suppose that  $c \in H$ . Then,  $h_n$  is consistent with  $D = \{(x_k, c(x_k))\}_{k=1,2,\dots,n}$ .
- **Limitations:**
  - Cannot tell whether Find-S has learnt target concept.
  - Cannot tell when training examples are inconsistent (contain error or noise).
  - Picks a maximally specific  $h$ .
  - Depending on  $H$ , there might be many maximally specific  $h$ .

### 1.4. List-Then-Eliminate Algorithm

1.  $VersionSpace \leftarrow$  a list containing every hypothesis in  $H$
2. For each training example  $(x, c(x))$ :

Remove from  $VersionSpace$  any hypothesis  $h$  for which  $h(x) \neq c(x)$

3. Output the list of hypotheses in  $VersionSpace$

- **General idea:** Start with all hypotheses in  $H$ . Then eliminate any hypothesis found inconsistent with any training example.
- **Limitations:**
  - Prohibitively expensive to exhaustively enumerate all hypotheses in finite  $H$ .

### 1.5. Candidate-Elimination Algorithm

1.  $G \leftarrow$  maximally general hypotheses in  $H$
2.  $S \leftarrow$  maximally specific hypotheses in  $H$
3. For each training example  $d$ :
  - If  $d$  is a +ve example:
    - Remove from  $G$  any hypothesis inconsistent with  $d$
    - For each  $s \in S$  not consistent with  $d$ :
      - Remove  $s$  from  $S$
      - Add to  $S$  all minimal generalisations  $h$  of  $s$  such that  $h$  is consistent with  $d$  and some member of  $G$  is more general than or equal to  $h$
      - Remove from  $S$  any hypothesis that is more general than another hypothesis in  $S$
  - If  $d$  is a -ve example:
    - Remove from  $S$  any hypothesis inconsistent with  $d$
    - For each  $g \in G$  not consistent with  $d$ :
      - Remove  $g$  from  $G$
      - Add to  $G$  all minimal specialisations  $h$  of  $g$  such that  $h$  is consistent with  $d$  and some member of  $S$  is more specific than or equal to  $h$
      - Remove from  $G$  any hypothesis that is more specific than another hypothesis in  $G$

- **Limitations:**
  - $S$  and  $G$  may reduce to  $\emptyset$  with error/noise in training data.
  - $S$  and  $G$  may reduce to  $\emptyset$  with insufficiently expressive hypothesis representation (biased hypothesis space where  $c \notin H$ ).
- Active learner should query input instance that satisfies exactly half of hypotheses in version space, so that version space reduces by half with each training example and requires at least  $\lceil \log_2 |VS_{H,D}| \rceil$  examples to find target concept  $c$ .
- **Inductive Bias:**  $B = \{c \in H\}$

## 2. Decision Tree

### 2.1. Decision Tree Algorithm

- Number of distinct binary decision trees with  $m$  Boolean attributes:  $2^{2^m}$

**function** LEARN(*examples, attributes, parent\_examples*):  
**if** *examples* is empty **then return**  
 PLURALITY\_VALUE(*parent\_examples*)  
**else if** all *examples* have the same classification **then return** the classification  
**else if** *attributes* is empty **then return**  
 PLURALITY\_VALUE(*examples*)  
**else:**  
 $A \leftarrow \text{argmax}_{a \in \text{attributes}} \text{Importance}(a, \text{examples})$   
 $tree \leftarrow$  a new decision tree with root test  $A$   
**for** each value  $v_k$  of  $A$  **do:**  
 $exs \leftarrow \{e: e \in \text{examples} \text{ and } e.A = v_k\}$   
 $subtree \leftarrow \text{Learn}(exs, \text{attributes} - A, \text{examples})$   
 add a branch to  $tree$  with label  $(A = v_k)$  and subtree  $subtree$   
**return tree**

2.2. Entropy and Information Gain

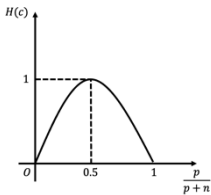
- Entropy: Entropy measures uncertainty of  $c$  in  $\{c_1, \dots, c_k\}$ :

$$H(c) = E(\log_2 P(c_i)) = - \sum_{i=1}^k P(c_i) \log_2 P(c_i)$$

- Entropy of Boolean function that is true with probability  $q$ :  
 $B(q) = -q \log_2 q - (1-q) \log_2 (1-q)$
- Entropy of target concept  $c$  with a training set containing  $p$  +ve examples and  $n$  -ve examples:

$$H(c) = B\left(\frac{p}{p+n}\right) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

For example,  
 if  $p = n \neq 0$ ,  $H(c) = 1$  (maximum uncertainty)  
 if  $p = 0, n \neq 0$ ,  $H(c) = 0$  (no uncertainty)



- Information Gain: Expected reduction in entropy from the attribute test on  $A$ :

$$\text{Gain}(c, A) = B\left(\frac{p}{p+n}\right) - H(c|A)$$

- Expected remaining entropy after testing  $A$ :

$$H(c|A) = E\left(B\left(\frac{p_i}{p_i+n_i}\right)\right) = \sum_{i=1}^d \frac{p_i+n_i}{p+n} B\left(\frac{p_i}{p_i+n_i}\right)$$

[0,1]

- Continuous-valued attributes: Partition into a discrete set of intervals.
- Attributes with many values: Use *GainRatio*:

$$\text{GainRatio}(C, A) = \frac{\text{Gain}(C, A)}{\text{SplitInformation}(C, A)}$$

$$\text{SplitInformation}(C, A) = - \sum_{i=1}^d \frac{|E_i|}{|E|} \log_2 \frac{|E_i|}{|E|}$$

- where  $E_i$  is the subsets of  $E$  divided by  $A$ .
- Attributes with differing costs: Use 
$$\frac{\text{Gain}^2(C, A)}{\text{Cost}(A)}$$
  

$$\frac{2^{\text{Gain}(C, A)} - 1}{(\text{Cost}(A) + 1)^w}$$
 where  $w \in [0,1]$  denotes importance of cost.
- Attributes with missing values:
  - Assign most common value
  - Assign most common value with same value of target concept
  - Assign probability  $p_i$  to each possible values
- Inductive Bias of Decision Tree Learning:
  - (a) Shorter trees are preferred.
  - (b) Trees that place high information gain attributes close to the root are preferred.
    - Occam's Razor: Prefer shortest/simplest hypothesis that fits the data.

- cases for**  $\rightarrow$ 
  - Short hypothesis unlikely to be coincidence
  - Many ways to define small set of hypotheses
- cases against**  $\rightarrow$ 
  - Can be obtained using different hypothesis representations

2.4. Overfitting

- Hypothesis  $h \in H$  overfits the set  $D$  of training examples if and only if 
$$\exists h' \in H \setminus \{h\} \text{ s.t. } (error_D(h) < error_D(h')) \wedge (error_{D_X}(h) > error_{D_X}(h'))$$
 where  $error_D(h)$  and  $error_{D_X}(h)$  denotes errors of  $h$  over  $D$  and set  $D_X$  of examples corresponding to instance space  $X$ .
- How to select best decision tree?
  - Measure performance over training data.
  - Measure performance over a separate validation data set.
  - Minimise  $size(tree)$  &  $size(misclassifications(tree))$
- Avoiding overfitting:
  - Stop growing decision tree when expanding a node is not statistically significant.
  - Allow decision tree to grow and overfit the data, then post-prune it.
- Pruning:
  - Reduced-Error Pruning

Partition data into training and validation sets  
 Do until further pruning is harmful:  
 1. Evaluate impact on validation set of pruning each possible node

2. Greedily remove the one that most improves validation set accuracy

- Rule Post-Pruning

- Convert learning decision tree to an equivalent set of rules
- Prune each rule by removing any precondition that improves its estimated accuracy
- Sort pruned rules by estimated accuracy into desired sequence for use when classifying unobserved input instances

3. Neural Network

3.1. Perceptron Training Rule **bias weight makes hypothesis space more general**

- Perceptron Unit: 
$$o(x_1, \dots, x_n) = \begin{cases} 1, & \text{if } w_0 + w_1x_1 + \dots + w_nx_n > 0 \\ -1, & \text{otherwise} \end{cases}$$
 **orthogonal**

$w_0 + w_1x_1 + \dots + w_nx_n$  can also be written as  $w \cdot x$ .

- Perceptron Training Rule: Initialise  $w$  randomly, apply perceptron training rule to every training example, and iterate through all training examples till  $w$  is consistent: 
$$w_i \leftarrow w_i + \Delta w_i, \Delta w_i = \eta(t - o)x_i$$
 for  $i = 0, 1, 2, \dots, n$  where
  - $t = c(x)$  is the target output;
  - $o = o(x)$  is the perceptron output;
  - $\eta$  is some small +ve constant called learning rate.

It is guaranteed to converge if training examples are linearly separable and  $\eta$  is sufficiently small.

3.2. Linear Unit Training Rule with Gradient Descent

- Loss function:  $L_D(w) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$
- Training rule: Initialise  $w$  randomly, and then repeatedly updating it in the direction of steepest descent:

$$w \leftarrow w + \Delta w, \Delta w = -\eta \nabla L_D(w)$$

$$\text{Gradient: } \nabla L_D(w) = \left[ \frac{\partial L_D}{\partial w_0}, \frac{\partial L_D}{\partial w_1}, \dots, \frac{\partial L_D}{\partial w_n} \right]$$

- Initialise each  $w_i$  to some small random value
- Until termination condition is met, do:
  - Initialise each  $\Delta w_i$  to zero
  - For each  $d \in D$ , do:
    - Input instance  $x_d$  to linear unit and compute  $o$
    - For each linear unit weight  $w_i$ , do 
$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_{id}$$
  - For each linear unit weight  $w_i$ , do 
$$w_i \leftarrow w_i + \Delta w_i$$

**decision surface is linear  $\rightarrow$  multi-layer**