

---

## An Optimization Approach to Learning Falling Rule Lists (Supplementary Material)

---

### 8 Algorithm FRL

In this section, we present Algorithm FRL in detail. Given an instance  $(D, A, w, C)$  of Program 2.9, the algorithm searches through the space of falling rule lists that are compatible with  $D$  and outputs a compatible falling rule list that respects the constraints of Program 2.9, and whose objective value is the smallest among all the falling rule lists that the algorithm explores. It does so by iterating over  $T$  steps, in each of which the algorithm constructs a compatible falling rule list  $d$ , while keeping track of the falling rule list  $d^*$  that has the smallest objective value  $L_{\text{best}} = L(d^*, D, 1/(1+w), w, C)$  among all the falling rule lists that the algorithm has constructed so far. At the end of  $T$  iterations, the algorithm outputs the falling rule list that has the smallest objective value out of the  $T$  lists it has constructed.

In the process of constructing a falling rule list  $d$ , the algorithm chooses the antecedents successively: first for the antecedent  $a_0^{(d)}$  in the top rule, then for the antecedent  $a_1^{(d)}$  in the next rule, and so forth. For each antecedent  $a_j^{(d)}$  chosen, the algorithm also computes its empirical positive proportion  $\alpha_j^{(d,D)}$ . After  $p$  rules have been constructed so that  $d$  currently holds the prefix  $e = \{(a_0^{(d)}, \alpha_0^{(d,D)}), (a_1^{(d)}, \alpha_1^{(d,D)}), \dots, (a_{p-1}^{(d)}, \alpha_{p-1}^{(d,D)})\}$ , the algorithm either: (1) terminates the construction of  $d$  by computing the empirical positive proportion after  $e$ ,  $\tilde{\alpha}_{e,D}$ , and then adding to  $d$  the final else clause with probability estimate  $\tilde{\alpha}_{e,D}$ , or (2) randomly picks an antecedent from a candidate set  $S$  of possible next antecedents, computes its empirical positive proportion, and uses these as the next rule  $(a_p^{(d)}, \alpha_p^{(d,D)})$  for  $d$ .

The algorithm uses various properties of Program 2.9, which are presented in Section 4, to prune the search space. More specifically, the algorithm terminates the construction of  $d$  if Inequality (9) in Theorem 4.6 holds. Otherwise it either terminates the construction of  $d$  with some probability, or proceeds to construct a candidate set  $S$  of possible next antecedents, as follows. For every antecedent  $A_l \in A$  that has not been chosen before, it constructs a candidate next rule  $(a_p^{(d)}, \alpha_p^{(d,D)})$  by setting  $a_p^{(d)} = A_l$  and computing  $\alpha_p^{(d,D)}$  using Definition 2.5. The algorithm then checks if the monotonicity constraint  $\alpha_p^{(d,D)} \leq \alpha_{p-1}^{(d,D)}$  and the necessary condition for optimality  $\alpha_p^{(d,D)} > 1/(1+w)$  (Corollary 4.5) are satisfied, if the prefix  $e' = \{e, (a_p^{(d)}, \alpha_p^{(d,D)})\}$  is feasible under Program 2.9 (i.e. whether there exists a compatible falling rule list that begins with the prefix  $e'$ ) using Proposition 4.2, and if the best possible objective value  $L^*(e', D, w, C)$  achievable by any falling rule list that begins with  $e'$  and is compatible with  $D$  (Theorem 4.6) is less than the current best objective value  $L_{\text{best}} = L(d^*, D, 1/(1+w), w, C)$ . If all of the above conditions are satisfied, the algorithm adds  $A_l$  to  $S$ . Once the construction of  $S$  is complete, the algorithm randomly chooses an antecedent  $A_l \in S$  with probability  $P(A_l|S, e, D)$  and uses this antecedent, together with its empirical positive proportion, as the next rule  $(a_p^{(d)}, \alpha_p^{(d,D)})$  for  $d$ . If  $S$  is empty, the algorithm terminates the construction of  $d$ .

In practice, we define the probability  $P(A_l|S, e, D)$  for  $A_l \in S$  by first defining a curiosity function  $f_{S,e,D} : S \rightarrow \mathbb{R}_{\geq 0}$  and then normalizing it:

$$P(A_l|S, e, D) = \frac{f_{S,e,D}(A_l)}{\sum_{A_{l'}} f_{S,e,D}(A_{l'})}.$$

A possible choice of the curiosity function  $f_{S,e,D}$  for use in Algorithm FRL is given by

$$f_{S,e,D}(A_l) = \lambda \alpha(A_l, e, D) + (1 - \lambda) \frac{n^+(A_l, e, D)}{\tilde{n}_{e,D}^+}, \quad (11)$$

where  $\alpha(A_l, e, D)$  is the empirical positive proportion of  $A_l$ , and  $n^+(A_l, e, D)$  is the number of positive training inputs captured by  $A_l$ , should  $A_l$  be chosen as the next antecedent after the prefix  $e$ . The curiosity function  $f_{S, e, D}$  given by (11) is a weighted sum of  $\alpha(A_l, e, D)$  and  $n^+(A_l, e, D)/\tilde{n}_{e, D}^+$  for each  $A_l \in S$ : the former encourages the algorithm to choose antecedents that have large empirical positive proportions, and the latter encourages the algorithm to choose antecedents that have large positive supports in the training data not captured by  $e$ . We used this curiosity function for Algorithm FRL in our experiments.

The pseudocode of Algorithm FRL is shown in Algorithm 1.

```

Input: an instance  $(D, A, w, C)$  of Program 2.9
Result: a falling rule list  $d^*$  that are compatible with  $D$  and whose antecedents come from  $A$ 
initialize  $d^* = \emptyset$ ,  $L_{\text{best}} = \infty$ ;
for  $t = 1, \dots, T$  do
    set  $p = -1$ ,  $\alpha_p = 1$ ,  $d = e = \emptyset$ ;
    while Inequality (9) in Theorem 4.6 does not hold do
        go to Terminate with some probability;
        set  $p = p + 1$ ,  $S = \emptyset$ ;
        for every antecedent  $A_l \in A$  that is not in  $d$  do
            set  $a_p^{(d)} = A_l$ , compute  $\alpha_p^{(d, D)}$ , and let  $e' = \{e, (a_p^{(d)}, \alpha_p^{(d, D)})\}$ ;
            if  $\alpha_p^{(d, D)} \leq \alpha_{p-1}^{(d, D)}$ ,  $\alpha_p^{(d, D)} > 1/(1+w)$ , and  $e'$  is feasible under Program 2.9 then
                compute  $L^*(e', D, w, C)$  using Theorem 4.6;
                if  $L^*(e', D, w, C) < L(d^*, D, 1/(1+w), w, C)$  then
                    add  $A_l$  to  $S$ ;
                end
            end
        end
        end
        if  $S \neq \emptyset$  then
            choose an antecedent  $A_l \in S$  with probability  $P(A_l|S, e, D)$  according to a discrete probability distribution over  $S$ ;
            set  $a_p^{(d)} = A_l$  and add  $(a_p^{(d)}, \alpha_p^{(d, D)})$  to  $d$ ;
            set  $e = d$ ;
            // save the partially constructed list  $d$  as the prefix  $e$ 
        else
            go to Terminate
        end
    end
    Terminate: terminate the construction of  $d$ , and compute  $L(d, D, 1/(1+w), w, C)$ ;
    if  $L(d, D, 1/(1+w), w, C) < L_{\text{best}}$  then
        set  $d^* = d$ ,  $L_{\text{best}} = L(d, D, 1/(1+w), w, C)$ ;
    end
end

```

Algorithm 1: Algorithm FRL

## 9 Algorithm softFRL

In this section, we present Algorithm softFRL in detail. Given an instance  $(D, A, w, C, C_1)$  of Program 5.1, the algorithm searches through the space of rule lists that are compatible with  $D$  and finds a compatible rule list whose antecedents come from  $A$ , and whose objective value is the smallest among all the rule lists that the algorithm explores. It does so by iterating over  $T$  steps, in each of which the algorithm constructs a compatible rule list  $d$ , while keeping track of the rule list  $d^*$  that has the smallest objective value  $\tilde{L}_{\text{best}} = \tilde{L}(d^*, D, 1/(1+w), w, C, C_1)$  among all the rule lists that the algorithm has constructed so far. At the end of  $T$  iterations, the algorithm transforms the rule list  $d^*$  that has the smallest objective value out of

the  $T$  lists it has constructed, into a falling rule list by setting  $\hat{\alpha}_j^{(d^*)} = \min_{k \leq j} \alpha_k^{(d^*, D)}$ .

In the process of constructing a rule list  $d$ , the algorithm chooses the antecedents successively: first for the antecedent  $a_0^{(d)}$  in the top rule, then for the antecedent  $a_1^{(d)}$  in the next rule, and so forth. For each antecedent  $a_j^{(d)}$  chosen, the algorithm also computes its empirical positive proportion  $\alpha_j^{(d, D)}$ . After  $p$  rules have been constructed so that  $d$  currently holds the prefix  $e = \{(a_0^{(d)}, \alpha_0^{(d, D)}), (a_1^{(d)}, \alpha_1^{(d, D)}), \dots, (a_{p-1}^{(d)}, \alpha_{p-1}^{(d, D)})\}$ , the algorithm either: (1) terminates the construction of  $d$  by computing the empirical positive proportion after  $e$ ,  $\tilde{\alpha}_{e, D}$ , and then adding to  $d$  the final else clause with probability estimate  $\tilde{\alpha}_{e, D}$ , or (2) randomly picks an antecedent from a candidate set  $S$  of possible next antecedents, computes its empirical positive proportion, and use these as the next rule  $(a_p^{(d)}, \alpha_p^{(d, D)})$  for  $d$ .

The algorithm uses Theorem 5.2 to prune the search space. More specifically, the algorithm terminates the construction of  $d$  if  $\tilde{L}^*(e, D, w, C, C_1)$  defined by Equation (10) in Theorem 5.2 is equal to  $\tilde{L}(\bar{e}, D, 1/(1+w), w, C, C_1)$ , where  $\bar{e} = \{e, \tilde{\alpha}_{e, D}\}$  is the compatible rule list in which the prefix  $e$  is followed directly by the final else clause. The condition  $\tilde{L}^*(e, D, w, C, C_1) = \tilde{L}(\bar{e}, D, 1/(1+w), w, C, C_1)$  implies that  $\bar{e}$  is an optimal compatible rule list that begins with  $e$ . If we have  $\tilde{L}^*(e, D, w, C, C_1) < \tilde{L}(\bar{e}, D, 1/(1+w), w, C, C_1)$  instead, the algorithm either terminates the construction of  $d$  with some probability, or it proceeds to construct a candidate set  $S$  of possible next antecedents, as follows. For every antecedent  $A_l \in A$  that has not been chosen before, it constructs a candidate next rule  $(a_p^{(d)}, \alpha_p^{(d, D)})$  by setting  $a_p^{(d)} = A_l$  and computing  $\alpha_p^{(d, D)}$  using Definition 2.5. The algorithm then checks if the best possible objective value  $\tilde{L}^*(e', D, w, C, C_1)$  achievable by any rule list that begins with  $e' = \{e, (a_p^{(d)}, \alpha_p^{(d, D)})\}$  and is compatible with  $D$  (Theorem 5.2) is less than the current best objective value  $\tilde{L}_{\text{best}} = \tilde{L}(d^*, D, 1/(1+w), w, C, C_1)$ . If so, the algorithm adds  $A_l$  to  $S$ . Once the construction of  $S$  is complete, the algorithm randomly chooses an antecedent  $A_l \in S$  with probability  $P(A_l|S, e, D)$  and uses this antecedent, together with its empirical positive proportion, as the next rule  $(a_p^{(d)}, \alpha_p^{(d, D)})$  for  $d$ . If  $S$  is empty, the algorithm terminates the construction of  $d$ .

In practice, we define the probability  $P(A_l|S, e, D)$  for  $A_l \in S$  by first defining a curiosity function  $f_{S, e, D} : S \rightarrow \mathbb{R}_{\geq 0}$  and then normalizing it:

$$P(A_l|S, e, D) = \frac{f_{S, e, D}(A_l)}{\sum_{A_{l'}} f_{S, e, D}(A_{l'})}.$$

A possible choice of the curiosity function  $f_{S, e, D}$  for use in Algorithm softFRL is given by

$$f_{S, e, D}(A_l) = \lambda [\min(\alpha(A_l, e, D), \frac{1.01}{0.01} \alpha_{\min}^{(e, D)} - \frac{1}{0.01} \alpha(A_l, e, D))]_+ + (1 - \lambda) \frac{n^+(A_l, e, D)}{\tilde{n}_{e, D}^+}, \quad (12)$$

where  $\alpha_{\min}^{(e, D)} = \min_{k < |e|} \alpha_k^{(e, D)}$  is the minimum empirical positive proportion of the antecedents in the prefix  $e$ ,  $\alpha(A_l, e, D)$  is the empirical positive proportion of  $A_l$ , and  $n^+(A_l, e, D)$  is the number of positive training inputs captured by  $A_l$ , should  $A_l$  be chosen as the next antecedent after the prefix  $e$ . The curiosity function  $f_{S, e, D}$  given by (12) is a weighted sum of  $[\min(\alpha(A_l, e, D), (1.01/0.01)\alpha_{|e|-1}^{(e, D)} - (1/0.01)\alpha(A_l, e, D))]_+$  and  $n^+(A_l, e, D)/\tilde{n}_{e, D}^+$  for each  $A_l \in S$ : the former encourages the algorithm to choose antecedents that have large empirical positive proportions but do not violate the monotonicity constraint  $\alpha(A_l, e, D) \leq \alpha_{\min}^{(e, D)}$  by more than 1%, and the latter encourages the algorithm to choose antecedents that have large positive supports in the training data not captured by  $e$ . We used this curiosity function for Algorithm softFRL in our experiments.

The pseudocode of Algorithm softFRL is shown in Algorithm 2.

## 10 Proofs of Theorem 2.8, Proposition 4.2, Lemma 4.4, Corollary 4.5, and Theorem 4.6

Theorem 2.8. Given the training data  $D$ , a rule list  $d$  that is compatible with  $D$ , and the weight  $w$  for the positive class, we have

$$R(d, D, 1/(1+w), w) \leq R(d, D, \tau, w)$$

---

Input: an instance  $(D, A, w, C, C_1)$  of Program 5.1  
Result: a falling rule list  $d^*$  whose antecedents come from  $A$   
initialize  $d^* = \emptyset, \tilde{L}_{\text{best}} = \infty$ ;  
for  $t = 1, \dots, T$  do  
    set  $p = -1, \alpha_p = 1, d = e = \emptyset$ ;  
    while  $\tilde{L}^*(e, D, w, C, C_1) < \tilde{L}(\bar{e}, D, 1/(1+w), w, C, C_1)$  do  
        go to Terminate with some probability;  
        set  $p = p + 1, S = \emptyset$ ;  
        for every antecedent  $A_l \in A$  that is not in  $d$  do  
            set  $a_p^{(d)} = A_l$ , compute  $\alpha_p^{(d,D)}$ , and let  $e' = \{e, (a_p^{(d)}, \alpha_p^{(d,D)})\}$ ;  
            compute  $\tilde{L}^*(e', D, w, C, C_1)$  using Theorem 5.2;  
            if  $\tilde{L}^*(e', D, w, C, C_1) < \tilde{L}(d^*, D, 1/(1+w), w, C, C_1)$  then  
                | add  $A_l$  to  $S$ ;  
            end  
        end  
        if  $S \neq \emptyset$  then  
            choose an antecedent  $A_l \in S$  with probability  $P(A_l|S, e, D)$  according to a discrete probability distribution over  $S$ ;  
            set  $a_p^{(d)} = A_l$  and add  $(a_p^{(d)}, \alpha_p^{(d,D)})$  to  $d$ ;  
            set  $e = d$ ;  
            // save the partially constructed list  $d$  as the prefix  $e$   
        else  
            | go to Terminate  
        end  
    end  
    Terminate: terminate the construction of  $d$ , and compute  $\tilde{L}(d, D, 1/(1+w), w, C, C_1)$ ;  
    if  $\tilde{L}(d, D, 1/(1+w), w, C, C_1) < \tilde{L}_{\text{best}}$  then  
        | set  $d^* = d, \tilde{L}_{\text{best}} = \tilde{L}(d, D, 1/(1+w), w, C, C_1)$ ;  
    end  
end  
transform  $d^*$  into a falling rule list by setting  $\hat{\alpha}_j^{(d^*)} = \min_{k \leq j} \alpha_k^{(d^*, D)}$ ;  
Algorithm 2: Algorithm softFRL

for all  $\tau \geq 0$ .

*Proof.* Suppose  $\tau > 1/(1+w)$ . Consider the  $j$ -th rule  $(a_j^{(d)}, \alpha_j^{(d,D)})$  in  $d$ , whose antecedent captures  $\alpha_j^{(d,D)} n_{j,d,D}$  positive training inputs and  $(1 - \alpha_j^{(d,D)}) n_{j,d,D}$  negative training inputs. Let  $R_j(d, D, \tau, w)$  denote the contribution by the  $j$ -th rule to  $R(d, D, \tau, w)$ , i.e.

$$R_j(d, D, \tau, w) = \frac{1}{n} \left( w \sum_{\substack{i: y_i=1 \wedge \\ \text{capt}(\mathbf{x}_i, d)=j}} [\alpha_j^{(d,D)} \leq \tau] + \sum_{\substack{i: y_i=-1 \wedge \\ \text{capt}(\mathbf{x}_i, d)=j}} [\alpha_j^{(d,D)} > \tau] \right) = \begin{cases} \frac{1}{n} n_{j,d,D}^- & \text{if } \alpha_j^{(d,D)} > \tau \\ \frac{w}{n} n_{j,d,D}^+ & \text{otherwise.} \end{cases} \quad (13)$$

Case 1.  $1/(1+w) < \alpha_j^{(d,D)} \leq \tau$ . In this case, we have

$$\begin{aligned} R_j(d, D, 1/(1+w), w) &= \frac{1}{n} n_{j,d,D}^- \quad (\text{by the definition of } R_j \text{ in Equation (13)}) \\ &= \frac{1}{n} (n_{j,d,D} - n_{j,d,D}^+) \quad (\text{by the definition of } n_{j,d,D}^+, n_{j,d,D}^-, n_{j,d,D} \text{ in Definition 2.5}) \\ &= \frac{1}{n} (n_{j,d,D} - \alpha_j^{(d,D)} n_{j,d,D}) \quad (\text{by the definition of } \alpha_j^{(d,D)} \text{ in Definition 2.5}) \\ &= \frac{1}{n} (1 - \alpha_j^{(d,D)}) n_{j,d,D} \\ &< \frac{1}{n} \left( 1 - \frac{1}{1+w} \right) n_{j,d,D} \\ &= \frac{w}{n} \frac{1}{1+w} n_{j,d,D} \\ &< \frac{w}{n} \alpha_j^{(d,D)} n_{j,d,D} \\ &= \frac{w}{n} n_{j,d,D}^+ \quad (\text{by the definition of } \alpha_j^{(d,D)} \text{ in Definition 2.5}) \\ &= R_j(d, D, \tau, w). \quad (\text{by the definition of } R_j \text{ in Equation (13)}) \end{aligned}$$

Case 2.  $\alpha_j^{(d,D)} > \tau$ . In this case, both  $R_j(d, D, 1/(1+w), w)$  and  $R_j(d, D, \tau, w)$  are equal to  $\frac{1}{n} n_{j,d,D}^-$ .

Case 3.  $\alpha_j^{(d,D)} \leq 1/(1+w)$ . In this case, both  $R_j(d, D, 1/(1+w), w)$  and  $R_j(d, D, \tau, w)$  are equal to  $\frac{w}{n} n_{j,d,D}^+$ .

Hence, given  $\tau > 1/(1+w)$ , we have

$$R(d, D, 1/(1+w), w) = \sum_{j=0}^{|d|} R_j(d, D, 1/(1+w), w) \leq \sum_{j=0}^{|d|} R_j(d, D, \tau, w) = R(d, D, \tau, w).$$

The proof for  $R(d, D, 1/(1+w), w) \leq R(d, D, \tau, w)$  given  $\tau < 1/(1+w)$  is similar.  $\square$

**Proposition 4.2.** Given the training data  $D$ , the set of antecedents  $A$ , and a prefix  $e$  that is compatible with  $D$  and satisfies  $a_j^{(e)} \in A$  for all  $j \in \{0, 1, \dots, |e| - 1\}$  and  $\alpha_k^{(e,D)} \geq \alpha_k^{(e,D)}$  for all  $k \in \{1, 2, \dots, |e| - 1\}$ , the following statements are equivalent: (1)  $e$  is feasible for Program 2.9 under  $D$  and  $A$ ; (2)  $\tilde{\alpha}_{e,D} \leq \alpha_{|e|-1}^{(e,D)}$  holds; (3)  $\tilde{n}_{e,D}^- \geq ((1/\alpha_{|e|-1}^{(e,D)}) - 1) \tilde{n}_{e,D}^+$  holds.

*Proof.* (1)  $\Rightarrow$  (3): Suppose that Statement (1) holds. Then there exists a falling rule list

$$d = \{e, (a_{|e|}^{(d)}, \alpha_{|e|}^{(d,D)}), \dots, (a_{|d|-1}^{(d)}, \alpha_{|d|-1}^{(d,D)}), \alpha_{|d|}^{(d,D)}\}$$

that is compatible with  $D$ , and we have

$$\begin{aligned}
\tilde{n}_{e,D}^- &= \tilde{n}_{e,D} - \tilde{n}_{e,D}^+ \\
&= n_{|e|,d,D} + \dots + n_{|d|,d,D} - \tilde{n}_{e,D}^+ \\
&= \frac{1}{\alpha_{|e|}^{(d,D)}} n_{|e|,d,D}^+ + \dots + \frac{1}{\alpha_{|d|}^{(d,D)}} n_{|d|,d,D}^+ - \tilde{n}_{e,D}^+ \quad (\text{by Definition 2.5}) \\
&\geq \frac{1}{\alpha_{|e|-1}^{(d,D)}} n_{|e|,d,D}^+ + \dots + \frac{1}{\alpha_{|e|-1}^{(d,D)}} n_{|d|,d,D}^+ - \tilde{n}_{e,D}^+ \quad (\text{by the monotonicity constraint}) \\
&= \frac{1}{\alpha_{|e|-1}^{(d,D)}} (n_{|e|,d,D}^+ + \dots + n_{|d|,d,D}^+) - \tilde{n}_{e,D}^+ \\
&= \frac{1}{\alpha_{|e|-1}^{(d,D)}} \tilde{n}_{e,D}^+ - \tilde{n}_{e,D}^+ \\
&= ((1/\alpha_{|e|-1}^{(d,D)}) - 1) \tilde{n}_{e,D}^+ \\
&= ((1/\alpha_{|e|-1}^{(e,D)}) - 1) \tilde{n}_{e,D}^+.
\end{aligned}$$

(3)  $\Rightarrow$  (2): Suppose that Statement (3) holds. Then we have

$$\begin{aligned}
\tilde{\alpha}_{e,D} &= \frac{\tilde{n}_{e,D}^+}{\tilde{n}_{e,D}} \quad (\text{by Definition 2.5}) \\
&= \frac{\tilde{n}_{e,D}^+}{\tilde{n}_{e,D}^+ + \tilde{n}_{e,D}^-} \\
&\leq \frac{\tilde{n}_{e,D}^+}{\tilde{n}_{e,D}^+ + ((1/\alpha_{|e|-1}^{(d,D)}) - 1) \tilde{n}_{e,D}^+} \quad (\text{by Statement (3)}) \\
&= \frac{\tilde{n}_{e,D}^+}{(1 + (1/\alpha_{|e|-1}^{(d,D)}) - 1) \tilde{n}_{e,D}^+} = \alpha_{|e|-1}^{(d,D)}.
\end{aligned}$$

(2)  $\Rightarrow$  (1): Suppose that Statement (2) holds. Then the falling rule list  $d = \{e, \tilde{\alpha}_{e,D}\}$  begins with  $e$  and is compatible with  $D$ . By Definition 4.1,  $e$  is feasible for Program 2.9 under the training data  $D$ .  $\square$

Before we proceed with proving Lemma 4.4, we make the following observation.

Observation 10.1 For any rule list

$$d' = \{e, (a_{|e|}^{(d')}, \hat{\alpha}_{|e|}^{(d')}), \dots, (a_{|d'|-1}^{(d')}, \hat{\alpha}_{|d'|-1}^{(d')}), \hat{\alpha}_{|d'|}^{(d')}\}$$

that begins with a given prefix  $e$ , we have

$$\tilde{n}_{e,D}^+ = n_{|e|,d',D}^+ + \dots + n_{|d'|,d',D}^+, \quad (14)$$

$$\tilde{n}_{e,D}^- = n_{|e|,d',D}^- + \dots + n_{|d'|,d',D}^-, \quad (15)$$

and

$$\tilde{n}_{e,D} = n_{|e|,d',D} + \dots + n_{|d'|,d',D}. \quad (16)$$

*Proof.* Any positive training input  $\mathbf{x}_i$  that is not captured by the prefix  $e$  must be captured by some antecedent  $a_j^{(d')}$  with  $|e| \leq j < |d'|$  in  $d'$ , or the final else clause in  $d'$ . Conversely, any positive training input  $\mathbf{x}_i$  that is captured by some antecedent  $a_j^{(d')}$  with  $|e| \leq j < |d'|$  in  $d'$ , or the final else clause in  $d'$ , must not satisfy

any antecedent in the prefix  $e$  and is consequently not captured by the prefix  $e$ . This means that the set of positive training inputs that are not captured by  $e$  is exactly the set of positive training inputs that are captured by some antecedent  $a_j^{(d')}$  with  $|e| \leq j < |d'|$  in  $d'$ , or the final else clause in  $d'$ . It then follows that these two sets have the same number of elements. The former set has  $\tilde{n}_{e,D}^+$  number of elements, and the latter has  $n_{|e|,d',D}^+ + \dots + n_{|d'|,d',D}^+$  number of elements. This establishes Equation (14).

We can establish Equations (15) and (16) using essentially the same argument.  $\square$

We now prove Lemma 4.4.

Lemma 4.4. Suppose that we are given an instance  $(D, A, w, C)$  of Program 2.9, a prefix  $e$  that is feasible for Program 2.9 under the training data  $D$  and the set of antecedents  $A$ , and a (possibly hypothetical) falling rule list  $d$  that begins with  $e$  and is compatible with  $D$ . Then there exists a falling rule list  $d'$ , possibly hypothetical with respect to  $A$ , such that  $d'$  begins with  $e$ , has at most one more rule (excluding the final else clause) following  $e$ , is compatible with  $D$ , and satisfies

$$L(d', D, 1/(1+w), w, C) \leq L(d, D, 1/(1+w), w, C).$$

Moreover, if either  $\alpha_j^{(d,D)} > 1/(1+w)$  holds for all  $j \in \{|e|, |e|+1, \dots, |d|\}$ , or  $\alpha_j^{(d,D)} \leq 1/(1+w)$  holds for all  $j \in \{|e|, |e|+1, \dots, |d|\}$ , then the falling rule list  $\bar{e} = \{e, \tilde{\alpha}_{e,D}\}$  (i.e. the falling rule list in which the final else clause immediately follows the prefix  $e$ , and the probability estimate of the final else clause is  $\tilde{\alpha}_{e,D}$ ) is compatible with  $D$  and satisfies  $L(\bar{e}, D, 1/(1+w), w, C) \leq L(d, D, 1/(1+w), w, C)$ .

*Proof.* Case 1. There exists some  $k \in \{|e|+1, \dots, |d|\}$  that satisfies  $\alpha_{k-1}^{(d,D)} > 1/(1+w)$  but  $\alpha_k^{(d,D)} \leq 1/(1+w)$ . For any  $j \in \{|e|, \dots, k-1\}$ , we have  $\alpha_j^{(d,D)} > 1/(1+w)$ , and the contribution  $R_j(d, D, 1/(1+w), w)$  by the  $j$ -th rule to  $R(d, D, 1/(1+w), w)$ , defined by Equation (13) with  $\tau = 1/(1+w)$ , is given by

$$R_j(d, D, 1/(1+w), w) = \frac{1}{n} n_{j,d,D}^- \quad (17)$$

For any  $j \in \{k, \dots, |d|\}$ , we have  $\alpha_j^{(d,D)} \leq 1/(1+w)$ , and the contribution  $R_j(d, D, 1/(1+w), w)$  by the  $j$ -th rule to  $R(d, D, 1/(1+w), w)$  is given by

$$R_j(d, D, 1/(1+w), w) = \frac{w}{n} n_{j,d,D}^+ \quad (18)$$

The rest of the proof for this case proceeds in three steps.

Step 1. Construct a hypothetical falling rule list  $d'$  that begins with  $e$ , has exactly one more rule (excluding the final else clause) following  $e$ , and is compatible with  $D$ . In later steps, we shall show that the falling rule list  $d'$  constructed in this step satisfies  $L(d', D, 1/(1+w), w, C) \leq L(d, D, 1/(1+w), w, C)$ .

Let  $d' = \{e, (a_{|e|}^{(d')}, \hat{\alpha}_{|e|}^{(d')}), \hat{\alpha}_{|e|+1}^{(d')}\}$  be the falling rule list of size  $|d'| = |e| + 1$  that is compatible with  $D$ , such that

$$a_{|e|}^{(d')} = a_{|e|}^{(d)} \vee \dots \vee a_{k-1}^{(d)}$$

is the antecedent given by the logical OR's of the antecedents  $a_{|e|}^{(d)}$  through  $a_{k-1}^{(d)}$  in  $d$ .

Step 2. Show that the empirical risk of misclassification by the falling rule list  $d'$  is the same as that by the falling rule list  $d$ .

To see this, we observe that the training instances captured by  $a_{|e|}^{(d')}$  in  $d'$  are exactly those captured by the antecedents  $a_{|e|}^{(d)}$  through  $a_{k-1}^{(d)}$  in  $d$ , and the training instances captured by  $a_{|e|+1}^{(d')}$  (i.e. the final else clause) in  $d'$  are exactly those captured by the antecedents  $a_k^{(d)}$  through  $a_{|d|}^{(d)}$  in  $d$ . This observation implies

$$n_{|e|,d',D}^+ = n_{|e|,d,D}^+ + \dots + n_{k-1,d,D}^+, \quad (19)$$

$$n_{|e|,d',D}^- = n_{|e|,d,D}^- + \dots + n_{k-1,d,D}^-, \quad (20)$$

$$n_{|e|,d',D} = n_{|e|,d,D} + \dots + n_{k-1,d,D}, \quad (21)$$

$$n_{|e|+1,d',D}^+ = n_{k,d,D}^+ + \dots + n_{|d|,d,D}^+, \quad (22)$$

and

$$n_{|e|+1,d',D} = n_{k,d,D} + \dots + n_{|d|,d,D}. \quad (23)$$

Since  $d'$  is compatible with  $D$ , using the definition of a compatible rule list in Definition 2.6 and the definition of the empirical positive proportion in Definition 2.5, together with (19), (21), (22), and (23), we must have

$$\begin{aligned} \hat{\alpha}_{|e|}^{(d')} &= \alpha_{|e|}^{(d',D)} = \frac{n_{|e|,d',D}^+}{n_{|e|,d',D}} = \frac{n_{|e|,d,D}^+ + \dots + n_{k-1,d,D}^+}{n_{|e|,d,D} + \dots + n_{k-1,d,D}} \\ &= \frac{\alpha_{|e|}^{(d,D)} n_{|e|,d,D} + \dots + \alpha_{k-1}^{(d,D)} n_{k-1,d,D}}{n_{|e|,d,D} + \dots + n_{k-1,d,D}} > \frac{1}{1+w}, \end{aligned}$$

and

$$\begin{aligned} \hat{\alpha}_{|e|+1}^{(d')} &= \alpha_{|e|+1}^{(d',D)} = \frac{n_{|e|+1,d',D}^+}{n_{|e|+1,d',D}} = \frac{n_{k,d,D}^+ + \dots + n_{|d|,d,D}^+}{n_{k,d,D} + \dots + n_{|d|,d,D}} \\ &= \frac{\alpha_k^{(d,D)} n_{k,d,D} + \dots + \alpha_{|d|}^{(d,D)} n_{|d|,d,D}}{n_{k,d,D} + \dots + n_{|d|,d,D}} \leq \frac{1}{1+w}. \end{aligned}$$

This means that the contribution  $R_{|e|}(d', D, 1/(1+w), w)$  by the  $|e|$ -th rule to  $R(d', D, 1/(1+w), w)$  is given by

$$R_{|e|}(d', D, 1/(1+w), w) = \frac{1}{n} n_{|e|,d',D}^- = \frac{1}{n} (n_{|e|,d,D}^- + \dots + n_{k-1,d,D}^-),$$

where we have used (20), and the contribution  $R_{|e|+1}(d', D, 1/(1+w), w)$  by the  $(|e|+1)$ -st “rule” (i.e. the final else clause) to  $R(d', D, 1/(1+w), w)$  is given by

$$R_{|e|+1}(d', D, 1/(1+w), w) = \frac{w}{n} n_{|e|+1,d',D}^+ = \frac{w}{n} (n_{k,d,D}^+ + \dots + n_{|d|,d,D}^+),$$

where we have used (22).

It then follows that the empirical risk of misclassification by the rule list  $d'$  is the same as that by the rule list  $d$ :

$$\begin{aligned} &R(d', D, 1/(1+w), w) \\ &= R(e, D, 1/(1+w), w) + R_{|e|}(d', D, 1/(1+w), w) + R_{|e|+1}(d', D, 1/(1+w), w) \\ &= R(e, D, 1/(1+w), w) + \frac{1}{n} (n_{|e|,d,D}^- + \dots + n_{k-1,d,D}^-) + \frac{w}{n} (n_{k,d,D}^+ + \dots + n_{|d|,d,D}^+) \\ &= R(e, D, 1/(1+w), w) + \sum_{j=|e|}^{|d|} R_j(d, D, 1/(1+w), w) \\ &= R(d, D, 1/(1+w), w). \end{aligned} \quad (24)$$

Step 3. Put everything together.

Using (24), together with the observation  $|d'| = |e| + 1 \leq |d|$ , we must also have

$$\begin{aligned} L(d', D, 1/(1+w), w, C) &= R(d', D, 1/(1+w), w) + C|d'| \\ &\leq R(d, D, 1/(1+w), w) + C|d| = L(d, D, 1/(1+w), w, C), \end{aligned}$$



as desired.

Case 2.  $\alpha_j^{(d,D)} > 1/(1+w)$  holds for all  $j \in \{|e|, |e|+1, \dots, |d|\}$ . Then the contribution  $R_j(d, D, 1/(1+w), w)$  by the  $j$ -th rule to  $R(d, D, 1/(1+w), w)$ , for all  $j \in \{|e|, |e|+1, \dots, |d|\}$ , is given by Equation (17). Let  $d' = \{e, \hat{\alpha}_{|e|}^{(d')}\}$  be the falling rule list of size  $|d'| = |e|$  that is compatible with  $D$ . Then the instances captured by  $a_{|e|}^{(d')}$  (i.e. the final else clause) in  $d'$  are exactly those that are not captured by  $e$ , or equivalently, those that are captured by  $a_{|e|}^{(d)}$  through  $a_{|d|}^{(d)}$ . This implies

$$n_{|e|,d',D}^+ = n_{|e|,d,D}^+ + \dots + n_{|d|,d,D}^+ \quad (25)$$

$$n_{|e|,d',D}^- = n_{|e|,d,D}^- + \dots + n_{|d|,d,D}^- \quad (26)$$

and

$$n_{|e|,d',D} = n_{|e|,d,D} + \dots + n_{|d|,d,D}. \quad (27)$$

Since  $d'$  is compatible with  $D$ , using the definition of a compatible rule list in Definition 2.6 and the definition of the empirical positive proportion in Definition 2.5, together with (25) and (27), we must have

$$\begin{aligned} \hat{\alpha}_{|e|}^{(d')} &= \alpha_{|e|}^{(d',D)} = \frac{n_{|e|,d',D}^+}{n_{|e|,d',D}} \\ &= \frac{n_{|e|,d,D}^+ + \dots + n_{|d|,d,D}^+}{n_{|e|,d,D} + \dots + n_{|d|,d,D}} \end{aligned} \quad (28)$$

$$= \frac{\alpha_{|e|}^{(d,D)} n_{|e|,d,D} + \dots + \alpha_{|d|}^{(d,D)} n_{|d|,d,D}}{n_{|e|,d,D} + \dots + n_{|d|,d,D}} > \frac{1}{1+w}. \quad (29)$$

Note that the right-hand side of Equality (28) is equal to  $\tilde{n}_{e,D}^+/\tilde{n}_{e,D} = \tilde{\alpha}_{e,D}$ , by Equations (14) and (16) in Observation 10.1. Therefore, we also have  $\hat{\alpha}_{|e|}^{(d')} = \tilde{\alpha}_{e,D}$ .

Inequality (29) implies that the contribution  $R_{|e|}(d', D, 1/(1+w), w)$  by the  $|e|$ -th ‘‘rule’’ (i.e. the final else clause) to  $R(d', D, 1/(1+w), w)$  is given by

$$R_{|e|}(d', D, 1/(1+w), w) = \frac{1}{n} n_{|e|,d',D}^- = \frac{1}{n} (n_{|e|,d,D}^- + \dots + n_{|d|,d,D}^-),$$

where we have used (26).

It then follows that the empirical risk of misclassification by the rule list  $d'$  is the same as that by the rule list  $d$ :

$$\begin{aligned} &R(d', D, 1/(1+w), w) \\ &= R(e, D, 1/(1+w), w) + R_{|e|}(d', D, 1/(1+w), w) \\ &= R(e, D, 1/(1+w), w) + \frac{1}{n} (n_{|e|,d,D}^- + \dots + n_{|d|,d,D}^-) \\ &= R(e, D, 1/(1+w), w) + \sum_{j=|e|}^{|d|} R_j(d, D, 1/(1+w), w) \\ &= R(d, D, 1/(1+w), w). \end{aligned}$$

Since we clearly have  $|d'| = |e| \leq |d|$ , we must also have

$$\begin{aligned} L(d', D, 1/(1+w), w, C) &= R(d', D, 1/(1+w), w) + C|d'| \\ &\leq R(d, D, 1/(1+w), w) + C|d| = L(d, D, 1/(1+w), w, C), \end{aligned}$$

as desired.

Case 3.  $\alpha_j^{(d,D)} \leq 1/(1+w)$  holds for all  $j \in \{|e|, |e|+1, \dots, |d|\}$ . The proof is similar to Case 2, with  $R_j(d, D, 1/(1+w), w)$  for all  $j \in \{|e|, |e|+1, \dots, |d|\}$  given by Equation (18), the “greater than” in Inequality 29 replaced by “less than or equal to”, and  $R_{|e|}(d', D, 1/(1+w), w)$  given by

$$R_{|e|}(d', D, 1/(1+w), w) = \frac{w}{n} n_{|e|,d',D}^+ = \frac{w}{n} (n_{|e|,d,D}^+ + \dots + n_{|d|,d,D}^+).$$

□

Corollary 4.5. If  $d^*$  is an optimal solution for a given instance  $(D, A, w, C)$  of Program 2.9, then we must have  $\alpha_j^{(d^*,D)} > 1/(1+w)$  for all  $j \in \{0, 1, \dots, |d^*| - 1\}$ .

*Proof.* Suppose that  $d^*$  were an optimal solution for a given instance  $(D, A, w, C)$  of Program 2.9, such that  $\alpha_k^{(d^*,D)} \leq 1/(1+w)$  for some  $k \in \{0, 1, \dots, |d^*| - 1\}$ . Let

$$e = \{(a_0^{(d^*)}, \alpha_0^{(d^*,D)}), \dots, (a_{k-1}^{(d^*)}, \alpha_{k-1}^{(d^*,D)})\}$$

be a prefix consisting of the top  $k$  rules in  $d^*$ . By Lemma 4.4, the falling rule list  $\bar{e} = \{e, \tilde{\alpha}_{e,D}\}$  satisfies  $L(\bar{e}, D, 1/(1+w), w, C) \leq L(d^*, D, 1/(1+w), w, C)$ . In fact, the inequality is strict because the size of  $\bar{e}$  is strictly less than that of  $d^*$ . This contradicts the optimality of  $d^*$ . □

Before we proceed with proving Theorem 4.6, we make two other observations.

Observation 10.2. For any rule list  $d'$ , we have

$$n_{|e|,d',D}^- = \left( \frac{1}{\alpha_{|e|}^{(d',D)}} - 1 \right) n_{|e|,d',D}^+, \quad (30)$$

*Proof.* By Definition 2.5, we have

$$\alpha_{|e|}^{(d',D)} = n_{|e|,d',D}^+ / n_{|e|,d',D}.$$

Since  $n_{|e|,d',D}$  denotes the total number of training inputs captured by the  $|e|$ -th antecedent in  $d'$ , which is exactly the sum of the number of positive training inputs captured by that antecedent (denoted  $n_{|e|,d',D}^+$ ), and the number of negative training inputs captured by the same antecedent (denoted  $n_{|e|,d',D}^-$ ), we have

$$\alpha_{|e|}^{(d',D)} = \frac{n_{|e|,d',D}^+}{n_{|e|,d',D}^+ + n_{|e|,d',D}^-}.$$

The desired equation follows from rearranging the terms. □

Observation 10.3. For any rule list

$$d' = \{e, (a_{|e|}^{(d')}, \hat{\alpha}_{|e|}^{(d')}), \hat{\alpha}_{|e|+1}^{(d')}\}$$

that has exactly one rule (excluding the final else clause) following a given prefix  $e$ , we have

$$n_{|e|+1,d',D}^+ = \tilde{n}_{e,D}^+ - n_{|e|,d',D}^+, \quad (31)$$

$$n_{|e|+1,d',D}^- = \tilde{n}_{e,D}^- - n_{|e|,d',D}^-, \quad (32)$$

and

$$n_{|e|+1,d',D} = \tilde{n}_{e,D} - n_{|e|,d',D}. \quad (33)$$

Note that since  $n_{|e|+1,d',D}^+$ ,  $n_{|e|+1,d',D}^-$ , and  $n_{|e|+1,d',D}$  are non-negative, Equations (63), (64), and (65) imply  $n_{|e|,d',D}^+ \leq \tilde{n}_{e,D}^+$ ,  $n_{|e|,d',D}^- \leq \tilde{n}_{e,D}^-$ , and  $n_{|e|,d',D} \leq \tilde{n}_{e,D}$ .

*Proof.* Applying Observation 10.1 with  $|d'| = |e| + 1$ , we have

$$\tilde{n}_{e,D}^+ = n_{|e|,d',D}^+ + n_{|e|+1,d',D}^+,$$

$$\tilde{n}_{e,D}^- = n_{|e|,d',D}^- + n_{|e|+1,d',D}^-,$$

and

$$\tilde{n}_{e,D} = n_{|e|,d',D} + n_{|e|+1,d',D}.$$

Equations (31), (32), and (33) follow from rearranging the terms in the above equations.  $\square$

We now prove Theorem 4.6.

Theorem 4.6. Suppose that we are given an instance  $(D, A, w, C)$  of Program 2.9 and a prefix  $e$  that is feasible for Program 2.9 under the training data  $D$  and the set of antecedents  $A$ . Then any falling rule list  $d$  that begins with  $e$  and is compatible with  $D$  satisfies

$$L(d, D, 1/(1+w), w, C) \geq L^*(e, D, w, C),$$

where

$$L^*(e, D, w, C) = L(e, D, 1/(1+w), w, C) + \min \left( \frac{1}{n} \left( \frac{1}{\alpha_{|e|-1}^{(e,D)}} - 1 \right) \tilde{n}_{e,D}^+ + C, \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right)$$

is a lower bound on the objective value of any compatible falling rule list that begins with  $e$ , which we call a prefix bound for  $e$ , under the instance  $(D, A, w, C)$  of Program 2.9. Furthermore, if

$$C \geq \min \left( \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right) - \frac{1}{n} \left( \frac{1}{\alpha_{|e|-1}^{(e,D)}} - 1 \right) \tilde{n}_{e,D}^+ \quad (34)$$

holds, then the falling rule list  $\bar{e} = \{e, \tilde{\alpha}_{e,D}\}$  satisfies  $L(\bar{e}, D, 1/(1+w), w, C) = L^*(e, D, w, C)$ .

*Proof.* Let  $\mathcal{F}(\mathcal{X}, D, e)$  be the set of (hypothetical and non-hypothetical) falling rule lists that begin with  $e$  and are compatible with  $D$ , and let  $\mathcal{F}(\mathcal{X}, D, e, k)$  be the subset of  $\mathcal{F}(\mathcal{X}, D, e)$ , consisting of those falling rule lists in  $\mathcal{F}(\mathcal{X}, D, e)$  that have exactly  $k$  rules (excluding the final else clause) following the prefix  $e$ .

Let  $d \in \mathcal{F}(\mathcal{X}, D, e)$ .

Case 1.  $\alpha_{|e|-1}^{(e,D)} > 1/(1+w)$ .

In this case, Lemma 4.4 implies

$$L(d, D, 1/(1+w), w, C) \geq \inf_{d' \in \mathcal{F}(\mathcal{X}, D, e, 1) \cup \mathcal{F}(\mathcal{X}, D, e, 0)} L(d', D, 1/(1+w), w, C). \quad (35)$$

Note that we have  $\mathcal{F}(\mathcal{X}, D, e, 0) = \{\bar{e}\}$ , where  $\bar{e} = \{e, \tilde{\alpha}_{e,D}\}$  is the falling rule list in which the final else clause immediately follows the prefix  $e$ , and the probability estimate of the final else clause is  $\tilde{\alpha}_{e,D}$ . To see this, we first observe  $\bar{e} \in \mathcal{F}(\mathcal{X}, D, e, 0)$ . This is because:

- (i)  $\bar{e}$  clearly begins with  $e$ , and has no additional rules (excluding the final else clause) following the prefix  $e$ ;
- (ii) the feasibility of  $e$  implies  $\alpha_{k-1}^{(e,D)} \geq \alpha_k^{(e,D)}$  for all  $k \in \{1, 2, \dots, |e|-1\}$  (otherwise we could not possibly have a falling rule list that begins with  $e$ , and we would violate Definition 4.1), and  $\tilde{\alpha}_{e,D} \leq \alpha_{|e|-1}^{(e,D)}$  (by Proposition 4.2), which together imply that  $\bar{e}$  is indeed a falling rule list; and
- (iii) we have

$$\begin{aligned} \tilde{\alpha}_{e,D} &= \frac{\tilde{n}_{e,D}^+}{\tilde{n}_{e,D}} \quad (\text{by the definition of } \tilde{\alpha}_{e,D} \text{ in Definition 2.5}) \\ &= \frac{n_{|\bar{e}|, \bar{e}, D}^+}{n_{|\bar{e}|, \bar{e}, D}} \quad (\text{by Equations (14) and (16) in Observation 10.1, applied to } \bar{e}) \\ &= \alpha_{|\bar{e}|}^{(\bar{e}, D)}, \quad (\text{by the definition of the empirical positive proportion in Definition 2.5}) \end{aligned}$$

which implies that  $\bar{e}$  is indeed compatible with  $D$ .

Conversely, for any  $d_0 = \{e, \hat{\alpha}_{|e|}^{(d_0)}\} \in \mathcal{F}(\mathcal{X}, D, e, 0)$ , we must have

$$\begin{aligned} \hat{\alpha}_{|e|}^{(d_0)} &= \alpha_{|e|}^{(d_0, D)} \quad (\text{because } d_0 \text{ must be compatible with } D) \\ &= \frac{n_{|e|, d_0, D}^+}{n_{|e|, d_0, D}} \quad (\text{by the definition of the empirical positive proportion in Definition 2.5}) \\ &= \frac{\tilde{n}_{e, D}^+}{\tilde{n}_{e, D}} \quad (\text{by Equations (14) and (16) in Observation 10.1, applied to } d_0 \text{ here}) \\ &= \tilde{\alpha}_{e, D}, \end{aligned}$$

which implies  $d_0 = \bar{e}$ . This establishes  $\mathcal{F}(\mathcal{X}, D, e, 0) = \{\bar{e}\}$ .

Let  $\mathcal{F}'(\mathcal{X}, D, e, 1)$  be the subset of  $\mathcal{F}(\mathcal{X}, D, e, 1)$ , consisting of those falling rule lists

$$d' = \{e, (a_{|e|}^{(d')}, \alpha_{|e|}^{(d', D)}), \alpha_{|e|+1}^{(d', D)}\} \in \mathcal{F}(\mathcal{X}, D, e, 1)$$

with  $\alpha_{|e|}^{(d', D)} > 1/(1+w)$  and  $\alpha_{|e|+1}^{(d', D)} \leq 1/(1+w)$ . Note that for any  $d_1 = \{e, (a_{|e|}^{(d_1)}, \alpha_{|e|}^{(d_1, D)}), \alpha_{|e|+1}^{(d_1, D)}\} \in \mathcal{F}(\mathcal{X}, D, e, 1) - \mathcal{F}'(\mathcal{X}, D, e, 1)$ , we have either  $\alpha_{|e|}^{(d_1, D)} \geq \alpha_{|e|+1}^{(d_1, D)} > 1/(1+w)$  or  $\alpha_{|e|+1}^{(d_1, D)} \leq \alpha_{|e|}^{(d_1, D)} \leq 1/(1+w)$ , and Lemma 4.4 implies  $L(d_1, D, 1/(1+w), w, C) \geq L(\bar{e}, D, 1/(1+w), w, C)$ . This means

$$\inf_{d' \in \mathcal{F}(\mathcal{X}, D, e, 1) - \mathcal{F}'(\mathcal{X}, D, e, 1)} L(d', D, 1/(1+w), w, C) \geq L(\bar{e}, D, 1/(1+w), w, C). \quad (36)$$

Using  $\mathcal{F}(\mathcal{X}, D, e, 0) = \{\bar{e}\}$  and (36), we can write the right-hand side of (35) as

$$\begin{aligned} &\inf_{d' \in \mathcal{F}(\mathcal{X}, D, e, 1) \cup \mathcal{F}(\mathcal{X}, D, e, 0)} L(d', D, 1/(1+w), w, C) \\ &= \inf_{d' \in \mathcal{F}'(\mathcal{X}, D, e, 1) \cup (\mathcal{F}(\mathcal{X}, D, e, 1) - \mathcal{F}'(\mathcal{X}, D, e, 1)) \cup \{\bar{e}\}} L(d', D, 1/(1+w), w, C) \\ &= \min \left( \inf_{d' \in \mathcal{F}'(\mathcal{X}, D, e, 1)} L(d', D, 1/(1+w), w, C), \right. \\ &\quad \left. \inf_{d' \in \mathcal{F}(\mathcal{X}, D, e, 1) - \mathcal{F}'(\mathcal{X}, D, e, 1)} L(d', D, 1/(1+w), w, C), L(\bar{e}, D, 1/(1+w), w, C) \right) \\ &= \min \left( \inf_{d' \in \mathcal{F}'(\mathcal{X}, D, e, 1)} L(d', D, 1/(1+w), w, C), L(\bar{e}, D, 1/(1+w), w, C) \right). \quad (37) \end{aligned}$$

The rest of the proof for this case proceeds in three steps.

Step 1. Compute  $L(\bar{e}, D, 1/(1+w), w, C)$ .

Since the contribution by the final else clause to  $L(\bar{e}, D, 1/(1+w), w, C)$  is given by

$$R_{|e|}(\bar{e}, D, 1/(1+w), w) = \begin{cases} \frac{1}{n} n_{|e|, \bar{e}, D}^- & \text{if } \tilde{\alpha}_{e, D} > 1/(1+w) \\ \frac{w}{n} n_{|e|, \bar{e}, D}^+ & \text{otherwise,} \end{cases}$$

where we have used Equation (13), and since Observation 10.1 implies  $\tilde{n}_{e, D}^+ = n_{|e|, \bar{e}, D}^+$  and  $\tilde{n}_{e, D}^- = n_{|e|, \bar{e}, D}^-$ , it is not difficult to see

$$L(\bar{e}, D, 1/(1+w), w, C) = \begin{cases} L(e, D, 1/(1+w), w, C) + \frac{1}{n} \tilde{n}_{e, D}^- & \text{if } \tilde{\alpha}_{e, D} > 1/(1+w) \\ L(e, D, 1/(1+w), w, C) + \frac{w}{n} \tilde{n}_{e, D}^+ & \text{otherwise.} \end{cases}$$

Since  $\tilde{\alpha}_{e,D} > 1/(1+w)$  is equivalent to  $\tilde{n}_{e,D}^+ / (\tilde{n}_{e,D}^+ + \tilde{n}_{e,D}^-) > 1/(1+w)$ , or  $w\tilde{n}_{e,D}^+ > \tilde{n}_{e,D}^-$ , and similarly  $\tilde{\alpha}_{e,D} \leq 1/(1+w)$  is equivalent to  $w\tilde{n}_{e,D}^+ \leq \tilde{n}_{e,D}^-$ , we can write

$$L(\bar{e}, D, 1/(1+w), w, C) = \min \left( L(e, D, 1/(1+w), w, C) + \frac{1}{n} \tilde{n}_{e,D}^-, \right. \\ \left. L(e, D, 1/(1+w), w, C) + \frac{w}{n} \tilde{n}_{e,D}^+ \right). \quad (38)$$

Step 2. Determine a lower bound of  $L(d', D, 1/(1+w), w, C)$  for all  $d' \in \mathcal{F}'(\mathcal{X}, D, e, 1)$ .

Let  $d' = \{e, (a_{|e|}^{(d')}, \alpha_{|e|}^{(d',D)}), \alpha_{|e|+1}^{(d',D)}\} \in \mathcal{F}'(\mathcal{X}, D, e, 1)$ . Since the contribution by both the  $|e|$ -th rule and the final else clause to  $L(d', D, 1/(1+w), w, C)$  is given by  $R_{|e|}(d', D, 1/(1+w), w) + R_{|e|+1}(d', D, 1/(1+w), w) + C$ , where  $R_{|e|}(d', D, 1/(1+w), w)$  and  $R_{|e|+1}(d', D, 1/(1+w), w)$  are defined by Equation (13) and are given by

$$R_{|e|}(d', D, 1/(1+w), w) = \frac{1}{n} n_{|e|,d',D}^- \quad \text{and} \quad R_{|e|+1}(d', D, 1/(1+w), w) = \frac{w}{n} n_{|e|+1,d',D}^+$$

(because we have  $\alpha_{|e|}^{(d',D)} > 1/(1+w)$  and  $\alpha_{|e|+1}^{(d',D)} \leq 1/(1+w)$  for  $d' \in \mathcal{F}'(\mathcal{X}, D, e, 1)$ ), it is not difficult to see

$$L(d', D, 1/(1+w), w, C) = L(e, D, 1/(1+w), w, C) + \frac{1}{n} n_{|e|,d',D}^- + \frac{w}{n} n_{|e|+1,d',D}^+ + C. \quad (39)$$

Substituting (30) in Observation 10.2 and (31) in Observation 10.3 into Equation (39), we have

$$L(d', D, 1/(1+w), w, C) \\ = L(e, D, 1/(1+w), w, C) + \frac{1}{n} \left( \frac{1}{\alpha_{|e|}^{(d',D)}} - 1 \right) n_{|e|,d',D}^+ + \frac{w}{n} (\tilde{n}_{e,D}^+ - n_{|e|,d',D}^+) + C \\ = L(e, D, 1/(1+w), w, C) + \frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|}^{(d',D)}} - 1 - w \right) n_{|e|,d',D}^+ + w \tilde{n}_{e,D}^+ \right) + C. \quad (40)$$

Note that Equation (40) shows that given the prefix  $e$ ,  $L(d', D, 1/(1+w), w, C)$  is a function of  $\alpha_{|e|}^{(d',D)}$  and of  $n_{|e|,d',D}^+$ . Since we have

$$\frac{\partial L(d', D, 1/(1+w), w, C)}{\partial n_{|e|,d',D}^+} = \frac{1}{n} \left( \frac{1}{\alpha_{|e|}^{(d',D)}} - 1 - w \right) < 0$$

because  $\alpha_{|e|}^{(d',D)} > 1/(1+w)$  holds for any  $d' \in \mathcal{F}'(\mathcal{X}, D, e, 1)$ , and

$$\frac{\partial L(d', D, 1/(1+w), w, C)}{\partial \alpha_{|e|}^{(d',D)}} = -\frac{n_{|e|,d',D}^+}{n} \frac{1}{(\alpha_{|e|}^{(d',D)})^2} \leq 0,$$

we see that  $L(d', D, 1/(1+w), w, C)$  is indeed a monotonically decreasing function of both  $n_{|e|,d',D}^+$  and  $\alpha_{|e|}^{(d',D)}$ . Thus, we can obtain a lower bound of  $L(d', D, 1/(1+w), w, C)$  by substituting  $n_{|e|,d',D}^+$  and  $\alpha_{|e|}^{(d',D)}$  with their respective upper bound. The inequality  $n_{|e|,d',D}^+ \leq \tilde{n}_{e,D}^+$  in Observation 10.3 gives an upper bound for  $n_{|e|,d',D}^+$ , and the inequality  $\alpha_{|e|}^{(d',D)} \leq \alpha_{|e|-1}^{(d',D)} = \alpha_{|e|-1}^{(e,D)}$  from  $d'$  being a falling rule list gives an upper bound for  $\alpha_{|e|}^{(d',D)}$ . Substituting these upper bounds into (40), we obtain the following inequality, which gives

a lower bound of  $L(d', D, 1/(1+w), w, C)$ :

$$\begin{aligned}
& L(d', D, 1/(1+w), w, C) \\
& \geq L(e, D, 1/(1+w), w, C) + \frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(e,D)}} - 1 - w \right) \tilde{n}_{e,D}^+ + w \tilde{n}_{e,D}^+ \right) + C \\
& = L(e, D, 1/(1+w), w, C) + \frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(e,D)}} - 1 \right) \tilde{n}_{e,D}^+ \right) + C.
\end{aligned}$$

This means

$$\inf_{d' \in \mathcal{F}'(\mathcal{X}, D, e, 1)} L(d', D, 1/(1+w), w, C) \geq L(e, D, 1/(1+w), w, C) + \frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(e,D)}} - 1 \right) \tilde{n}_{e,D}^+ \right) + C. \quad (41)$$

Step 3. Put everything together.

Using (35), (37), (38), and (41), we have

$$\begin{aligned}
& L(d, D, 1/(1+w), w, C) \\
& \geq \min \left( \inf_{d' \in \mathcal{F}'(\mathcal{X}, D, e, 1)} L(d', D, 1/(1+w), w, C), L(\bar{e}, D, 1/(1+w), w, C) \right) \\
& \geq \min \left( L(e, D, 1/(1+w), w, C) + \frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(e,D)}} - 1 \right) \tilde{n}_{e,D}^+ \right) + C, \right. \\
& \quad \left. \min \left( L(e, D, 1/(1+w), w, C) + \frac{1}{n} \tilde{n}_{e,D}^-, L(e, D, 1/(1+w), w, C) + \frac{w}{n} \tilde{n}_{e,D}^+ \right) \right) \\
& = L(e, D, 1/(1+w), w, C) + \min \left( \frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(e,D)}} - 1 \right) \tilde{n}_{e,D}^+ \right) + C, \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right),
\end{aligned}$$

as desired.

Case 2.  $\alpha_{|e|-1}^{(e,D)} \leq 1/(1+w)$ .

This implies  $\alpha_j^{(d,D)} \leq 1/(1+w)$  for all  $j \in \{|e|, \dots, |d|\}$ . By Lemma 4.4, we have

$$L(d, D, 1/(1+w), w, C) \geq L(\bar{e}, D, 1/(1+w), w, C).$$

Since  $L(\bar{e}, D, 1/(1+w), w, C)$  is given by Equation (38), we have

$$L(d, D, 1/(1+w), w, C) \geq L(e, D, 1/(1+w), w, C) + \min \left( \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right). \quad (42)$$

Given  $\alpha_{|e|-1}^{(e,D)} \leq 1/(1+w)$ , we must also have

$$\frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(d',D)}} - 1 \right) \tilde{n}_{e,D}^+ \right) + C \geq \frac{w}{n} \tilde{n}_{e,D}^+ + C \geq \frac{w}{n} \tilde{n}_{e,D}^+,$$

which means

$$\min \left( \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right) = \min \left( \frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(d',D)}} - 1 \right) \tilde{n}_{e,D}^+ \right) + C, \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right). \quad (43)$$

Substituting (43) into (42) completes the proof for Case 2.

Finally, if Inequality (34) holds, then we have

$$\frac{1}{n} \left( \left( \frac{1}{\alpha_{|e|-1}^{(d',D)}} - 1 \right) \tilde{n}_{e,D}^+ \right) + C \geq \min \left( \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right),$$

which implies

$$L^*(e, D, w, C) = L(e, D, 1/(1+w), w, C) + \min \left( \frac{w}{n} \tilde{n}_{e,D}^+, \frac{1}{n} \tilde{n}_{e,D}^- \right) = L(\bar{e}, D, 1/(1+w), w, C).$$

□

## 11 Proof of Theorem 5.2

Theorem 5.2. Suppose that we are given an instance  $(D, A, w, C, C_1)$  of Program 5.1 and a prefix  $e$  that is compatible with  $D$ . Then any rule list  $d$  that begins with  $e$  and is compatible with  $D$  satisfies

$$\tilde{L}(d, D, 1/(1+w), w, C, C_1) \geq \tilde{L}^*(e, D, w, C, C_1),$$

where

$$\begin{aligned} \tilde{L}^*(e, D, w, C, C_1) &= \tilde{L}(e, D, 1/(1+w), w, C, C_1) \\ &+ \min \left( \frac{1}{n} \left( \frac{1}{\alpha_{\min}^{(e,D)}} - 1 \right) \tilde{n}_{e,D}^+ + C + C_1 [\tilde{\alpha}_{e,D} - \alpha_{\min}^{(e,D)}]_+ + \frac{w}{n} \tilde{n}_{e,D}^+ [\tilde{\alpha}_{e,D} \geq \alpha_{\min}^{(e,D)}], \right. \\ &\left. \inf_{\beta: \zeta < \beta \leq 1} g(\beta), \frac{w}{n} \tilde{n}_{e,D}^+ + C_1 [\tilde{\alpha}_{e,D} - \alpha_{\min}^{(e,D)}]_+, \frac{1}{n} \tilde{n}_{e,D}^- + C_1 [\tilde{\alpha}_{e,D} - \alpha_{\min}^{(e,D)}]_+ \right) \end{aligned} \quad (44)$$

is a lower bound on the objective value of any compatible rule list that begins with  $e$ , under the instance  $(D, A, w, C, C_1)$  of Program 5.1. In Equation (44),  $\alpha_{\min}^{(e,D)}$ ,  $\zeta$ , and  $g$  are defined by

$$\alpha_{\min}^{(e,D)} = \min_{k < |e|} \alpha_k^{(e,D)}, \quad \zeta = \max(\alpha_{\min}^{(e,D)}, \tilde{\alpha}_{e,D}, 1/(1+w)),$$

$$g(\beta) = \frac{1}{n} \left( \frac{1}{\beta} - 1 \right) \tilde{n}_{e,D}^+ + C + C_1 (\beta - \alpha_{\min}^{(e,D)}).$$

Note that  $\inf_{\beta: \zeta < \beta \leq 1} g(\beta)$  can be computed analytically:  $\inf_{\beta: \zeta < \beta \leq 1} g(\beta) = g(\beta^*)$  if  $\beta^* = \sqrt{\tilde{n}_{e,D}^+ / (C_1 n)}$  satisfies  $\zeta < \beta^* \leq 1$ , and  $\inf_{\beta: \zeta < \beta \leq 1} g(\beta) = \min(g(\zeta), g(1))$  otherwise.

To prove Theorem 5.2, we need the following lemma:

Lemma. Suppose that we are given an instance  $(D, A, w, C, C_1)$  of Program 5.1, a prefix  $e$  that is compatible with  $D$ , and a (possibly hypothetical) rule list  $d$  that begins with  $e$  and is compatible with  $D$ . Then there exists a rule list  $d'$ , possibly hypothetical with respect to  $A$ , such that  $d'$  begins with  $e$ , has at most one more rule (excluding the final else clause) following  $e$ , is compatible with  $D$ , and satisfies

$$\tilde{L}(d', D, 1/(1+w), w, C, C_1) \leq \tilde{L}(d, D, 1/(1+w), w, C, C_1). \quad (45)$$

Moreover, if either  $\alpha_j^{(d,D)} > 1/(1+w)$  holds for all  $j \in \{|e|, |e|+1, \dots, |d|\}$ , or  $\alpha_j^{(d,D)} \leq 1/(1+w)$  holds for all  $j \in \{|e|, |e|+1, \dots, |d|\}$ , then the rule list  $\bar{e} = \{e, \tilde{\alpha}_{e,D}\}$  (i.e. the rule list in which the final else clause follows immediately the prefix  $e$ , and the probability estimate of the final else clause is  $\tilde{\alpha}_{e,D}$ ) is compatible with  $D$  and satisfies  $\tilde{L}(\bar{e}, D, 1/(1+w), w, C, C_1) \leq \tilde{L}(d, D, 1/(1+w), w, C)$ .

*Proof.* Case 1. There exists some  $k \in \{|e|, \dots, |d|\}$  that satisfies  $\alpha_k^{(d,D)} > 1/(1+w)$  and some  $k' \in \{|e|, \dots, |d|\}$  that satisfies  $\alpha_{k'}^{(d,D)} \leq 1/(1+w)$ . For any  $j \in \{|e|, \dots, |d|\}$  with  $\alpha_j^{(d,D)} > 1/(1+w)$ , the contribution  $R_j(d, D, 1/(1+w), w)$  by the  $j$ -th rule to  $R(d, D, 1/(1+w), w)$ , defined by the right-hand side of Equation (13) with  $\tau = 1/(1+w)$ , is given by

$$R_j(d, D, 1/(1+w), w) = \frac{1}{n} n_{j,d,D}^-.$$

For any  $j \in \{|e|, \dots, |d|\}$  with  $\alpha_j^{(d,D)} \leq 1/(1+w)$ , the contribution  $R_j(d, D, 1/(1+w), w)$  by the  $j$ -th rule to  $R(d, D, 1/(1+w), w)$  is given by

$$R_j(d, D, 1/(1+w), w) = \frac{w}{n} n_{j,d,D}^+.$$

The rest of the proof for this case proceeds in four steps.

Step 1. Construct a hypothetical rule list  $d'$  that begins with  $e$ , has exactly one more rule (excluding the final else clause) following  $e$ , and is compatible with  $D$ . In later steps, we shall show that the rule list  $d'$  constructed in this step satisfies (45).

Let  $d' = \{e, (a_{|e|}^{(d')}, \hat{\alpha}_{|e|}^{(d')}), \hat{\alpha}_{|e|+1}^{(d')}\}$  be the hypothetical rule list of size  $|d'| = |e| + 1$  that is compatible with  $D$ , and whose  $|e|$ -th antecedent  $a_{|e|}^{(d')}$  is defined by

$$a_{|e|}^{(d')}(\mathbf{x}) = [\alpha_{\text{capt}(\mathbf{x}, d)}^{(d,D)} > 1/(1+w)] \cdot [|e| \leq \text{capt}(\mathbf{x}, d) \leq |d|].$$

Step 2. Show that the empirical risk of misclassification by the rule list  $d'$  is the same as that by the rule list  $d$ .

To see this, we observe that the training instances in  $D$  captured by  $a_{|e|}^{(d')}$  in  $d'$  are exactly those captured by the antecedents  $a_j^{(d)}$ ,  $|e| \leq j \leq |d|$ , in  $d$  whose empirical positive proportion satisfies  $\alpha_j^{(d,D)} > 1/(1+w)$ , and the training instances in  $D$  captured by  $a_{|e|+1}^{(d')}$  (i.e. the final else clause) in  $d'$  are exactly those captured by the antecedents  $a_j^{(d)}$ ,  $|e| \leq j \leq |d|$ , in  $d$  whose empirical positive proportion satisfies  $\alpha_j^{(d,D)} \leq 1/(1+w)$ . This observation implies

$$n_{|e|, d', D}^+ = \sum_{j: |e| \leq j \leq |d| \wedge \alpha_j^{(d,D)} > 1/(1+w)} n_{j,d,D}^+, \quad (46)$$

$$n_{|e|, d', D}^- = \sum_{j: |e| \leq j \leq |d| \wedge \alpha_j^{(d,D)} > 1/(1+w)} n_{j,d,D}^-, \quad (47)$$

$$n_{|e|, d', D} = \sum_{j: |e| \leq j \leq |d| \wedge \alpha_j^{(d,D)} > 1/(1+w)} n_{j,d,D}, \quad (48)$$

$$n_{|e|+1, d', D}^+ = \sum_{j: |e| \leq j \leq |d| \wedge \alpha_j^{(d,D)} \leq 1/(1+w)} n_{j,d,D}^+ \quad (49)$$

and

$$n_{|e|+1, d', D} = \sum_{j: |e| \leq j \leq |d| \wedge \alpha_j^{(d,D)} \leq 1/(1+w)} n_{j,d,D}. \quad (50)$$

Since  $d'$  is compatible with  $D$ , using the definition of a compatible rule list in Definition 2.6 and the definition



of the empirical positive proportion in Definition 2.5, together with (46), (48), (49), and (50), we must have

$$\begin{aligned}\hat{\alpha}_{|e|}^{(d')} &= \alpha_{|e|}^{(d',D)} = \frac{n_{|e|,d',D}^+}{n_{|e|,d',D}} = \frac{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}>1/(1+w)} n_{j,d,D}^+}{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}>1/(1+w)} n_{j,d,D}} \\ &= \frac{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}>1/(1+w)} \alpha_j^{(d,D)} n_{j,d,D}}{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}>1/(1+w)} n_{j,d,D}} > \frac{1}{1+w},\end{aligned}$$

and

$$\begin{aligned}\hat{\alpha}_{|e|+1}^{(d')} &= \alpha_{|e|+1}^{(d',D)} = \frac{n_{|e|+1,d',D}^+}{n_{|e|+1,d',D}} = \frac{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}\leq 1/(1+w)} n_{j,d,D}^+}{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}\leq 1/(1+w)} n_{j,d,D}} \\ &= \frac{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}\leq 1/(1+w)} \alpha_j^{(d,D)} n_{j,d,D}}{\sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}\leq 1/(1+w)} n_{j,d,D}} \leq \frac{1}{1+w}.\end{aligned}$$

This means that the contribution  $R_{|e|}(d', D, 1/(1+w), w)$  by the  $|e|$ -th rule to  $R(d', D, 1/(1+w), w)$  is given by

$$R_{|e|}(d', D, 1/(1+w), w) = \frac{1}{n} n_{|e|,d',D}^- = \frac{1}{n} \sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}>1/(1+w)} n_{j,d,D}^-,$$

where we have used (47), and the contribution  $R_{|e|+1}(d', D, 1/(1+w), w)$  by the  $(|e|+1)$ -st “rule” (i.e. the final else clause) to  $R(d', D, 1/(1+w), w)$  is given by

$$R_{|e|+1}(d', D, 1/(1+w), w) = \frac{w}{n} n_{|e|+1,d',D}^+ = \frac{w}{n} \sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}\leq 1/(1+w)} n_{j,d,D}^+,$$

where we have used (49).

It then follows that the empirical risk of misclassification by the rule list  $d'$  is the same as that by the rule list  $d$ :

$$\begin{aligned}&R(d', D, 1/(1+w), w) \\ &= R(e, D, 1/(1+w), w) + R_{|e|}(d', D, 1/(1+w), w) + R_{|e|+1}(d', D, 1/(1+w), w) \\ &= R(e, D, 1/(1+w), w) \\ &\quad + \frac{1}{n} \sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}>1/(1+w)} n_{j,d,D}^- + \frac{w}{n} \sum_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}\leq 1/(1+w)} n_{j,d,D}^+ \\ &= R(e, D, 1/(1+w), w) + \sum_{j=|e|}^{|d|} R_j(d, D, 1/(1+w), w) \\ &= R(d, D, 1/(1+w), w).\end{aligned}\tag{51}$$

Step 3. Show that the monotonicity penalty of the rule list  $d'$  is at most that of  $d$ .

Let  $S(d, D) = \sum_{j=0}^{|d|} [\alpha_j^{(d,D)} - \min_{k<j} \alpha_k^{(d,D)}]_+$  be the monotonicity penalty of the rule list  $d$ . We now show  $S(d', D) \leq S(d, D)$ . Let  $S_j(d, D) = [\alpha_j^{(d,D)} - \min_{k<j} \alpha_k^{(d,D)}]_+$  be the monotonicity penalty for the  $j$ -th rule in  $d$ .

Let  $l \in \{|e|, \dots, |d|\}$  be any integer with

$$\alpha_l^{(d,D)} = \max_{j:|e|\leq j\leq |d|\wedge\alpha_j^{(d,D)}>1/(1+w)} \alpha_j^{(d,D)}.\tag{52}$$

Then the total monotonicity penalty for all the rules  $(a_j^{(d)}; j^{(d;D)})$  in  $d$  with  $j \in j$  and  $j^{(d;D)} > 1=(1+w)$  satisfies

$$\begin{aligned} & \sum_{j: j \in j, j^{(d;D)} > 1=(1+w)} S_j(d; D) - S_j(d; D) \quad (\text{because } S_j(d; D) \text{ is included in the sum on the left}) \\ &= b_l^{(d;D)} \min_{k < l} (d;D)_k c_+ \\ &= b_l^{(d;D)} \min_{k < j_{ej}} (d;D)_k c_+ : \end{aligned} \quad (53)$$

On the other hand, the monotonicity penalty for the  $j_{ej}$ -th rule in  $d^0$  satisfies

$$S_{j_{ej}}(d^0, D) = b_{j_{ej}}^{(d^0;D)} \min_{k < j_{ej}} (d^0;D)_k c_+ - b_l^{(d;D)} \min_{k < j_{ej}} (d;D)_k c_+ ; \quad (54)$$

because we have  $\min_{k < j_{ej}} (d^0;D)_k = \min_{k < j_{ej}} (d;D)_k$  ( $d$  and  $d^0$  begin with the same prefix), and

$$\begin{aligned} \frac{(d^0;D)_{j_{ej}}}{b_{j_{ej}}^{(d^0;D)}} &= \frac{n_{j_{ej};d^0;D}^+}{p_{j_{ej};d^0;D}} \quad (\text{by the definition of the empirical positive proportion in Definition 2.5}) \\ &= \frac{p_{j: j \in j, j^{(d;D)} > 1=(1+w)} n_{j;d;D}^+}{p_{j: j \in j, j^{(d;D)} > 1=(1+w)} n_{j;d;D}} \quad (\text{by Equations (46) and (48)}) \\ &= \frac{p_{j: j \in j, j^{(d;D)} > 1=(1+w)} (d;D)_j n_{j;d;D}}{p_{j: j \in j, j^{(d;D)} > 1=(1+w)} n_{j;d;D}} \quad (\text{by the definition of } (d;D)_j \text{ in Definition 2.5}) \\ &= \frac{p_{j: j \in j, j^{(d;D)} > 1=(1+w)} (d;D)_l n_{j;d;D}}{p_{j: j \in j, j^{(d;D)} > 1=(1+w)} n_{j;d;D}} \quad (\text{by the definition of } l \text{ in (52)}) \\ &= (d;D)_l. \end{aligned}$$

Combining (53) and (54), we have

$$S_{j_{ej}}(d^0, D) - \sum_{j: j \in j, j^{(d;D)} > 1=(1+w)} S_j(d; D) : \quad (55)$$

A similar argument will show

$$S_{j_{ej+1}}(d^0, D) - \sum_{j: j \in j, j^{(d;D)} > 1=(1+w)} S_j(d; D) : \quad (56)$$

It then follows from (55) and (56) that the monotonicity penalty of  $d^0$  is at most that of  $d$ :

$$\begin{aligned} S(d^0, D) &= \sum_{j=0}^{\infty} S_j(d^0, D) A^j + S_{j_{ej}}(d^0, D) + S_{j_{ej+1}}(d^0, D) \\ &= \sum_{j=0}^{\infty} S_j(d; D) A^j + \sum_{j: j \in j, j^{(d;D)} > 1=(1+w)} S_j(d; D) \end{aligned} \quad (57)$$

$$\begin{aligned} &+ \sum_{j: j \in j, j^{(d;D)} > 1=(1+w)} S_j(d; D) \\ &= S(d; D) : \end{aligned} \quad (58)$$

Step 4. Put everything together.

Using (51) and (58), together with the observation  $jd^0_j = je_j + 1 - j d_j$ , we must also have

$$\begin{aligned} \Gamma(d^0; D; 1=(1+w); w; C; C_1) &= R(d^0; D; 1=(1+w); w) + Cjd^0_j + C_1S(d^0; D) \\ &= R(d; D; 1=(1+w); w) + Cjd_j + C_1S(d; D) \\ &= \Gamma(d; D; 1=(1+w); w; C; C_1); \end{aligned}$$

Case 2. Either  $\binom{(d;D)}{j} > 1=(1+w)$  holds for all  $j \geq j_e; \dots; j_djg$ , or  $\binom{(d;D)}{j} = 1=(1+w)$  holds for all  $j \geq j_e; \dots; j_djg$ . The construction of  $d^0 = e$  and the proof for  $R(d^0; D; 1=(1+w); w) = R(d; D; 1=(1+w); w)$  is similar to those given in the proof of Lemma 4.4. The proof for  $S(d^0; D) = S(d; D)$  is similar to that in Case 1. The desired inequality then follows from  $jd^0_j = je_j - j d_j$ .  $\square$

Before we proceed with proving Theorem 5.2, we make the following four observations. Observations 11.1, 11.2, and 11.3 are the same as Observations 10.1, 10.2 and 10.3. They are repeated here for convenience.

Observation 11.1 For any rule list

$$d^0 = f e; (a_{je_j}^{(d^0)}; \wedge_{je_j}^{(d^0)}); \dots; (a_{jd^0_j-1}^{(d^0)}; \wedge_{jd^0_j-1}^{(d^0)}); \wedge_{jd^0_j}^{(d^0)} g$$

that begins with a given pre x e, we have

$$\mathfrak{R}_{e;D}^+ = n_{je_j;d^0;D}^+ + \dots; n_{jd^0_j;d^0;D}^+; \quad (59)$$

$$\mathfrak{R}_{e;D} = n_{je_j;d^0;D} + \dots; n_{jd^0_j;d^0;D}; \quad (60)$$

and

$$\mathfrak{R}_{e;D} = n_{je_j;d^0;D} + \dots; n_{jd^0_j;d^0;D}; \quad (61)$$

Proof. Same as Observation 10.1.  $\square$

Observation 11.2. For any rule list  $d^0$ , we have

$$n_{je_j;d^0;D} = \binom{(d^0;D)}{je_j} n_{je_j;d^0;D}^+; \quad (62)$$

Proof. Same as Observation 10.2.  $\square$

Observation 11.3. For any rule list

$$d^0 = f e; (a_{je_j}^{(d^0)}; \wedge_{je_j}^{(d^0)}); \wedge_{je_{j+1}}^{(d^0)} g$$

that has exactly one rule (excluding the final else clause) following a given pre x e, we have

$$n_{je_{j+1};d^0;D}^+ = \mathfrak{R}_{e;D}^+ - n_{je_j;d^0;D}^+; \quad (63)$$

$$n_{je_{j+1};d^0;D} = \mathfrak{R}_{e;D} - n_{je_j;d^0;D}; \quad (64)$$

and

$$n_{je_{j+1};d^0;D} = \mathfrak{R}_{e;D} - n_{je_j;d^0;D}; \quad (65)$$

Note that since  $n_{je_{j+1};d^0;D}^+$ ,  $n_{je_{j+1};d^0;D}$ , and  $n_{je_{j+1};d^0;D}$  are non-negative, Equations (63), (64), and (65) imply  $n_{je_j;d^0;D}^+ \leq \mathfrak{R}_{e;D}^+$ ,  $n_{je_j;d^0;D} \leq \mathfrak{R}_{e;D}$ , and  $n_{je_j;d^0;D} \leq \mathfrak{R}_{e;D}$ .

Proof. Same as Observation 10.3. □

Observation 11.4. For any rule list

$$d^0 = f e; (a_{j_{ej}}^{(d^0)}; \wedge_{j_{ej}}^{(d^0)}); \wedge_{j_{ej+1}}^{(d^0)} g$$

that has exactly one rule (excluding the final else clause) following a given prefix e, we have

$$n_{j_{ej+1}}^{(d^0; D)} = \frac{r_{e; D}^+ n_{j_{ej}; d^0; D}^+}{r_{e; D}^+ + r_{e; D} \frac{1}{n_{j_{ej}}^{(d^0; D)}} n_{j_{ej}; d^0; D}^+} : \quad (66)$$

Proof. By Definition 2.5, we have

$$n_{j_{ej+1}}^{(d^0; D)} = \frac{n_{j_{ej+1}; d^0; D}^+}{n_{j_{ej+1}; d^0; D}^+ + n_{j_{ej+1}; d^0; D}^+} = \frac{n_{j_{ej+1}; d^0; D}^+}{n_{j_{ej+1}; d^0; D}^+ + n_{j_{ej+1}; d^0; D}^+} :$$

Applying Equations (63) and (64) in Observation 11.3, we have

$$\begin{aligned} n_{j_{ej+1}}^{(d^0; D)} &= \frac{r_{e; D}^+ n_{j_{ej}; d^0; D}^+}{(r_{e; D}^+ n_{j_{ej}; d^0; D}^+) + (r_{e; D} n_{j_{ej}; d^0; D}^+)} \\ &= \frac{r_{e; D}^+ n_{j_{ej}; d^0; D}^+}{r_{e; D}^+ + r_{e; D} \frac{1}{n_{j_{ej}}^{(d^0; D)}} n_{j_{ej}; d^0; D}^+} : \end{aligned}$$

Applying Equation (62) in Observation 11.2, we have

$$\begin{aligned} n_{j_{ej+1}}^{(d^0; D)} &= \frac{r_{e; D}^+ n_{j_{ej}; d^0; D}^+}{r_{e; D}^+ + r_{e; D} \frac{1}{n_{j_{ej}}^{(d^0; D)}} n_{j_{ej}; d^0; D}^+} \\ &= \frac{r_{e; D}^+ n_{j_{ej}; d^0; D}^+}{r_{e; D}^+ + r_{e; D} \frac{1}{n_{j_{ej}}^{(d^0; D)}} n_{j_{ej}; d^0; D}^+} : \end{aligned}$$

□

We are now ready to prove Theorem 5.2.

Proof of Theorem 5.2. Let  $D(X; D; e)$  be the set of (hypothetical and non-hypothetical) rule lists that begin with e and are compatible with D, and let  $D(X; D; e; k)$  be the subset of  $D(X; D; e)$ , consisting of those rule lists in  $D(X; D; e)$  that have exactly k rules (excluding the final else clause) following the prefix e. Let  $S(X; D; e; 1)$  be the subset of  $D(X; D; e; 1)$ , consisting of those rule lists

$$d^0 = f e; (a_{j_{ej}}^{(d^0)}; \wedge_{j_{ej}}^{(d^0)}); \wedge_{j_{ej+1}}^{(d^0)} g \in D(X; D; e; 1)$$

with  $n_{j_{ej}}^{(d^0; D)} > 1 = (1 + w)$  and  $n_{j_{ej+1}}^{(d^0; D)} = 1 = (1 + w)$ .

Note that we have  $D(X; D; e; 0) = f e g$ , where  $e = f e; \sim_{e; D} g$  is the rule list in which the final else clause immediately follows the prefix e, and the probability estimate of the final else clause is  $\sim_{e; D}$ , by a similar argument as that given in the proof of Theorem 4.6 for  $F(X; D; e; 0) = f e g$ .

Let  $d \in D(X; D; e)$ .

The lemma that we have proved in this section, along with its proof, implies

$$\mathbb{E}(d; D; 1 = (1 + w); w; C; C_1) \geq \inf_{d \in D(X; D; e; 1)} \mathbb{E}(d^0; D; 1 = (1 + w); w; C; C_1) : \quad (67)$$

This is because if  $d$  obeys Case 1 in the proof of the lemma, then using the same argument as in the proof of the lemma we can construct a rule  $\text{list}_{d_1} = f e; (a_{j_{ej}}^{(d_1)}; \frac{(d_1; D)}{j_{ej}}); \frac{(d_1; D)}{j_{ej}+1} g \geq S(X; D; e; 1)$  that satisfies

$$\Gamma(d; D; 1=(1+w); w; C; C_1) \quad \Gamma(d_1; D; 1=(1+w); w; C; C_1): \quad (68)$$

Since  $d_1$  must also obey

$$\Gamma(d_1; D; 1=(1+w); w; C; C_1) \quad \inf_{d^0 \in S(X; D; e; 1)} \Gamma(d^0; D; 1=(1+w); w; C; C_1) \\ \inf_{d^0 \in S(X; D; e; 1)} \inf_{D(X; D; e; 0)} \Gamma(d^0; D; 1=(1+w); w; C; C_1); \quad (69)$$

combining the inequalities in (68) and (69) gives us (67). On the other hand, if  $d$  obeys Case 2 in the proof of the lemma, then by the lemma itself we know

$$\Gamma(d; D; 1=(1+w); w; C; C_1) \quad \Gamma(e; D; 1=(1+w); w; C; C_1): \quad (70)$$

Since we have  $D(X; D; e; 0) = f e g$ , it is straightforward to see

$$\Gamma(e; D; 1=(1+w); w; C; C_1) = \inf_{d^0 \in D(X; D; e; 0)} \Gamma(d^0; D; 1=(1+w); w; C; C_1) \\ \inf_{d^0 \in S(X; D; e; 1)} \inf_{D(X; D; e; 0)} \Gamma(d^0; D; 1=(1+w); w; C; C_1): \quad (71)$$

Combining the inequalities in (70) and (71) again gives us (67).

Note that if  $S(X; D; e; 1)$  is not empty, then the right-hand side of (67) can be expressed as

$$\inf_{d^0 \in S(X; D; e; 1)} \inf_{D(X; D; e; 0)} \Gamma(d^0; D; 1=(1+w); w; C; C_1) \\ = \inf_{d^0 \in S(X; D; e; 1)} \inf_{S_{f e g}} \Gamma(d^0; D; 1=(1+w); w; C; C_1) \\ = \min_{d^0 \in S(X; D; e; 1)} \inf \Gamma(d^0; D; 1=(1+w); w; C; C_1); \Gamma(e; D; 1=(1+w); w; C; C_1) : \quad (72)$$

The rest of the proof proceeds in six steps.

Step 1. Compute  $\Gamma(e; D; 1=(1+w); w; C; C_1)$ .

Since the contribution by the final else clause to  $\Gamma(e; D; 1=(1+w); w; C; C_1)$  is given by  $R_{j_{ej}}(e; D; 1=(1+w); w) + b_{\sim e; D} \min_{(e; D)} c_+$ , where  $R_{j_{ej}}(e; D; 1=(1+w); w)$  is defined by Equation (13) and is given by

$$R_{j_{ej}}(e; D; 1=(1+w); w) = \begin{cases} \frac{1}{n} n_{j_{ej}; e; D} & \text{if } \sim e; D > 1=(1+w) \\ \frac{w}{n} n_{j_{ej}; e; D}^+ & \text{otherwise,} \end{cases}$$

and since Observation 11.1 implies  $n_{e; D}^+ = n_{j_{ej}; e; D}^+$  and  $n_{e; D} = n_{j_{ej}; e; D}$ , it is not difficult to see

$$\Gamma(e; D; 1=(1+w); w; C; C_1) \\ = \begin{cases} \Gamma(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} n_{e; D} + C_1 b_{\sim e; D} & \min_{(e; D)} c_+ \quad \text{if } \sim e; D > 1=(1+w) \\ \Gamma(e; D; 1=(1+w); w; C; C_1) + \frac{w}{n} n_{e; D}^+ + C_1 b_{\sim e; D} & \min_{(e; D)} c_+ \quad \text{otherwise.} \end{cases}$$

Since  $\sim e; D > 1=(1+w)$  is equivalent to  $n_{e; D}^+ = (n_{e; D}^+ + n_{e; D}) > 1=(1+w)$ , or  $w n_{e; D}^+ > n_{e; D}$ , and similarly  $\sim e; D \leq 1=(1+w)$  is equivalent to  $w n_{e; D}^+ \leq n_{e; D}$ , we can write

$$\Gamma(e; D; 1=(1+w); w; C; C_1) \\ = \min \Gamma(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} n_{e; D} + C_1 b_{\sim e; D} \min_{(e; D)} c_+; \\ \Gamma(e; D; 1=(1+w); w; C; C_1) + \frac{w}{n} n_{e; D}^+ + C_1 b_{\sim e; D} \min_{(e; D)} c_+ \quad (73) \\ = \Gamma(e; D; 1=(1+w); w; C; C_1) + \min \left\{ \frac{w}{n} n_{e; D}^+; \frac{1}{n} n_{e; D} \right\} + C_1 b_{\sim e; D} \min_{(e; D)} c_+ :$$

Step 2. Partition the set  $S(X; D; e; 1)$  into three subsets based on how the softly falling objective is computed.

For any  $d^0 = f e; (a_{j_{ej}}^{(d^0)}; \frac{(d^0; D)}{j_{ej}}); \frac{(d^0; D)}{j_{ej+1}} \geq 2 S(X; D; e; 1)$ , the softly falling objective is given by

$$\begin{aligned} & \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \\ &= \mathcal{L}(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} n_{j_{ej}; d^0; D} + \frac{w}{n} n_{j_{ej+1}; d^0; D}^+ + C \\ &+ C_1 b_{j_{ej}}^{(d^0; D)} \frac{(e; D)}{\min} c_+ + C_1 b_{j_{ej+1}}^{(d^0; D)} \frac{(e; D)}{\min} c_+ : \end{aligned} \quad (74)$$

This is because for any  $d^0 \geq 2 S(X; D; e; 1)$ , the contribution by both the  $j_{ej}$ -th rule and the final else clause to  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  is given by

$$R_{j_{ej}}(d^0; D; 1=(1+w); w) + R_{j_{ej+1}}(d^0; D; 1=(1+w); w) + C + C_1 b_{j_{ej}}^{(d^0; D)} \frac{(e; D)}{\min} c_+ + C_1 b_{j_{ej+1}}^{(d^0; D)} \frac{(e; D)}{\min} c_+ ;$$

where  $R_{j_{ej}}(d^0; D; 1=(1+w); w)$  and  $R_{j_{ej+1}}(d^0; D; 1=(1+w); w)$  are defined by Equation (13) and are given by

$$R_{j_{ej}}(d^0; D; 1=(1+w); w) = \frac{1}{n} n_{j_{ej}; d^0; D} \quad \text{and} \quad R_{j_{ej+1}}(d^0; D; 1=(1+w); w) = \frac{w}{n} n_{j_{ej+1}; d^0; D}^+$$

(because we have  $\frac{(d^0; D)}{j_{ej}} > 1=(1+w)$  and  $\frac{(d^0; D)}{j_{ej+1}} \leq 1=(1+w)$  for  $d^0 \geq 2 S(X; D; e; 1)$ ).

Let

$$S_1(X; D; e; 1) = f d^0 = f e; (a_{j_{ej}}^{(d^0)}; \frac{(d^0; D)}{j_{ej}}); \frac{(d^0; D)}{j_{ej+1}} \geq 2 S(X; D; e; 1) : \frac{(e; D)}{\min} \frac{(d^0; D)}{j_{ej}} > \frac{(d^0; D)}{j_{ej+1}} g;$$

$$S_2(X; D; e; 1) = f d^0 = f e; (a_{j_{ej}}^{(d^0)}; \frac{(d^0; D)}{j_{ej}}); \frac{(d^0; D)}{j_{ej+1}} \geq 2 S(X; D; e; 1) : \frac{(d^0; D)}{j_{ej}} > \frac{(e; D)}{\min} \frac{(d^0; D)}{j_{ej+1}} g;$$

and

$$S_3(X; D; e; 1) = f d^0 = f e; (a_{j_{ej}}^{(d^0)}; \frac{(d^0; D)}{j_{ej}}); \frac{(d^0; D)}{j_{ej+1}} \geq 2 S(X; D; e; 1) : \frac{(d^0; D)}{j_{ej}} > \frac{(d^0; D)}{j_{ej+1}} > \frac{(e; D)}{\min} g;$$

It is easy to see

$$S(X; D; e; 1) = S_3(X; D; e; 1) [ S_1(X; D; e; 1) [ S_2(X; D; e; 1) ]$$

We observe here that given the prefix, we can write  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  as a function of  $n_{j_{ej}; d^0; D}^+$  and  $\frac{(d^0; D)}{j_{ej}}$ , by substituting (62), (63), and (66) in Observations 11.1, 11.2, and 11.3 into (74).

Step 3. Determine a lower bound of  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  for all  $d^0 \geq 2 S_1(X; D; e; 1)$ .

Let  $d^0 = f e; (a_{j_{ej}}^{(d^0)}; \frac{(d^0; D)}{j_{ej}}); \frac{(d^0; D)}{j_{ej+1}} \geq 2 S_1(X; D; e; 1)$ . By the definition of  $S_1(X; D; e; 1)$ , we have

$$\frac{(e; D)}{\min} \frac{(d^0; D)}{j_{ej}} > \frac{1}{1+w} \frac{(d^0; D)}{j_{ej+1}} \quad (75)$$

We first prove the following inequality

$$\frac{(e; D)}{\min} \frac{(d^0; D)}{j_{ej}} > \max(1=(1+w); \tilde{e}; D); \quad (76)$$

which will be useful later.

To prove (76), we use Definition 2.5 as well as (63) and (65) in Observation 11.3 to obtain

$$\tilde{e}; D = \frac{n_{e; D}^+}{n_{e; D}} = \frac{n_{j_{ej}; d^0; D}^+ + n_{j_{ej+1}; d^0; D}^+}{n_{j_{ej}; d^0; D} + n_{j_{ej+1}; d^0; D}} = \frac{\frac{(d^0; D)}{j_{ej}} n_{j_{ej}; d^0; D} + \frac{(d^0; D)}{j_{ej+1}} n_{j_{ej+1}; d^0; D}}{n_{j_{ej}; d^0; D} + n_{j_{ej+1}; d^0; D}} : \quad (77)$$

Substituting  $\frac{(d^0; D)}{j_{ej+1}} < \frac{(d^0; D)}{j_{ej}}$  from (75) into (77), we obtain  $\tilde{e}; D < \frac{(d^0; D)}{j_{ej}}$ . Combining this inequality with  $\frac{(e; D)}{\min} \frac{(d^0; D)}{j_{ej}} > \frac{1}{1+w}$  from (75), we obtain (76), as desired.

Note that since (76) has to hold for any  $d^0 \in S_1(X; D; e; 1)$ , if  $\min_{(e;D)} \max(1=(1+w); \sim_{e;D})$  is true for the given pre  $x \in e$ , then  $S_1(X; D; e; 1)$  is empty.

We now show that given the pre  $x \in e$ , the softly falling objective  $\Gamma(d^0; D; 1=(1+w); w; C; C_1)$  for  $d^0$  is a monotonically decreasing function of both  $n_{je; d^0; D}^+$  and  $\frac{(d^0; D)}{je}$ .

To do so, we substitute (62) and (63) in Observations 11.1 and 11.2 into (74) to obtain

$$\begin{aligned} & \Gamma(d^0; D; 1=(1+w); w; C; C_1) \\ &= \Gamma(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} \frac{1}{\frac{(d^0; D)}{je}} \left( 1 - wA n_{je; d^0; D}^+ + wR_{e; D}^+ A + C \right) \end{aligned} \quad (78)$$

Note that Equation (78) shows that given the pre  $x \in e$ ,  $\Gamma(d^0; D; 1=(1+w); w; C; C_1)$  is a function of  $n_{je; d^0; D}^+$  and  $\frac{(d^0; D)}{je}$ . Since we have

$$\frac{\partial \Gamma(d^0; D; 1=(1+w); w; C; C_1)}{\partial n_{je; d^0; D}^+} = \frac{1}{n} \frac{\partial}{\partial \frac{(d^0; D)}{je}} \left( 1 - wA \right) < 0$$

because  $\frac{(d^0; D)}{je} > 1=(1+w)$  holds for any  $d^0 \in S_1(X; D; e; 1)$ , and

$$\frac{\partial \Gamma(d^0; D; 1=(1+w); w; C; C_1)}{\partial \frac{(d^0; D)}{je}} = \frac{n_{je; d^0; D}^+}{n} \frac{1}{\left( \frac{(d^0; D)}{je} \right)^2} < 0;$$

we see that  $\Gamma(d^0; D; 1=(1+w); w; C; C_1)$  is indeed a monotonically decreasing function of both  $n_{je; d^0; D}^+$  and  $\frac{(d^0; D)}{je}$ . Thus, we can obtain a lower bound of  $\Gamma(d^0; D; 1=(1+w); w; C; C_1)$  by substituting  $n_{je; d^0; D}^+$  and  $\frac{(d^0; D)}{je}$  with their respective upper bound. The inequality  $n_{je; d^0; D}^+ \leq R_{e; D}^+$  in Observation 11.3 gives an upper bound for  $n_{je; d^0; D}^+$ , and the inequality  $\frac{(d^0; D)}{je} \leq \frac{(e; D)}{\min_{je+1}}$  from (75) gives an upper bound for  $\frac{(d^0; D)}{je}$ . Substituting these upper bounds into (78), we obtain the following inequality, which gives a lower bound of  $\Gamma(d^0; D; 1=(1+w); w; C; C_1)$ :

$$\begin{aligned} & \Gamma(d^0; D; 1=(1+w); w; C; C_1) \\ & \geq \Gamma(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} \frac{1}{\frac{(e; D)}{\min_{je+1}}} \left( 1 - w R_{e; D}^+ + w R_{e; D}^+ A + C \right) \\ & = \Gamma(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} \frac{1}{\frac{(e; D)}{\min_{je+1}}} \left( 1 - R_{e; D}^+ \right) + C \end{aligned}$$

Step 4. Determine a lower bound of  $\Gamma(d^0; D; 1=(1+w); w; C; C_1)$  for all  $d^0 \in S_2(X; D; e; 1)$ .

Let  $d^0 = f e; (a_{je}^{(d^0)}, \frac{(d^0; D)}{je}); \frac{(d^0; D)}{je+1} \in S_2(X; D; e; 1)$ . By the definition of  $S_2(X; D; e; 1)$ , we have

$$\frac{(d^0; D)}{je} > \frac{1}{1+w} \quad (79)$$

and

$$\frac{(d^0; D)}{je} > \frac{(e; D)}{\min_{je+1}} \frac{(d^0; D)}{je+1} \quad (80)$$

We first prove the following inequality

$$\frac{(d^0; D)}{je} > \max\left( \frac{(e; D)}{\min_{je+1}}; \sim_{e; D}; 1=(1+w) \right) = ; \quad (81)$$

which will be useful later.

To prove (81), we use Definition 2.5 as well as (63) and (65) in Observation 11.3 to obtain (77). Substituting  $\frac{(d^0;D)}{j_{ej}+1} < \frac{(d^0;D)}{j_{ej}}$  from (75) into (77), we obtain  $\frac{(d^0;D)}{j_{ej}} < \frac{(d^0;D)}{j_{ej}}$ . Combining this inequality with (79) and  $\frac{(d^0;D)}{j_{ej}} > \frac{(e;D)}{\min}$  from (80), we obtain (81), as desired.

We now show that given the pre x e and a particular value of  $\frac{(d^0;D)}{j_{ej}}$  that obeys (81), the softly falling objective  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  for  $d^0$  is a decreasing function of  $n_{j_{ej};d^0;D}^+$ .

To do so, we substitute (62) and (63) in Observations 11.1 and 11.2 into (74) to obtain

$$\begin{aligned} & \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \\ &= \mathcal{L}(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} \frac{1}{\frac{(d^0;D)}{j_{ej}}} \frac{1}{1} wA n_{j_{ej};d^0;D}^+ + wR_{e;D}^+ A + C \\ &+ C_1 \left( \frac{(d^0;D)}{j_{ej}} \frac{(e;D)}{\min} \right); \end{aligned} \tag{82}$$

Note that Equation (82) shows that given the pre x e,  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  is a function of  $n_{j_{ej};d^0;D}^+$  and  $\frac{(d^0;D)}{j_{ej}}$ . Differentiating  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  given in (82) with respect to  $n_{j_{ej};d^0;D}^+$ , we obtain

$$\frac{\partial \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)}{\partial n_{j_{ej};d^0;D}^+} = \frac{1}{n} \frac{1}{\frac{(d^0;D)}{j_{ej}}} \frac{1}{1} wA; \tag{83}$$

Since  $\frac{(d^0;D)}{j_{ej}}$  obeys (81), in particular, it obeys  $\frac{(d^0;D)}{j_{ej}} > 1=(1+w)$ , we have

$$\frac{1}{\frac{(d^0;D)}{j_{ej}}} \frac{1}{1} w < 0;$$

which then gives  $\frac{\partial \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)}{\partial n_{j_{ej};d^0;D}^+} < 0$ . This means that given the pre x e and a particular value of  $\frac{(d^0;D)}{j_{ej}}$  that obeys (81),  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  is a decreasing function of  $n_{j_{ej};d^0;D}^+$ .

Thus, given the pre x e and a particular value of  $\frac{(d^0;D)}{j_{ej}}$  that obeys (81), we can obtain a lower bound of  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  by substituting  $n_{j_{ej};d^0;D}^+$  with its upper bound. The inequality  $n_{j_{ej};d^0;D}^+ \leq R_{e;D}^+$  in Observation 11.3 gives an upper bound for  $n_{j_{ej};d^0;D}^+$ . Substituting  $n_{j_{ej};d^0;D}^+$  with its upper bound  $R_{e;D}^+$  into (82), we obtain a lower bound of  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$ , denoted by  $g\left(\frac{(d^0;D)}{j_{ej}}\right)$ , when  $\frac{(d^0;D)}{j_{ej}}$  is held constant:

$$\begin{aligned} g\left(\frac{(d^0;D)}{j_{ej}}\right) &= \mathcal{L}(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} \frac{1}{\frac{(d^0;D)}{j_{ej}}} \frac{1}{1} wA R_{e;D}^+ + C + C_1 \left( \frac{(d^0;D)}{j_{ej}} \frac{(e;D)}{\min} \right) \\ &= \mathcal{L}(e; D; 1=(1+w); w; C; C_1) + g\left(\frac{(d^0;D)}{j_{ej}}\right) \end{aligned}$$

where  $g$  is defined in the statement of the theorem. In other words, given the pre x e and a particular value of  $\frac{(d^0;D)}{j_{ej}}$  that obeys (81), we have  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \geq g\left(\frac{(d^0;D)}{j_{ej}}\right)$ . Since (81) is true for any  $d^0 \in S_2(X; D; e; 1)$ , we always have  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \geq g\left(\frac{(d^0;D)}{j_{ej}}\right)$  for any  $d^0 \in S_2(X; D; e; 1)$ . This implies

$$\begin{aligned} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) & \geq \inf_{\frac{(d^0;D)}{j_{ej}} : < \frac{(d^0;D)}{j_{ej}} > 1} g\left(\frac{(d^0;D)}{j_{ej}}\right) \\ &= \mathcal{L}(e; D; 1=(1+w); w; C; C_1) + \inf_{\frac{(d^0;D)}{j_{ej}} : < \frac{(d^0;D)}{j_{ej}} > 1} g\left(\frac{(d^0;D)}{j_{ej}}\right); \end{aligned}$$



Step 5. Determine a lower bound of  $\underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1)$  for all  $d^0 \in S_3(X; D; e; 1)$ .

Let  $d^0 = f_{ej}; (a_{je}; (d^0; D)); (d^0; D)$   $g \in S_3(X; D; e; 1)$ . By the definition of  $S_3(X; D; e; 1)$ , we have

$$\frac{(d^0; D)}{j_{ej}} > \frac{1}{1+w} \quad \frac{(d^0; D)}{j_{ej+1}} > \frac{(e; D)}{\min} \quad (84)$$

We first prove the following inequality

$$1 - \frac{(d^0; D)}{j_{ej}} > \max\left(\frac{(e; D)}{\min}; \sim_{e; D}; 1=(1+w)\right) = ; \quad (85)$$

which will be useful later.

To prove (85), we use Definition 2.5 as well as (63) and (65) in Observation 11.3 to obtain (77). Substituting  $\frac{(d^0; D)}{j_{ej+1}} < \frac{(d^0; D)}{j_{ej}}$  from (84) into (77), we obtain  $\sim_{e; D} < \frac{(d^0; D)}{j_{ej}}$ . Combining this inequality with  $\frac{(d^0; D)}{j_{ej}} > \frac{1}{1+w} > \frac{(e; D)}{\min}$  from (84), we obtain (85), as desired.

To determine a lower bound of  $\underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1)$ , we observe

$$\begin{aligned} & \underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1) \\ & \underline{\Gamma}(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} n_{j_{ej}; d^0; D} + \frac{w}{n} n_{j_{ej+1}; d^0; D} + C + C_1 b_{j_{ej}} \frac{(d^0; D)}{1} \frac{(e; D)}{\min} C_+ \quad (86) \\ & = \underline{\Gamma}(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} \frac{1}{(d^0; D)} \frac{1}{j_{ej}} \frac{wA}{n_{j_{ej}; d^0; D}} + wA_{e; D}^+ A + C \\ & + C_1 \left( \frac{(d^0; D)}{j_{ej}} \frac{(e; D)}{\min} \right) \quad (87) \end{aligned}$$

where the last equality follows by substituting (62) and (63) in Observations 11.1 and 11.2 into (86). Using (85) and applying the same argument as in Step 4, the quantity labeled (87) is also lower-bounded by

$$\underline{\Gamma}(e; D; 1=(1+w); w; C; C_1) + \inf_{\frac{(d^0; D)}{j_{ej}} < \frac{(d^0; D)}{j_{ej}}} g\left(\frac{(d^0; D)}{j_{ej}}\right);$$

so that we again have

$$\underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1) \geq \underline{\Gamma}(e; D; 1=(1+w); w; C; C_1) + \inf_{\frac{(d^0; D)}{j_{ej}} < \frac{(d^0; D)}{j_{ej}}} g\left(\frac{(d^0; D)}{j_{ej}}\right);$$

Step 6. Put everything together.

Suppose, first, that  $S(X; D; e; 1)$  is not empty.

In the case where  $S_1(X; D; e; 1)$  is not empty, we observe the following inequality

$$\inf_{d^0 \in S_1(X; D; e; 1)} \underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1) \geq \underline{\Gamma}(e; D; 1=(1+w); w; C; C_1) + \frac{1}{n} \frac{1}{\frac{(e; D)}{\min}} \frac{1}{1} n_{e; D}^+ + C; \quad (88)$$

which follows from the definition of  $\inf$  being the greatest lower bound, as well as the lower bound of  $\underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1)$  for  $d^0 \in S_1(X; D; e; 1)$ , which we have derived in Step 3.

In the case where  $S_2(X; D; e; 1) \cap S_3(X; D; e; 1)$  is not empty, we observe the following inequality

$$\inf_{d^0 \in S_2(X; D; e; 1) \cap S_3(X; D; e; 1)} \underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1) \geq \underline{\Gamma}(e; D; 1=(1+w); w; C; C_1) + \inf_{<} g\left(\frac{(d^0; D)}{j_{ej}}\right); \quad (89)$$

which follows from the definition of  $\inf$  being the greatest lower bound, as well as the lower bound of  $\underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1)$  for  $d^0 \in S_2(X; D; e; 1)$ , which we have derived in Step 4, and the lower bound of  $\underline{\Gamma}(d^0; D; 1=(1+w); w; C; C_1)$  for  $d^0 \in S_3(X; D; e; 1)$ , which we have derived in Step 5.

To derive a lower bound of  $\mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  for  $d^0 \geq S(X; D; e; 1)$ , we further observe that if  $\min_{(e;D)} \max(1=(1+w); \sim_{e;D})$  holds, then by our remark in Step 3,  $S_1(X; D; e; 1)$  is empty, and consequently, using (89), we have

$$\inf_{d^0 S_2(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) = \inf_{d^0 S_2(X; D; e; 1) \cap S_3(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) + \inf_{: < } g( ) \quad (90)$$

On the other hand, if  $\min_{(e;D)} > \max(1=(1+w); \sim_{e;D})$  holds, then  $S_1(X; D; e; 1)$  may or may not be empty. If, in addition, both  $S_1(X; D; e; 1)$  and  $S_2(X; D; e; 1) \cap S_3(X; D; e; 1)$  are not empty, then using (88) and (89), we have

$$\begin{aligned} & \inf_{d^0 S_2(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \\ = & \min_{d^0 S_1(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1); \inf_{d^0 S_2(X; D; e; 1) \cap S_3(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \\ & \mathcal{L}(e; D; 1=(1+w); w; C; C_1) + \min \frac{1}{n} \frac{1}{\min_{(e;D)}} 1 \mathcal{R}_{e;D}^+ + C; \inf_{: < } g( ) \quad (91) \end{aligned}$$

If either  $S_1(X; D; e; 1)$  or  $S_2(X; D; e; 1) \cap S_3(X; D; e; 1)$  is empty, then  $\inf_{d^0 S_2(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1)$  is given by either

$$\inf_{d^0 S_2(X; D; e; 1) \cap S_3(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \quad \text{or} \quad \inf_{d^0 S_1(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1);$$

both of which are lower-bounded by the quantity labeled (91) because of (89) and (88).

Putting these cases together, we have

$$\begin{aligned} & \inf_{d^0 S_2(X; D; e; 1)} \mathcal{L}(d^0; D; 1=(1+w); w; C; C_1) \\ & \mathcal{L}(e; D; 1=(1+w); w; C; C_1) \\ + & \begin{cases} \min \frac{1}{n} \frac{1}{\min_{(e;D)}} 1 \mathcal{R}_{e;D}^+ + C; \inf_{: < } g( ) & \text{if } \min_{(e;D)} > \max(1=(1+w); \sim_{e;D}); \\ \inf_{: < } g( ) & \text{otherwise.} \end{cases} \quad (92) \end{aligned}$$

Combining (67), (72), (73), and (92), we have

$$\begin{aligned} & \mathcal{L}(d; D; 1=(1+w); w; C; C_1) \\ & \mathcal{L}(e; D; 1=(1+w); w; C; C_1) \\ + & \begin{cases} \min \frac{1}{n} \frac{1}{\min_{(e;D)}} 1 \mathcal{R}_{e;D}^+ + C; \inf_{: < } g( ); \frac{w}{n} \mathcal{R}_{e;D}^+ + C_1 b_{\sim_{e;D}} \min_{(e;D)} C_+; \\ \frac{1}{n} \mathcal{R}_{e;D} + C_1 b_{\sim_{e;D}} \min_{(e;D)} C_+ & \text{if } \min_{(e;D)} > \max(1=(1+w); \sim_{e;D}); \\ \min \inf_{: < } g( ); \frac{w}{n} \mathcal{R}_{e;D}^+ + C_1 b_{\sim_{e;D}} \min_{(e;D)} C_+; \\ \frac{1}{n} \mathcal{R}_{e;D} + C_1 b_{\sim_{e;D}} \min_{(e;D)} C_+ & \text{otherwise.} \end{cases} \quad (93) \end{aligned}$$

Note that the quantity labeled (93) is precisely equal to  $\mathcal{L}(e; D; w; C; C_1)$  given by Equation (44) in the statement of the theorem, because:

(i) if  $\min_{(e;D)} > \max(1=(1+w); \sim_{e;D})$  holds, then the first term in the minimum on the right-hand side of Equation (44) is precisely  $\frac{1}{n} \frac{1}{\min_{(e;D)}} 1 \mathcal{R}_{e;D}^+ + C;$

(ii) if  $\frac{R_{e,D}^{(e;D)}}{\min} > \max(1=(1+w); \sim_{e,D})$  does not hold, then we have  $\frac{R_{e,D}^{(e;D)}}{\min} = 1=(1+w)$  or  $\frac{R_{e,D}^{(e;D)}}{\min} < \sim_{e,D}$  in the former case where  $\frac{R_{e,D}^{(e;D)}}{\min} = 1=(1+w)$  holds, we have

$$\frac{1}{n} - \frac{1}{\frac{R_{e,D}^{(e;D)}}{\min}} + \frac{1}{n} R_{e,D}^+ + \frac{w}{n} R_{e,D}^+;$$

which implies that the first term in the minimum on the right-hand side of Equation (44) is bounded below by  $\frac{w}{n} R_{e,D}^+ + C_1 b_{\sim_{e,D}} \frac{R_{e,D}^{(e;D)}}{\min} C_+$ , and thus has no influence over the computation of the minimum; in the latter case where  $\frac{R_{e,D}^{(e;D)}}{\min} < \sim_{e,D}$  holds, the first term in the minimum on the right-hand side of Equation (44) is clearly bounded below by  $\frac{w}{n} R_{e,D}^+ + C_1 b_{\sim_{e,D}} \frac{R_{e,D}^{(e;D)}}{\min} C_+$ , and again has no influence over the computation of the minimum.

This proves that  $\underline{L}(e; D; w; C; C_1)$  given by Equation (44) is indeed a lower bound of  $\underline{L}(d; D; 1=(1+w); w; C; C_1)$  for  $d \in D(X; D; e)$ , in the case where  $S(X; D; e; 1)$  is not empty. In the case where  $S(X; D; e; 1)$  is empty, using (67) and (73), along with the fact  $D(X; D; e; 0) = \emptyset$ , we have

$$\begin{aligned} \underline{L}(d; D; 1=(1+w); w; C; C_1) &= \inf_{d \in D(X; D; e; 0)} \underline{L}(d; D; 1=(1+w); w; C; C_1) \\ &= \underline{L}(e; D; 1=(1+w); w; C; C_1) \\ &= \underline{L}(e; D; 1=(1+w); w; C; C_1) + \min \left( \frac{w}{n} R_{e,D}^+; \frac{1}{n} R_{e,D}^+ + C_1 b_{\sim_{e,D}} \frac{R_{e,D}^{(e;D)}}{\min} C_+ \right); \end{aligned}$$

where the last quantity is clearly lower-bounded by  $\underline{L}(e; D; w; C; C_1)$  defined in Equation (44). We have now proven that  $\underline{L}(e; D; w; C; C_1)$  given by Equation (44) is a lower bound of  $\underline{L}(d; D; 1=(1+w); w; C; C_1)$  for  $d \in D(X; D; e)$ .

Finally, we compute  $\inf_{\alpha < 1} g(\alpha)$  analytically. Since the derivative of  $g$  is given by

$$g'(\alpha) = \frac{R_{e,D}^+}{n^2} + C_1;$$

and must be positive, the only stationary point of  $g$  that could satisfy the constraint  $\alpha < 1$  is given by

$$\alpha = \frac{R_{e,D}^+}{C_1 n};$$

and the second derivative test confirms that  $\alpha$  is a local minimum of  $g$ . It then follows that  $\inf_{\alpha < 1} g(\alpha)$  is given by

$$\inf_{\alpha < 1} g(\alpha) = \begin{cases} g(\alpha) & \text{if } \alpha < 1 \\ \min(g(\alpha); g(1)) & \text{otherwise} \end{cases}$$

□

## 12 Additional Rule Lists Demonstrating the Effect of Varying Parameter Values

In this section, we include some additional rule lists created using Algorithm FRL and Algorithm softFRL with varying parameter values. The default parameter values we used in creating these rule lists are  $w = 7$ ,  $C = 0.000001$ , and  $C_1 = 0.5$ . In each of the following subsections, the rule lists were created with default parameter values, other than the parameter that was being varied.

### 12.1 Effect of Varying $w$ on Algorithm FRL

Running Algorithm FRL with  $w = 1$  on the bank-full dataset produces the following falling rule list:

	antecedent		probability	positive support	negative support
IF	poutcome=success AND loan=no	THEN success prob. is	0.65	934	495
ELSE IF	poutcome=success AND marital=married	THEN success prob. is	0.62	31	19
ELSE IF	poutcome=success AND campaign=1	THEN success prob. is	0.56	9	7
ELSE		success prob. is	0.10	4315	39401

Table 2: Falling rule list for bank-full dataset, created using Algorithm FRL with  $w = 1$

Running Algorithm FRL with  $w = 3$  on the bank-full dataset produces the following falling rule list:

	antecedent		probability	positive support	negative support
IF	poutcome=success AND previous 2	THEN success prob. is	0.65	677	361
ELSE IF	poutcome=success AND campaign=1	THEN success prob. is	0.65	185	99
ELSE IF	poutcome=success AND loan=no	THEN success prob. is	0.63	111	65
ELSE IF	poutcome=success AND marital=married	THEN success prob. is	0.56	5	4
ELSE IF	60 age < 100 AND housing=no	THEN success prob. is	0.30	390	919
ELSE		success prob. is	0.09	3921	38474

Table 3: Falling rule list for bank-full dataset, created using Algorithm FRL with  $w = 3$

Running Algorithm FRL with  $w = 5$  on the bank-full dataset produces the following falling rule list:

	antecedent		probability	positive support	negative support
IF	poutcome=success AND default=no	THEN success prob. is	0.65	978	531
ELSE IF	60 age < 100 AND loan=no	THEN success prob. is	0.29	426	1030
ELSE IF	17 age < 30 AND housing=no	THEN success prob. is	0.25	504	1539
ELSE IF	previous 2 AND housing=no	THEN success prob. is	0.23	242	796
ELSE		success prob. is	0.08	3139	36026

Table 4: Falling rule list for bank-full dataset, created using Algorithm FRL with  $w = 5$

Running Algorithm FRL with  $w = 7$  on the bank-full dataset produces the following falling rule list:

	antecedent		probability	positive support	negative support
IF	poutcome=success AND default=no	THEN success prob. is	0.65	978	531
ELSE IF	60 age < 100 AND default=no	THEN success prob. is	0.28	434	1113
ELSE IF	17 age < 30 AND housing=no	THEN success prob. is	0.25	504	1539
ELSE IF	previous 2 AND housing=no	THEN success prob. is	0.23	242	794
ELSE IF	campaign=1 AND housing=no	THEN success prob. is	0.14	658	4092
ELSE IF	previous 2 AND education=tertiary	THEN success prob. is	0.13	108	707
ELSE		success prob. is	0.07	2365	31146

Table 5: Falling rule list for bank-full dataset, created using Algorithm FRL with  $w = 7$

As the positive class weight  $w$  increases, the falling rule list created using Algorithm FRL tends to have rules whose probability estimates are smaller. This is not surprising – a larger value of  $w$  means a smaller threshold  $\tau = 1/(1+w)$ , and by including rules whose probability estimates are not much larger than the threshold, the falling rule list produced by the algorithm will more likely predict positive, thereby reducing the (weighted) empirical risk of misclassification. Note that Algorithm FRL will never include rules whose probability estimates are less than the threshold (see Corollary 4.5).

## 12.2 Effect of Varying $w$ on Algorithm softFRL

Running Algorithm softFRL with  $w = 1$  on the bank-full dataset produces the following softly falling rule list:

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success AND campaign=1	THEN prob. is	0.67	0.67	557	280
ELSE IF	poutcome=success AND marital=married	THEN prob. is	0.65	0.65	263	143
ELSE IF	poutcome=success AND loan=no	THEN prob. is	0.61	0.61	154	98
ELSE		prob. is	0.10	0.10	4315	39401

Table 6: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $w = 1$

Note that there is an extra column `positive proportion` in a table showing a softly falling rule list. This column gives the empirical positive proportion of each antecedent in the softly falling rule list. When the probability estimate of a rule is less than the positive proportion of the antecedent in the same rule, we know that the softly falling rule list has been transformed from a non-falling compatible rule list, and that the monotonicity penalty has been incurred in the process of running Algorithm softFRL.

Running Algorithm softFRL with  $w = 3$  on the bank-full dataset produces the following softly falling rule list:

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success AND marital=married	THEN prob. is	0.65	0.65	547	289
ELSE IF	poutcome=success AND loan=no	THEN prob. is	0.65	0.65	418	225
ELSE IF	poutcome=success AND campaign=1	THEN prob. is	0.56	0.56	9	7
ELSE IF	poutcome=success AND previous 2	THEN prob. is	0.33	0.33	4	8
ELSE IF	60 age < 100 AND housing=no	THEN prob. is	0.30	0.30	390	919
ELSE IF	previous 2 AND campaign=1	THEN prob. is	0.15	0.15	281	1559
ELSE		prob. is	0.09	0.09	3640	36915

Table 7: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $w = 3$

Running Algorithm softFRL with  $w = 5$  on the bank-full dataset produces the following softly falling rule list:

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success	THEN prob. is	0.65	0.65	978	533
ELSE IF	60 age < 100 AND loan=no	THEN prob. is	0.29	0.29	426	1030
ELSE IF	poutcome=unknown AND contact=cellular	THEN prob. is	0.11	0.11	2380	18659
ELSE		prob. is	0.07	0.07	1505	19700

Table 8: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $w = 5$

Running Algorithm softFRL with  $w = 7$  on the bank-full dataset produces the following softly falling rule list:

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success	THEN prob. is	0.65	0.65	978	533
ELSE IF	60 age < 100	THEN prob. is	0.28	0.28	435	1120
ELSE IF	marital=single AND housing=no	THEN prob. is	0.18	0.18	970	4504
ELSE IF	contact=cellular AND default=no	THEN prob. is	0.10	0.10	2255	19970
ELSE		prob. is	0.05	0.05	651	13795

Table 9: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $w = 7$

As the positive class weight  $w$  increases, the softly falling rule list created using Algorithm softFRL also tends to have rules whose probability estimates are smaller. This is again not surprising a larger value of  $w$  means a smaller threshold  $\tau = 1/(1+w)$ , and by including rules whose probability estimates are not much larger than the threshold, the softly falling rule list produced by the algorithm will more likely predict positive, thereby reducing the (weighted) empirical risk of misclassification.

### 12.3 Effect of Varying C on Algorithm FRL

Running Algorithm FRL with  $C = 0:000001$  on the bank-full dataset produces the following falling rule list:

	antecedent		probability	positive support	negative support
IF	poutcome=success AND default=no	THEN success prob. is	0.65	978	531
ELSE IF	60 age< 100 AND default=no	THEN success prob. is	0.28	434	1113
ELSE IF	17 age< 30 AND housing=no	THEN success prob. is	0.25	504	1539
ELSE IF	previous 2 AND housing=no	THEN success prob. is	0.23	242	794
ELSE IF	campaign=1 AND housing=no	THEN success prob. is	0.14	658	4092
ELSE IF	previous 2 AND education=tertiary	THEN success prob. is	0.13	108	707
ELSE		success prob. is	0.07	2365	31146

Table 10: Falling rule list for bank-full dataset, created using Algorithm FRL with  $C = 0:000001$

Running Algorithm FRL with  $C = 0:01$  on the bank-full dataset produces the following falling rule list:

	antecedent		probability	positive support	negative support
IF	poutcome=success AND default=no	THEN success prob. is	0.65	978	531
ELSE IF	60 age< 100 AND loan=no	THEN success prob. is	0.29	426	1030
ELSE IF	17 age< 30 AND contact=cellular	THEN success prob. is	0.20	653	2621
ELSE IF	campaign=1 AND housing=no	THEN success prob. is	0.15	803	4634
ELSE		success prob. is	0.07	2429	31106

Table 11: Falling rule list for bank-full dataset, created using Algorithm FRL with  $C = 0:01$

Running Algorithm FRL with  $C = 0:1$  on the bank-full dataset produces the following falling rule list:

	antecedent		probability	positive support	negative support
IF	housing=no AND contact=cellular	THEN success prob. is	0.20	2883	11799
ELSE		success prob. is	0.08	2406	28123

Table 12: Falling rule list for bank-full dataset, created using Algorithm FRL with  $C = 0:1$

As the cost  $C$  of adding a rule increases, the size of the falling rule list created by Algorithm FRL decreases, as expected.

## 12.4 Effect of Varying $C$ on Algorithm softFRL

Running Algorithm softFRL with  $C = 0:000001$  on the bank-full dataset produces the following softly falling rule list:

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success	THEN prob. is	0.65	0.65	978	533
ELSE IF	60 age < 100	THEN prob. is	0.28	0.28	435	1120
ELSE IF	marital=single AND housing=no	THEN prob. is	0.18	0.18	970	4504
ELSE IF	contact=cellular AND default=no	THEN prob. is	0.10	0.10	2255	19970
ELSE		prob. is	0.05	0.05	651	13795

Table 13: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $C = 0:000001$

Running Algorithm softFRL with  $C = 0:01$  on the bank-full dataset produces the following softly falling rule list:

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success AND loan=no	THEN prob. is	0.65	0.65	934	495
ELSE IF	housing=no AND contact=cellular	THEN prob. is	0.16	0.16	2245	11535
ELSE IF	housing=yes AND default=no	THEN prob. is	0.07	0.07	1677	22591
ELSE		prob. is	0.07	0.08	433	5301

Table 14: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $C = 0:01$

Running Algorithm softFRL with  $C = 0:1$  on the bank-full dataset produces the following softly falling rule list:

	antecedent		probability	positive proportion	positive support	negative support
IF	housing=no AND contact=cellular	THEN prob. is	0.20	0.20	2883	11799
ELSE		prob. is	0.08	0.08	2406	28123

Table 15: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $C = 0:1$

As the cost  $C$  of adding a rule increases, the size of the softly falling rule list created by Algorithm softFRL decreases, as expected.

## 12.5 Effect of Varying $C_1$ on Algorithm softFRL

Running Algorithm softFRL with  $C_1$  of 0:005; 0:05; 0:5g on the bank-full dataset produces the softly falling rule lists shown in Tables 16, 17, and 18.

When the monotonicity penalty  $C_1$  is small, the softly falling rule list created by Algorithm softFRL exhibits the "pulling down" of the empirical positive proportion for a substantial number of rules, because with little monotonicity penalty the algorithm will more likely choose a rule list that frequently violates monotonicity



	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success	THEN prob. is	0.65	0.65	978	533
ELSE IF	60 ≤ age < 100	THEN prob. is	0.30	0.30	599	1177
	AND housing=no					
ELSE IF	marital=single	THEN prob. is	0.18	0.18	970	4504
	AND housing=no					
ELSE IF	marital=single	THEN prob. is	0.08	0.08	456	4936
	AND previous=0					
ELSE IF	campaign ≥ 3	THEN prob. is	0.06	0.06	323	5294
	AND education=secondary					
ELSE IF	30 ≤ age < 40	THEN prob. is	0.06	0.08	568	6849
	AND previous=0					
ELSE IF	education=tertiary	THEN prob. is	0.06	0.14	361	2237
	AND housing=no					
ELSE IF	loan=yes	THEN prob. is	0.05	0.05	106	1972
	AND previous=0					
ELSE IF	education=secondary	THEN prob. is	0.05	0.09	595	5779
	AND default=no					
ELSE IF	campaign=1	THEN prob. is	0.05	0.08	233	2564
ELSE IF	housing=no	THEN prob. is	0.05	0.05	68	1176
	AND previous=0					
ELSE IF	job=management	THEN prob. is	0.05	0.10	75	693
	AND contact=cellular					
ELSE IF	job=technician	THEN prob. is	0.05	0.07	10	143
	AND poutcome=unknown					
ELSE IF	marital=married	THEN prob. is	0.05	0.06	110	1841
ELSE IF	campaign ≥ 3	THEN prob. is	0.05	0.06	16	238
	AND housing=yes					
ELSE IF	marital=single	THEN prob. is	0.05	0.13	13	91
	AND housing=yes					
ELSE IF	housing=yes	THEN prob. is	0.05	0.10	8	69
	AND contact=cellular					
ELSE IF	job=blue-collar	THEN prob. is	0.05	0.16	4	21
	AND loan=no					
ELSE		prob. is	0.05	0.07	5	63

Table 16: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $C_1 = 0.005$

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success AND default=no	THEN prob. is	0.65	0.65	978	531
ELSE IF	housing=yes	THEN prob. is	0.07	0.07	1686	22974
ELSE IF	$50 \leq \text{age} < 60$ AND poutcome=unknown	THEN prob. is	0.07	0.09	367	3806
ELSE IF	contact=cellular AND default=no	THEN prob. is	0.07	0.18	1927	8961
ELSE IF	campaign=1 AND poutcome=unknown	THEN prob. is	0.07	0.08	126	1374
ELSE IF	campaign $\geq 3$ AND loan=no	THEN prob. is	0.07	0.08	93	1110
ELSE IF	campaign=2 AND education=tertiary	THEN prob. is	0.07	0.09	18	192
ELSE IF	loan=no	THEN prob. is	0.07	0.10	72	648
ELSE		prob. is	0.06	0.06	22	326

Table 17: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $C_1 = 0.05$

	antecedent		probability	positive proportion	positive support	negative support
IF	poutcome=success	THEN prob. is	0.65	0.65	978	533
ELSE IF	$60 \leq \text{age} < 100$	THEN prob. is	0.28	0.28	435	1120
ELSE IF	marital=single AND housing=no	THEN prob. is	0.18	0.18	970	4504
ELSE IF	contact=cellular AND default=no	THEN prob. is	0.10	0.10	2255	19970
ELSE		prob. is	0.05	0.05	651	13795

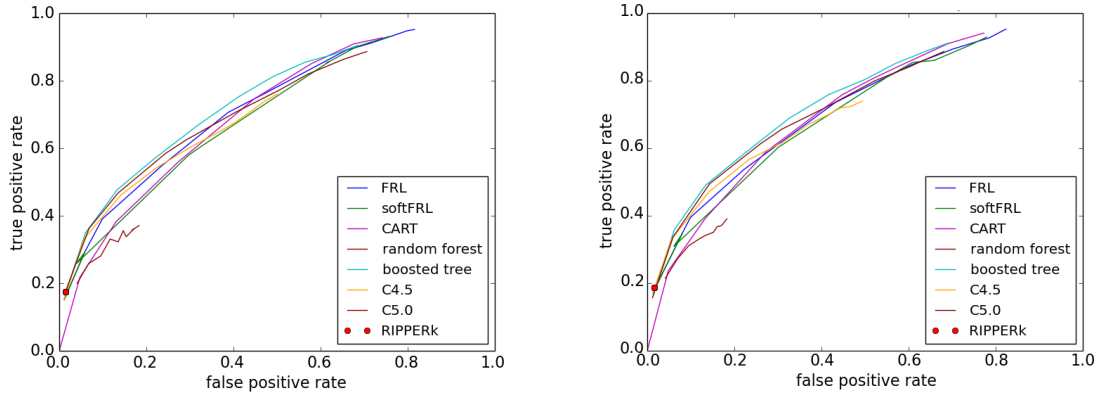
Table 18: Softly falling rule list for bank-full dataset, created using Algorithm softFRL with  $C_1 = 0.5$

but that has a small empirical risk on the training set, in the hope of getting more of the training instances “right”. This is also why the softly falling rule list tends to be longer when  $C_1$  is small: in minimizing the empirical risk on the training set with little regularization (the default  $C = 0.000001$  is very small), the algorithm tends to overfit the training data.

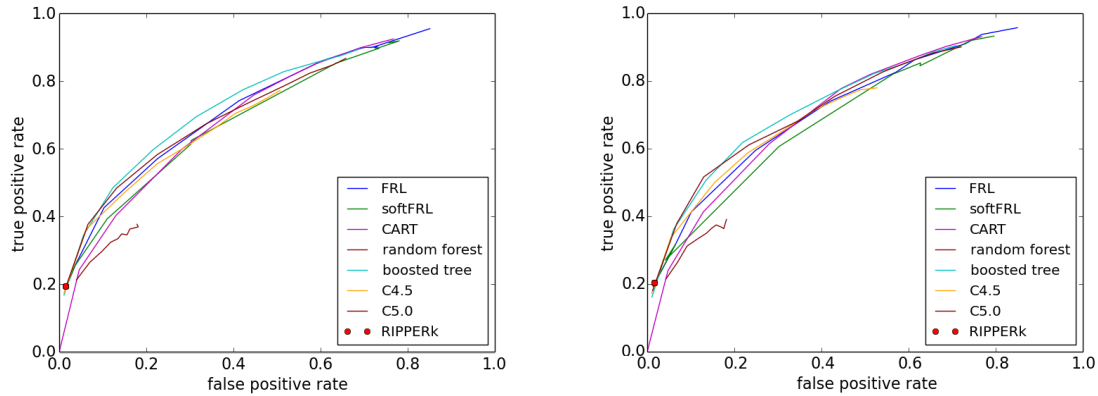
When  $C_1$  becomes larger, the softly falling rule list created by Algorithm softFRL exhibits less “pulling down” of the empirical positive proportion. This is consistent with our expectation that when  $C_1$  is larger, the penalty for violating monotonicity is higher and the algorithm will less likely choose a rule list that frequently violates monotonicity.

### 13 Additional Experiments Comparing Algorithm FRL and Algorithm softFRL to Other Classification Algorithms

Figure 3 shows the ROC curves on the test set using different values of  $w$ , for four additional training-test splits. As we can see, the curves in Figure 3 lie close to each other, again demonstrating the effectiveness of our algorithms in producing falling rule lists that, when used as classifiers, are comparable with classifiers produced by other widely used classification algorithms, in a cost-sensitive setting.



(a) ROC curves on the test set using different  $w$  values for the first additional training-test split (b) ROC curves on the test set using different  $w$  values for the second additional training-test split



(c) ROC curves on the test set using different  $w$  values for the third additional training-test split (d) ROC curves on the test set using different  $w$  values for the fourth additional training-test split

Figure 3: ROC curves on the test set using different  $w$  values for four additional training-test splits

## 14 Additional Experiments Comparing Bayesian Approach to Our Optimization Approach

We conducted a set of experiments comparing the Bayesian approach to our optimization approach. We trained falling rule lists on the entire bank-full dataset using both the Bayesian approach and our optimization approach (Algorithm FRL), and plotted the weighted training loss over real runtime. In particular, for each positive class weight  $w \in \{1, 3, 5, 7\}$ , we set the threshold to  $1/(1+w)$  (By Theorem 2.8, this is the threshold with the least weighted training loss for any given rule list), and computed the weighted training loss using this threshold. For the Bayesian approach, we recorded the runtime and computed the weighted training loss for every 100 iterations of Markov chain Monte-Carlo sampling with simulated annealing, up to 6000 iterations. For our optimization approach, we ran Algorithm FRL for 3000 iterations and recorded the runtime and the weighted training loss whenever the algorithm finds a falling rule list with a smaller (regularized) weighted training loss. Since we want to focus our experiments on the efficiency of searching the model space, the runtimes recorded do not include the time for mining the antecedents. Due to the random nature of both approaches, the experiments were repeated several times.

Figures 4 to 7 show the plots of the weighted training loss over real runtime for the Bayesian approach and our optimization approach (Algorithm FRL), for four additional runs of the same algorithms. Due to the random nature of both approaches, it is sometimes possible that our approach (Algorithm FRL) may find in 3000 iterations a falling rule list with a slightly larger weighted training loss, compared to the Bayesian

approach with 6000 iterations (see Figure 6d). However, in general, our approach tends to find a falling rule list with a smaller weighted training loss faster, due to aggressive pruning of the search space.

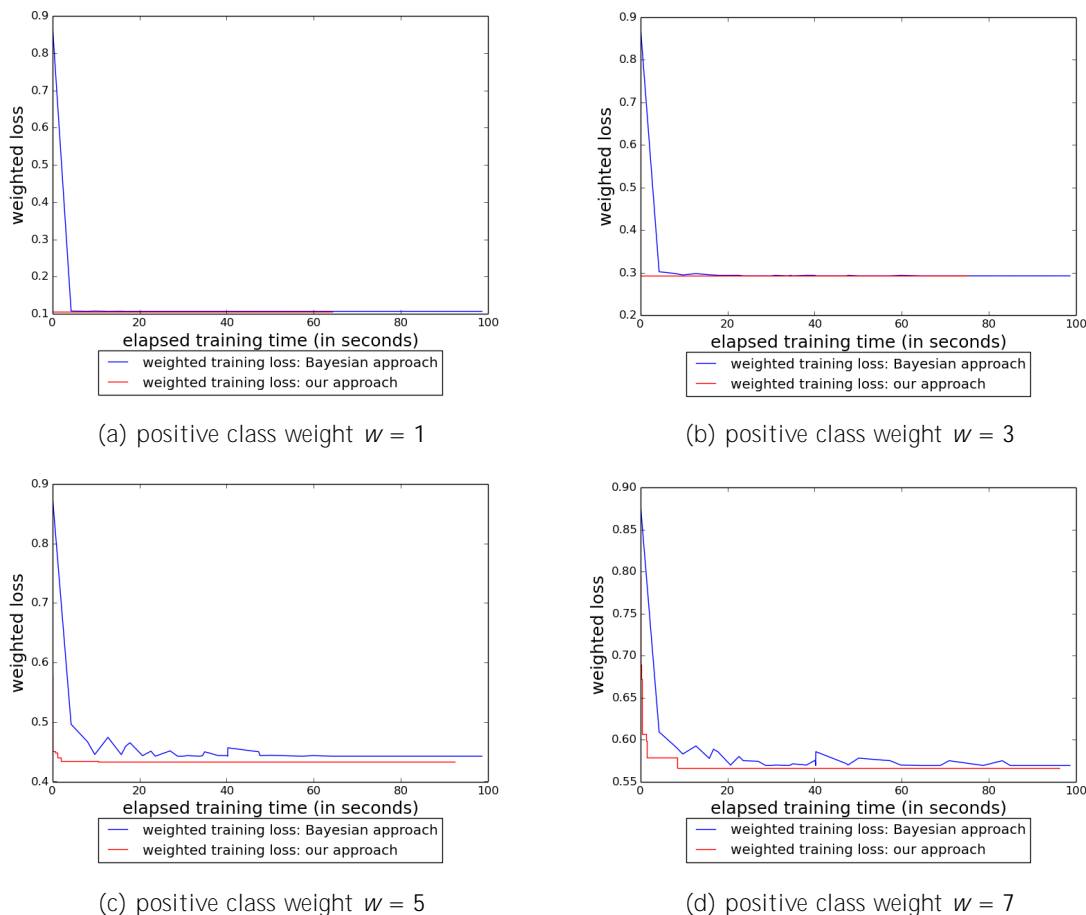


Figure 4: Plots of the weighted training loss over real runtime for the Bayesian approach and our optimization approach (Algorithm FRL): first additional run

It is worth pointing out that both the Bayesian approach and our optimization approach produce similar falling rule lists. Table 19 shows a falling rule list for the bank-full dataset, obtained in a particular run of the Bayesian approach with 6000 iterations. Table 20 shows a falling rule list for the same dataset, obtained in a particular run of Algorithm FRL with 3000 iterations and the positive class weight  $w = 7$ . As we can see, the top four rules in both falling rule lists are identical. Tables 21 and 22 show another pair of falling rule lists obtained using both approaches in different runs, and in this case, both approaches have identified some common rules for a high chance of marketing success. This means that both the Bayesian approach and our optimization approach tend to identify similar conditions that are significant, but our approach has the added advantage of faster training convergence over the Bayesian approach in general.

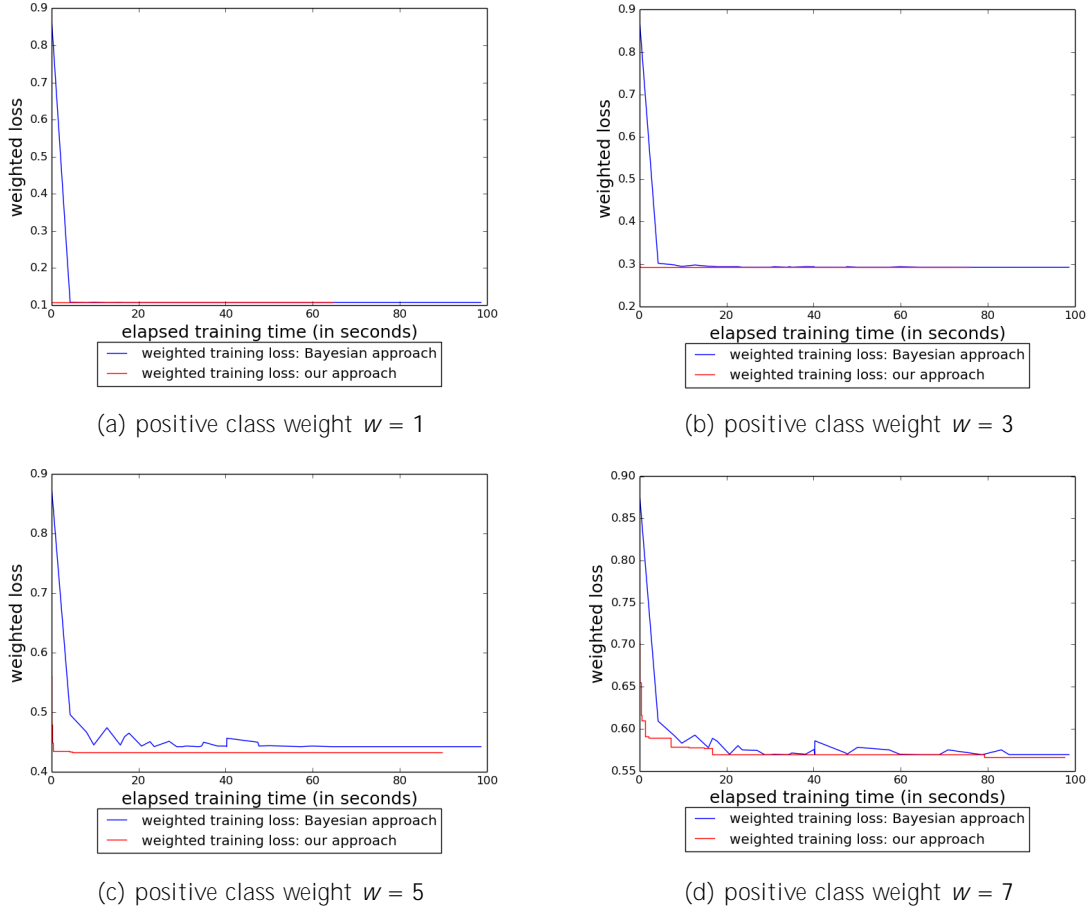


Figure 5: Plots of the weighted training loss over real runtime for the Bayesian approach and our optimization approach (Algorithm FRL): second additional run

	antecedent		probability	positive support	negative support
IF	poutcome=success AND default=no	THEN success prob. is	0.65	978	531
ELSE IF	$60 \leq \text{age} < 100$ AND loan=no	THEN success prob. is	0.29	426	1030
ELSE IF	$17 \leq \text{age} < 30$ AND housing=no	THEN success prob. is	0.25	504	1539
ELSE IF	campaign=1 AND housing=no	THEN success prob. is	0.15	787	4471
ELSE IF	education=tertiary AND housing=no	THEN success prob. is	0.12	460	3313
ELSE IF	marital=single AND contact=cellular	THEN success prob. is	0.11	550	4331
ELSE IF	contact=cellular	THEN success prob. is	0.08	1080	12709
ELSE		success prob. is	0.04	504	11998

Table 19: Falling rule list for bank-full dataset, trained using the Bayesian approach with 6000 iterations.

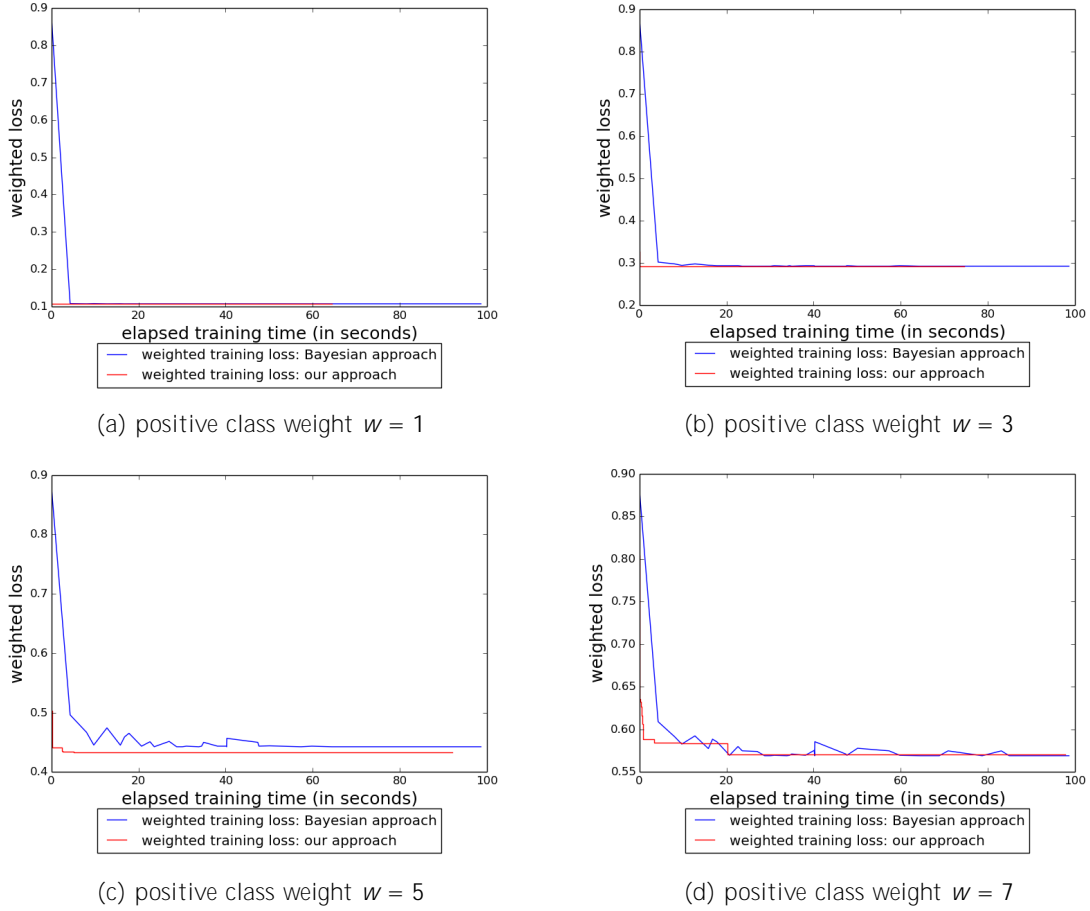
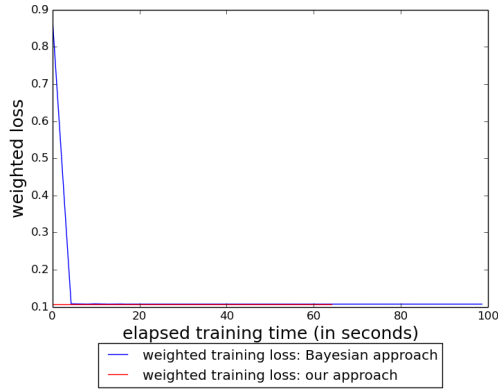


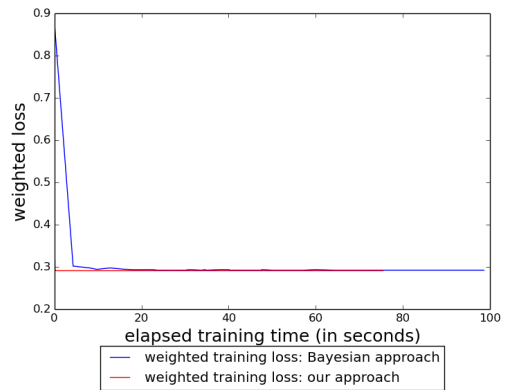
Figure 6: Plots of the weighted training loss over real runtime for the Bayesian approach and our optimization approach (Algorithm FRL): third additional run

	antecedent		probability	positive support	negative support
IF	poutcome=success AND default=no	THEN success prob. is	0.65	978	531
ELSE IF	$60 \leq \text{age} < 100$ AND loan=no	THEN success prob. is	0.29	426	1030
ELSE IF	$17 \leq \text{age} < 30$ AND housing=no	THEN success prob. is	0.25	504	1539
ELSE IF	campaign=1 AND housing=no	THEN success prob. is	0.15	787	4471
ELSE		success prob. is	0.07	2594	32351

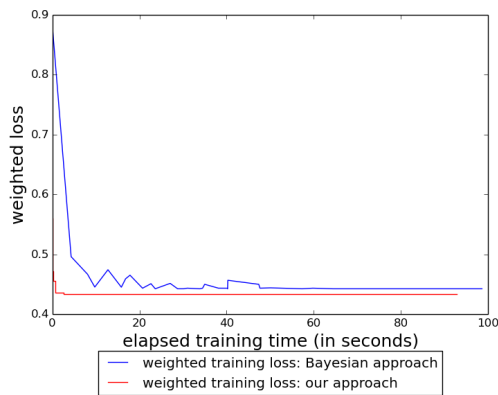
Table 20: Falling rule list for bank-full dataset, trained using the optimization approach (Algorithm FRL) with 3000 iterations and the positive class weight  $w = 7$ .



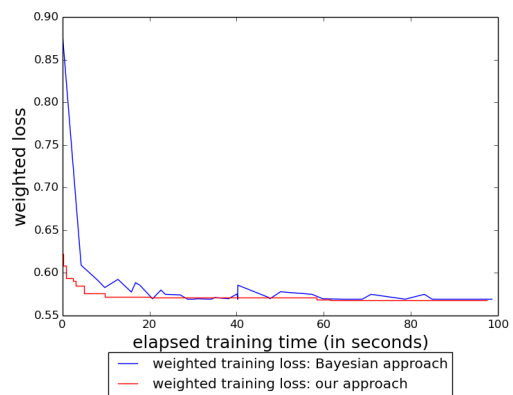
(a) positive class weight  $w = 1$



(b) positive class weight  $w = 3$



(c) positive class weight  $w = 5$



(d) positive class weight  $w = 7$

Figure 7: Plots of the weighted training loss over real runtime for the Bayesian approach and our optimization approach (Algorithm FRL): fourth additional run

---

	antecedent		probability	positive support	negative support
IF	poutcome=success AND housing=no	THEN success prob. is	0.70	729	311
ELSE IF	poutcome=success	THEN success prob. is	0.53	249	222
ELSE IF	60 ≤ age < 100 AND loan=no	THEN success prob. is	0.29	426	1030
ELSE IF	17 ≤ age < 30 AND housing=no	THEN success prob. is	0.25	504	1538
ELSE IF	education=tertiary AND housing=no	THEN success prob. is	0.14	790	4750
ELSE IF	marital=single AND contact=cellular	THEN success prob. is	0.12	648	4754
ELSE IF	1000 ≤ balance < 2000 AND housing=no	THEN success prob. is	0.11	135	1061
ELSE IF	campaign=1 AND contact=cellular	THEN success prob. is	0.10	571	4904
ELSE IF	contact=cellular AND loan=no	THEN success prob. is	0.08	587	6800
ELSE		success prob. is	0.04	650	14552

Table 21: Falling rule list for bank-full dataset, trained using the Bayesian approach with 6000 iterations.

	antecedent		probability	positive support	negative support
IF	poutcome=success AND housing=no	THEN success prob. is	0.70	729	311
ELSE IF	poutcome=success AND previous ≥ 2	THEN success prob. is	0.55	185	154
ELSE IF	poutcome=success AND default=no	THEN success prob. is	0.48	64	68
ELSE IF	60 ≤ age < 100 AND loan=no	THEN success prob. is	0.29	426	1030
ELSE IF	previous ≥ 2 AND housing=no	THEN success prob. is	0.25	302	921
ELSE IF	17 ≤ age < 30 AND housing=no	THEN success prob. is	0.24	444	1413
ELSE IF	education=tertiary AND housing=no	THEN success prob. is	0.13	671	4435
ELSE		success prob. is	0.07	2468	31590

Table 22: Falling rule list for bank-full dataset, trained using the optimization approach (Algorithm FRL) with 3000 iterations and the positive class weight  $w = 7$ .