

## A Parameter Setups

We set the parameters of the proposed training algorithm, CORELS, and SBRL, as follows.

**Proposed Algorithm** We tuned the penalty parameter  $\alpha$  in the training objective function (4) based on Algorithm 2. For all the data sets, we mined rules containing maximally two conditions. When the 16 GB memory on our machine was found to be insufficient for rule mining, we mined rules on a subset of the data set by reducing the number of observations. We searched for the optimal  $\alpha$  from the candidates ranging from  $10^{-4}$  to  $10^{-2}$ . We determined the optimal  $\alpha$  as the one that maximized the AUTAC on the training set  $D$ , under the constraint that number of conditions is less than 20.

---

### Algorithm 2 CRL Tuning Strategy

---

```

//Initial setting
 $I_{\max} \leftarrow 20$   $\triangleright$ max rules
 $C_{\max} \leftarrow 2$   $\triangleright$ max conditions
 $A \leftarrow [.01, .005, .001, .0008, .0005, .0002, .0001]$ 
 $\triangleright$ candidates of  $\alpha$ 
 $\alpha_{\text{opt}} \leftarrow 0$ ,  $\text{AUTAC}_{\text{opt}} \leftarrow 0$ 
FP-Growth minimal support  $\leftarrow 0.05$ 

//Tuning rule mining
 $b \leftarrow |D|$   $\triangleright$ # of observations for mining
while memory insufficient do
   $b \leftarrow \lfloor 0.9b \rfloor$ 
  mine rules on  $b$  observations with FP-Growth
end while

//Tuning  $\alpha$ 
for  $\alpha \in A$  do
   $M, \text{AUTAC} \leftarrow \text{train}(D, \alpha, C_{\max})$ 
  if  $M < I_{\max}$  then
    if  $\text{AUTAC} > \text{AUTAC}_{\text{opt}}$  then
       $\text{AUTAC}_{\text{opt}} \leftarrow \text{AUTAC}$ 
       $\alpha_{\text{opt}} \leftarrow \alpha$ 
    end if
  end if
end for
Output  $\alpha_{\text{opt}}$ 

```

---

**CORELS** For CORELS, we used an implementation publicly available at <https://github.com/fingoldin/pycorels>. We tuned the maximal number of iterations  $N$  and policy  $P$  based on Algorithm 3. We increment  $N$  by  $50k$  each time until  $30k$  and search the best policy  $P$  in list  $L$ . For all the data sets, we mined rules containing maximally two conditions. We determined the optimal  $N$  and  $P$  when predictive accuracy ACC is maximized on the testing set, under the constraint that condition number is less than 20.

---

### Algorithm 3 CORELS Tuning Strategy

---

```

//Initial setting
 $I_{\max} \leftarrow 20$   $\triangleright$ max rules
 $C_{\max} \leftarrow 2$   $\triangleright$ max conditions
 $N \leftarrow 100k$   $\triangleright$ maximum number of iterations
 $L \leftarrow [\text{curious}, \text{lowerbound}, \text{dfs}, \text{bfs}, \text{objective}]$ 
 $\triangleright$ candidates of  $P$ 
 $N_{\text{opt}} \leftarrow 0$ ,  $\text{ACC}_{\text{opt}} \leftarrow 0$ ,  $P_{\text{opt}} \leftarrow \text{None}$ 

//Tuning  $N$  and  $P$ 
for  $t = 1, 2, \dots, 5$  do
  for  $P \in L$  do
     $M, \text{ACC} \leftarrow \text{train}(D, N, P, C_{\max})$ 
    if  $M < I_{\max}$  then
      if  $\text{ACC} > \text{ACC}_{\text{opt}}$  then
         $\text{ACC}_{\text{opt}} \leftarrow \text{ACC}$ 
         $N_{\text{opt}} \leftarrow N$ 
         $P_{\text{opt}} \leftarrow P$ 
      end if
    end if
  end for
   $N \leftarrow N + 50k$ 
end for
output  $N_{\text{opt}}, P_{\text{opt}}$ 

```

---

**SBRL** For SBRL, we used an implementation publicly available at <https://github.com/Hongyuy/sbrrlmod>. We tuned number of chains  $Nc$  and the expected length of the rule list  $\lambda$  based on Algorithm 4. We mine rules maximally containing two rules on messidor, german, adult, magic and coupon. When the 16 GB memory on our machine was found to be insufficient for rule mining, we increment both positive and negative minimal support  $S_+$  and  $S_-$  by 0.05. When the minimal support is higher than an threshold  $T$ , we turn to mine rules maximally containing one rules. We mine rules maximally containing one rules on juvenile, frisk and recidivism. We determined the optimal  $\lambda$  and  $Nc$  as the ones that maximized predictive accuracy ACC on the testing set, under the constraint that number of conditions is less than 20.

---

**Algorithm 4** SBRL Tuning Strategy
 

---

```

//Initial setting
 $I_{\max} \leftarrow 20$   $\triangleright$ max rules
 $C_{\max} \leftarrow 2$   $\triangleright$ max conditions
 $T \leftarrow 0.7$   $\triangleright$ threshold to reduce  $C_{\max}$ 
 $\Lambda \leftarrow [1, 2, 5, 10, 15, 20, 25, 30]$ 
 $\triangleright$ candidates of  $\lambda$ 
 $Nc \leftarrow 0$   $\triangleright$ number of chains
 $\lambda_{\text{opt}} \leftarrow 0, Nc_{\text{opt}} \leftarrow 0, ACC_{\text{opt}} \leftarrow 0$ 
 $S_+ \leftarrow 0.05, S_- \leftarrow 0.05$ 
 $\triangleright$ minimal support for pos and neg rules

//Tuning minimal support
while memory insufficient do
   $S_+ \leftarrow S_+ + 0.05$ 
   $S_- \leftarrow S_- + 0.05$ 
  if  $S_+ > T$  then
     $C_{\max} \leftarrow 1$ 
  end if
end while

//Tuning  $\lambda$  and  $Nc$ 
for  $\lambda \in \Lambda$  do
  for  $t = 1 \dots 6$  do
     $M, ACC \leftarrow \text{train}(D, \lambda, Nc, S_+, S_-, C_{\max})$ 
     $Nc \leftarrow Nc + 5$ 
    if  $M < I_{\max}$  then
      if  $ACC > ACC_{\text{opt}}$  then
         $ACC_{\text{opt}} \leftarrow ACC$ 
         $\lambda_{\text{opt}} \leftarrow \lambda$ 
         $Nc_{\text{opt}} \leftarrow Nc$ 
      end if
    end if
  end for
end for
output  $\lambda_{\text{opt}}, Nc_{\text{opt}}$ 

```

---

## B Exhaustive Results

Here, we show the results for AdaBoost and XGBoost we omitted in Section 6 due to space limitation. Figures 5 and 6 show the trade-off curves, and Tables 4 and 5 show AUTACs. These results also confirm the validity of the proposed training algorithm.

Table 4: AUTACs on AdaBoost: The numbers in the parenthesis denote standard deviations. Bold fonts denote the best results (underlined), and the results which was not significantly different from the best result (t-test with the 5% significance level).

	CRL	CORELS	SBRL
messidor	<b>.675 (.010)</b>	<b>.674 (.013)</b>	<b>.660 (.020)</b>
german	<b>.754 (.010)</b>	.738 (.013)	<b>.749 (.005)</b>
adult	<b>.856 (.002)</b>	.837 (.004)	.850 (.004)
juvenile	<b>.892 (.005)</b>	<b>.885 (.013)</b>	.883 (.005)
frisk	<b>.685 (.003)</b>	.682 (.003)	.678 (.002)
recidivism	<b>.763 (.006)</b>	.744 (.007)	.759 (.006)
magic	<b>.864 (.007)</b>	.841 (.010)	.852 (.007)
coupon	<b>.740 (.012)</b>	.714 (.015)	<b>.743 (.017)</b>

Table 5: AUTACs on XGBoost: The numbers in the parenthesis denote standard deviations. Bold fonts denote the best results (underlined), and the results which was not significantly different from the best result (t-test with the 5% significance level).

	CRL	CORELS	SBRL
messidor	<b>.689 (.017)</b>	<b>.681 (.019)</b>	<b>.670 (.022)</b>
german	<b>.748 (.022)</b>	<b>.732 (.016)</b>	<b>.734 (.010)</b>
adult	<b>.856 (.002)</b>	.837 (.004)	.851 (.004)
juvenile	<b>.898 (.006)</b>	<b>.888 (.015)</b>	.884 (.005)
frisk	<b>.686 (.003)</b>	.683 (.003)	.679 (.003)
recidivism	<b>.766 (.007)</b>	.744 (.008)	.759 (.006)
magic	<b>.863 (.004)</b>	.840 (.007)	.853 (.004)
coupon	<b>.739 (.013)</b>	.713 (.015)	<b>.744 (.017)</b>

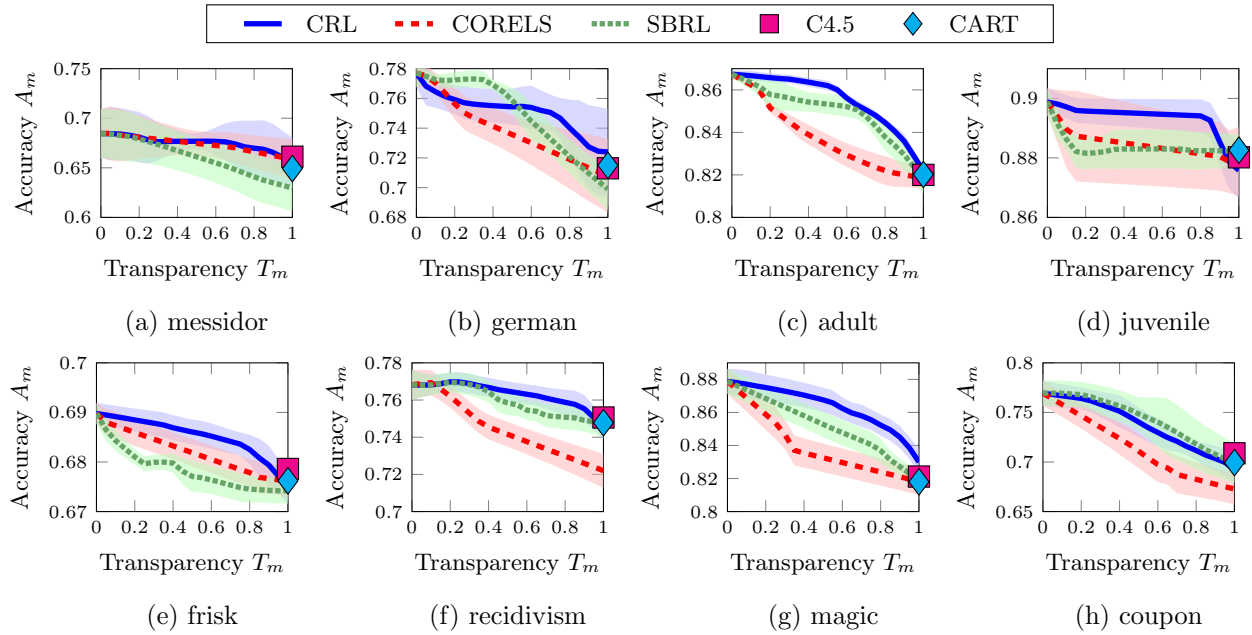


Figure 5: The transparency–accuracy trade-off (with AdaBoost as  $f_b$ ): The solid lines denote average trade-off curves, while the shaded regions denote  $\pm$  standard deviations evaluated via 5-fold cross validation.

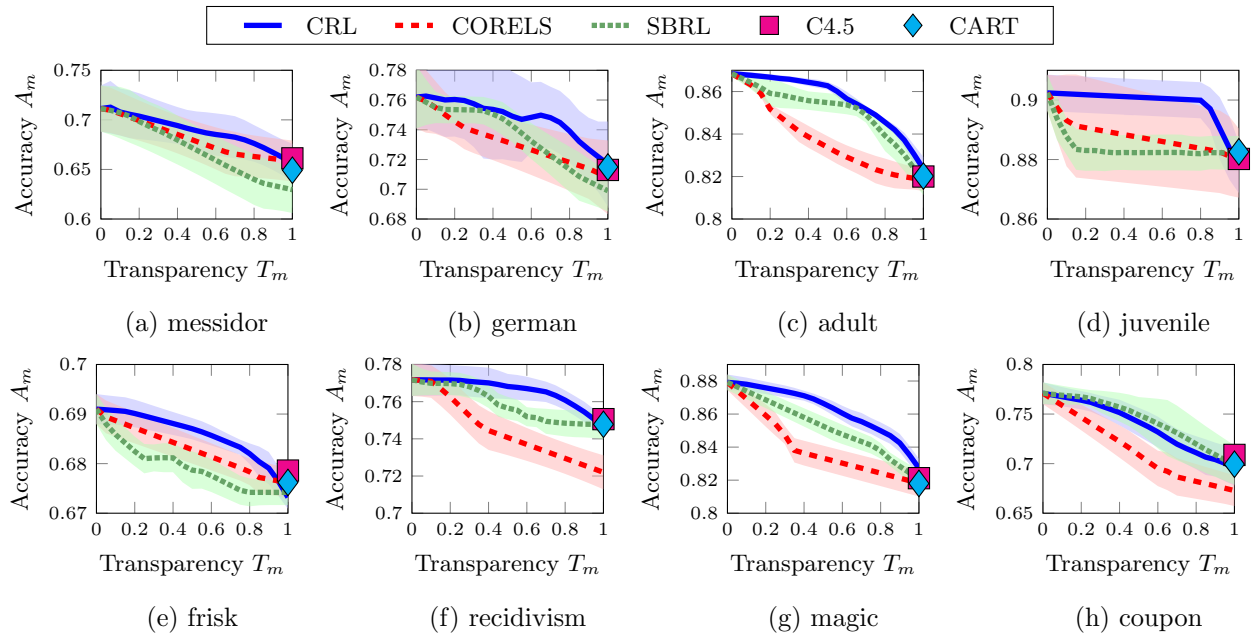


Figure 6: The transparency–accuracy trade-off (with XGBoost as  $f_b$ ): The solid lines denote average trade-off curves, while the shaded regions denote  $\pm$  standard deviations evaluated via 5-fold cross validation.

## C Survey

### C.1 Survey Question

Figure 7 is an example of our survey question.

### C.2 Additional Results

We analyzed the trade-off curve in different groups of users, by gender and by age, and summarized the results in Figures 8 and 9. The results suggest that gender does not seem to play any role in the preference for interpretability but the younger group (30 years old or younger) are more likely to choose rules over black-box models than the older group (older than 30 years old). These results suggest that the users indeed have their own preference on the trade-off between the accuracy and transparency. Thus, it is essential to provide the users a freedom of choosing the trade-off depending on their preference.

Now you are a judge and would like to use machine learning model to predict recidivism of criminals and use that to determine the sentence of the criminals.

The information about this person is provided below.

Number of conviction charges	Reason for being transferred	Weighted court assessment	Number of probation officers	Record of time to first hearing	Record of time from violation to hearing	Status of compliance with house arrest order	Risk score
5	Absconded	300	2	NA	NA	NA	10

We provide two options for you. One option is a rule which tells you why the prediction is such and the other is a black-box model which does not tell you how a prediction is generated. The estimated accuracy of both models are provided. Which one would you prefer to use and trust?

Model	Explanation	Estimated Accuracy
1	In collected data, 66.2% of people <b>whose number of probation officers is higher than 1 and the record of time from violation to hearing is NA</b> will not re-offend in the future	66.2%
2	Unknown	69.3%

Model 1

Model 2

Figure 7: An example of our survey questions.

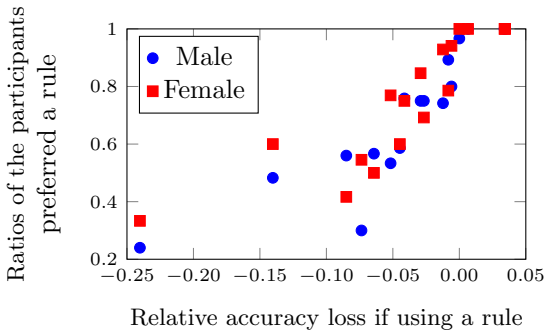


Figure 8: Human evaluated trade-off between model transparency and accuracy: Male vs. Female

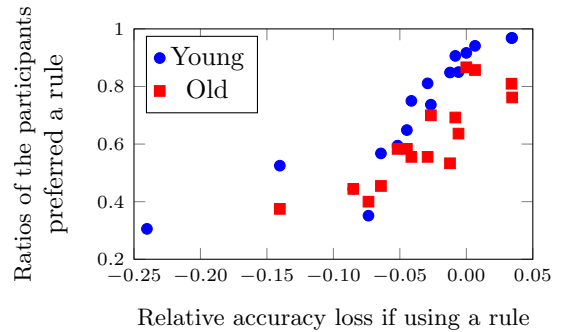


Figure 9: Human evaluated trade-off between model transparency and accuracy: Young vs. Old