# Interpretable Machine Learning for Recidivism Prediction

Cynthia Rudin

Department of Computer Science

Department of Electrical and Computer Engineering
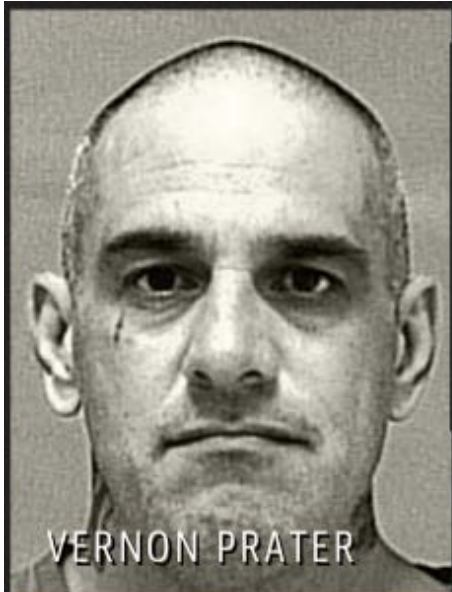
Duke University

joint work with Berk Ustun, Jiaming Zeng, Elaine Angelino, Daniel Alabi, Nicholas Larus-Stone, Margo Seltzer, and Hima Lakkaraju

# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016

COMPAS: Correctional Offender Management Profiling for Alternative Sanctions

# A computer program used for bail and sentencing decisions was labeled biased against blacks. It's actually not that clear.

By Sam Corbett-Davies, Emma Pierson, Avi Feller and Sharad Goel  October 17

## COMPAS may still be biased, but we can't tell.

Northpointe has refused to disclose the details of its proprietary algorithm, making it impossible to fully assess the extent to which it may be unfair, however inadvertently. That's understandable: Northpointe needs to protect its bottom line. But it raises questions about relying on for-profit companies to develop risk assessment tools.

I point if person
has social type
with below
average parole
violation rate

| SOCIAL TYPE | VIOLATION RATE |
|---|---|
| All persons.......................................................... | 26.5% |
| Ne'er-do-well........................................................ | 25.6 |
| Mean citizen......................................................... | 30.0 |
| Drunkard............................................................ | 38.9 |
| Gangster............................................................. | 23.2 |
| Recent immigrant.................................................... | 16.7 |
| Farm boy............................................................ | 10.2 |
| Drug addict......................................................... | 66.7 |

total score
over all 21
significant factors
predicts
success at parole

| POINTS FOR NUMBER OF FACTORS | Per Cent Non-violators of Parole |
|---|---|
| 16-21 | 98.5 |
| 14-15 | 97.8 |
| 13 | 91.2 |
| 12 | 84.9 |
| 11 | 77.3 |
| 10 | 65.9 |
| 7-9 | 56.1 |
| 5-6 | 32.9 |
| 2-4 | 24.0 |

Burgess. Factors determining success or failure on parole. 1928

| FACTOR | Score * |
|---|---|
| **Gender** | |
| Female | 0 |
| Male | 1 |
| **Age** | |
| Less than 24 | 3 |
| 24-29 | 2 |
| 30-49 | 1 |
| 50+ | 0 |
| **County** | |
| Rural counties | 0 |
| Smaller, urban count | 1 |
| Allegheny and | |
| Philadelphia | 2 |
| Counties | |
| **Total number of prior arrests** | |
| 0 | 0 |
| 1 | 1 |
| 2 to 4 | 2 |
| 5 to 12 | 3 |
| 13+ | 4 |
| **Prior property arrests** | |
| No | 0 |
| Yes | 1 |
| **Prior drug arrests** | |
| No | 0 |
| Yes | 1 |
| **Property offender** | |
| No | 0 |
| Yes | 1 |
| **Offense gravity score (OGS)** | |
| 4+ | 0 |

| Risk score | N | % Arrested |
|---|---|---|
| 0 | 3 | 0.0 |
| 1 | 47 | 17.0 |
| 2 | 181 | 9.9 |
| 3 | 436 | 23.6 |
| 4 | 737 | 24.8 |
| 5 | 1,036 | 32.4 |
| 6 | 1,067 | 40.7 |
| 7 | 1,434 | 47.2 |
| 8 | 1,934 | 55.5 |
| 9 | 2,103 | 62.3 |
| 10 | 1,829 | 69.9 |
| 11 | 1,098 | 72.2 |
| 12 | 278 | 79.1 |
| 13 | 25 | 80.0 |
| 14 | 3 | 66.7 |

Pennsylvania Commission on Sentencing, 2013

1. Lived with both biological parents to age 16 (except for death of parent):
Yes ............................................... -2
No ............................................... +3
Evidence:

2. Elementary School Maladjustment:
No Problems.............................................. -1
Slight (Minor discipline or attendance) or Moderate Problems............................. +2
Severe Problems (Frequent disruptive behavior and/or attendance or behavior resulting in expulsion or serious suspensions) ............................................. +5
(Same as CATS Item)

3. History of alcohol problems *(Check if present):*
˜ Parental Alcoholism        ˜ Teenage Alcohol Problem
˜ Adult Alcohol Problem      ˜ Alcohol involved in prior offense
˜ Alcohol involved in index offense
        No boxes checked..................................... -1
        1 or 2 boxes checked .............................. . 0
        3 boxes checked ...................................... +1
        4 or 5 boxes checked .............................. +2
        Evidence:

4. Marital status (at the time of or prior to index offense):
Ever married (or lived common law in the same home for at least six months) ......... -2
Never married........................................... +1
Evidence:

5. Criminal history score for nonviolent offenses prior to the index offense
Score 0 ....................................................... -2
Score 1 or 2.................................................   0
Score 3 or above ...................................... +3
(from the Cormier-Lang system, see below)

6. Failure on prior conditional release (includes parole or probation violation or revocation, failure to comply, bail violation, and any new arrest while on conditional release):
No............................................................0
Yes ......................................................... +3
Evidence:

7. Age at index offense
Enter Date of Index Offense: ___/___/_____
Enter Date of Birth: ___/___/_____
Subtract to get Age:
39 or over .................................................. -5
34 - 38 ...................................................... -2
28 - 33 ...................................................... -1
27 ...............................................................0
26 or less.................................................. +2

8. Victim Injury (for index offense; the most serious is scored):
Death............................................................ -2
Hospitalized.................................................0
Treated and released............................... +1
None or slight (includes no victim)........... +2
Note: admission for the gathering of forensic evidence only is NOT considered as either treated or hospitalized; ratings should be made based on the degree of injury.
Evidence:

9. Any female victim (for index offense)
Yes ............................................................ -1
No (includes no victim)............................. +1
Evidence:

10. Meets DSM criteria for any personality disorder (must be made by appropriately licensed or certified professional)
No................................................................ -2
Yes ............................................................ +3
Evidence:

11. Meets DSM criteria for schizophrenia (must be made by appropriately licensed or certified professional)
Yes ............................................................ -3
No .............................................................. +1
Evidence:

12. a. Psychopathy Checklist score (if available, otherwise use item 12.b. CATS score)........
4 or under ................................................. -3
5 – 9........................................................... -3
10-14 ......................................................... -1
15-24 ..........................................................   0
25-34 ......................................................... +4
35 or higher ........................................... +12
Note: If there are two or more PCL scores, average the scores.
Evidence:

12. b. CATS score (from the CATS worksheet)
0 or 1 ........................................................ -3
2 or 3 .........................................................0
4 ................................................................+2
5 or higher ............................................... +3

12. WEIGHT (Use the highest circled weight from 12 a. or 12 b.) ......................... _____

**TOTAL VRAG SCORE (SUM CIRCLED SCORES FOR ITEMS 1 – 11 PLUS THE WEIGHT FOR ITEM 12): _____**

| VRAG Score | Category of Risk |
|---|---|
|  |  |
| -24 | Low |
| -23 | Low |
| -22 | Low |
| -20 | Low |
| -19 | Low |
| -18 | Low |
| -17 | Low |
| -16 | Low |
| -15 | Low |
| -14 | Low |
| -13 | Low |
| -12 | Low |
| -11 | Low |
| -10 | Low |
| -9 | Low |
| -8 | Low |
| -7 | Medium |
| -6 | Medium |
| -5 | Medium |
| -4 | Medium |
| -3 | Medium |
| -2 | Medium |
| -1 | Medium |
| 0 | Medium |
| 1 | Medium |
| 2 | Medium |
| 3 | Medium |
| 4 | Medium |
| 5 | Medium |
| 6 | Medium |
| 7 | Medium |
| 8 | Medium |
| 9 | Medium |
| 10 | Medium |
| 11 | Medium |
| 2 | Medium |
| 13 | Medium |
| 14 | High |
| 15 | High |
| 16 | High |
| 17 | High |
| 18 | High |
| 19 | High |
| 20 | High |
| 21 | High |
| 22 | High |
| 23 | High |
| 24 | High |
| 25 | High |
| 26 | High |
| 28 | High |
| 32 | High |

Violence Risk Appraisal Guide (Quinsey et al, 2006)

**Is there a principled way to create scoring systems?**

Should we have experts create it and validate it afterwards?

Should we do manual feature selection and round logistic regression coefficients?

Should we actually solve it?

# Supersparse Linear Integer Models (SLIM)

$$\min_{\lambda \in \mathcal{L}} \left( C_+ \frac{1}{n_+} \sum_{i:y_i=1} 1_{(\mathbf{x}^T\lambda) \leq 0} + C_- \frac{1}{n_-} \sum_{i:y_i=-1} 1_{(\mathbf{x}^T\lambda) \geq 0} \right) + C_0 \| \lambda \|_0 + \epsilon \| \lambda \|_1$$

**Accuracy**

**Sparsity**

**Co-prime Coefficients**

(2,2,4,2,6) ➜ (1,1,2,1,3)

$\lambda \in \mathcal{L}$ means that $\forall j, \ \lambda_j \in \{-10, -9, ..., 0, ..., 9, 10\}$

**Meaningful Coefficients**

# Supersparse Linear Integer Models (SLIM)

$$\min_{\lambda \in \mathcal{L}} \left( \underbrace{C_+ \frac{1}{n_+} \sum_{i:y_i=1} 1_{(\mathbf{x}^T \lambda) \leq 0} + C_- \frac{1}{n_-} \sum_{i:y_i=-1} 1_{(\mathbf{x}^T \lambda) \leq 0}}_{\textbf{Accuracy}} \right) + \underbrace{C_0 \| \lambda \|_0}_{\textbf{Sparsity}} + \underbrace{\epsilon \| \lambda \|_1}_{\substack{\textbf{Co-prime}\\\textbf{Coefficients}}}$$

$\lambda \in \mathcal{L}$ means that $\forall j, \ \lambda_j \in \{-10, -9, ..., 0, ..., 9, 10\}$ **Meaningful Coefficients**

How much training accuracy do I sacrifice for one fewer term in the model? $C_0$

How much training accuracy do I trade for co-prime coefficients? Provably none.

Could there be a sparser model with equivalent training accuracy? Provably not.

# Supersparse Linear Integer Models (SLIM)

$$\min_{\lambda \in \mathcal{L}} \left( \underbrace{C_+ \frac{1}{n_+} \sum_{i:y_i=1} 1_{(\mathbf{x}^T\lambda) \leq 0} + C_- \frac{1}{n_-} \sum_{i:y_i=-1} 1_{(\mathbf{x}^T\lambda) \leq 0}}_{\textbf{Accuracy}} \right) + \underbrace{C_0 \| \lambda \|_0}_{\textbf{Sparsity}} + \underbrace{\epsilon \| \lambda \|_1}_{\substack{\textbf{Co-prime} \\ \textbf{Coefficients}}}$$

$$\lambda \in \mathcal{L} \text{ means that } \forall j, \ \lambda_j \in \{-10,-9,...,0,...,9,10\} \quad \substack{\textbf{Meaningful} \\ \textbf{Coefficients}}$$

Can I get a model that is optimal for a particular sensitivity/specificity (TP/FP) tradeoff?

# Supersparse Linear Integer Models (SLIM)

$$\min_{\lambda \in \mathcal{L}} \left( \underbrace{C_+ \frac{1}{n_+} \sum_{i:y_i=1} 1_{(\mathbf{x}^T \lambda) \le 0} + C_- \frac{1}{n_-} \sum_{i:y_i=-1} 1_{(\mathbf{x}^T \lambda) \le 0}}_{\textbf{Accuracy}} \right) + \underbrace{C_0 \| \lambda \|_0}_{\textbf{Sparsity}} + \underbrace{\epsilon \| \lambda \|_1}_{\substack{\textbf{Co-prime} \\ \textbf{Coefficients}}}$$

$\lambda \in \mathcal{L}$ means that $\forall j, \ \lambda_j \in \{-10, -9, ..., 0, ..., 9, 10\}$   **Meaningful Coefficients**

Does Lasso+rounding give the same result?

    No. Can be a lot worse.

# SLIM MIP

$$\min_{\boldsymbol{\lambda},\psi,\boldsymbol{\Phi},\alpha,\beta} \frac{1}{N}\sum_{i=1}^{N}\psi_i + \sum_{j=1}^{P}\Phi_j$$

s.t.

$$M_i\psi_i \geq \gamma - \sum_{j=0}^{P} y_i\lambda_j x_{i,j} \qquad i=1,...,N \text{ \textit{0-1 loss}}$$

$$\Phi_j = C_0\alpha_j + \epsilon\beta_j \qquad j=1,...,P \text{ \textit{int. penalty}}$$

$$-\Lambda_j\alpha_j \leq \lambda_j \leq \Lambda_j\alpha_j \qquad j=1,...,P \text{ \textit{$\ell_0$ norm}}$$

$$-\beta_j \leq \lambda_j \leq \beta_j \qquad j=1,...,P \text{ \textit{$\ell_1$ norm}}$$

$$\lambda_j \in \mathcal{L}_j \qquad j=0,...,P \text{ \textit{int. set}}$$

$$\psi_i \in \{0,1\} \qquad i=1,...,N \text{ \textit{loss variables}}$$

$$\Phi_j \in \mathbb{R}_+ \qquad j=1,...,P \text{ \textit{int. penalty variables}}$$

$$\alpha_j \in \{0,1\} \qquad j=1,...,P \text{ \textit{$\ell_0$ variables}}$$

$$\beta_j \in \mathbb{R}_+ \qquad j=1,...,P \text{ \textit{$\ell_1$ variables}}$$

(Code is publicly available)

# Recidivism Prediction Problems

*Recidivism of Prisoners Released in 1994* (Source: US DOJ BJS)

$N = 33{,}796$ prisoners tracked for 3 years after release from prison in 1994

$P = 49$ binary input variables

- *male, female*
- *prior_drug_abuse, prior_alcohol_abuse*
- *age_of_1$^{st}$_arrest, age_of_1$^{st}$_confinement, prior_arrests, prior_prison_time*
- *age_at_release, time_served, type of release, infraction_in_prison*

| Prediction Problem | $P(y_i = +1)$ | Outcome (rearrested in 3 year after release) |
|---|---|---|
| arrest | 59.0% | for any crime |
| drug | 20.0% | for drug crime (e.g. possession, trafficking, etc.) |
| general_violence | 19.1% | for violent crime (e.g. robbery, aggravated assault) |
| domestic_violence | 3.5% | for domestic violence crime |
| sexual_violence | 3.0% | for sexual violence crimes |
| fatal_violence | 0.7% | for murder or manslaughter |

# arrest



| Method | Algorithm-AUC |
|--------|---------------|
| SLIM | 0.72 |
| Boosting | 0.74 |
| SVM RBF | 0.72 |
| RF | 0.73 |
| Ridge | 0.73 |
| Lasso | 0.72 |
| C5.0R | 0.72 |
| C5.0T | 0.72 |
| CART | 0.68 |

# general violence



| Method | "AUC" |
|---|---|
| SLIM | 0.71 |
| Boosting | 0.72 |
| SVM RBF | 0.69 |
| RF | 0.71 |
| Ridge | 0.72 |
| Lasso | 0.72 |
| C5.0R | 0.56 |
| C5.0T | 0.57 |
| CART | 0.56 |

# sexual violence



| Method | "AUC" |
|--------|-------|
| SLIM | 0.70 |
| Boosting | 0.70 |
| SVM RBF | 0.70 |
| RF | 0.54 |
| Ridge | 0.72 |
| Lasso | 0.72 |
| C5.0R | 0.50 |
| C5.0T | 0.50 |
| CART | 0.51 |

True Positive Rate

False Positive Rate

# arrest



| Method | "AUC" |
|--------|-------|
| SLIM | 0.72 |
| Boosting | 0.74 |
| SVM RBF | 0.72 |
| RF | 0.73 |
| Ridge | 0.73 |
| Lasso | 0.72 |
| C5.0R | 0.72 |
| C5.0T | 0.72 |
| CART | 0.68 |

## PREDICT ARREST FOR ANY OFFENSE IF SCORE $> 1$

| | | | | |
|---|---|---|---|---|
| 1. | *age_at_release_18_to_24* | 2 points | | · · · · · · |
| 2. | *prior_arrests$\geq$5* | 2 points | + | · · · · · · |
| 3. | *prior_arrest_for_misdemeanor* | 1 point | + | · · · · · · |
| 4. | *no_prior_arrests* | -1 point | + | · · · · · · |
| 5. | *age_at_release$\geq$40* | -1 point | + | · · · · · · |
| **ADD POINTS FROM ROWS 1–5** | | **SCORE** | = | · · · · · · |

**PREDICT** `arrest` **if**

*age_at_release_18_to_24*

**OR** *prior_arrests $\geq 5$* **AND** *age_at_release $\leq 40$*

**OR** *prior_arrests $\geq 5$* **AND** *age_at_release $\geq 40$* AND *misdemeanor*

# domestic violence

**PREDICT ARREST FOR DOMESTIC VIOLENCE OFFENSE IF SCORE $> 3$**

| | | | | |
|---|---|---|---|---|
| 1. | *prior_arrest_for_misdemeanor* | 4 points | | · · · · · · |
| 2. | *prior_arrest_for_felony* | 3 points | + | · · · · · · |
| 3. | *prior_arrest_for_domestic_violence* | 2 points | + | · · · · · · |
| 4. | *age_1st_confinement_18_to_24* | 1 point | + | · · · · · · |
| 5. | *infraction_in_prison* | -5 points | + | · · · · · · |
| | **ADD POINTS FROM ROWS 1-5** | **SCORE** | = | · · · · · · |

Test TPR/FPR:          76.6/44.5%
Validation TPR/FPR:  81.4/48.0%

# general_violence

**PREDICT ARREST FOR GENERAL VIOLENCE OFFENSE IF SCORE > 7**

| | | | | |
|---|---|---|---|---|
| 1. | *prior_arrest_for_general_violence* | 8 points | | · · · · · · |
| 2. | *prior_arrest_for_misdemeanor* | 5 points | + | · · · · · · |
| 3. | *infraction_in_prison* | 3 points | + | · · · · · · |
| 4. | *prior_arrest_for_local_ord* | 3 points | + | · · · · · · |
| 5. | *prior_arrest_for_property* | 2 points | + | · · · · · · |
| 6. | *prior_arrest_for_fatal_violence* | 2 points | + | · · · · · · |
| 7. | *prior_arrest_with_firearms_involved* | 1 point | + | · · · · · · |
| 8. | *age_at_release≥40* | -7 points | + | · · · · · · |
| | **ADD POINTS FROM ROWS 1-8** | **SCORE** | = | · · · · · · |

**Test TPR/FPR:**         76.7/45.4%
**Validation TPR/FPR:**   76.8/47.6%

# sexual_violence

**PREDICT ARREST FOR SEXUAL VIOLENCE OFFENSE IF SCORE $> 2$**

| | | | | |
|---|---|---|---|---|
| 1. | *prior_arrest_for_sexual* | 3 points | | $\cdots\cdots$ |
| 2. | *prior_arrests$\geq$5* | 1 point | $+$ | $\cdots\cdots$ |
| 3. | *multiple_prior_jail_time* | 1 point | $+$ | $\cdots\cdots$ |
| 4. | *prior_arrest_for_multiple_types_of_crime* | -1 point | $+$ | $\cdots\cdots$ |
| 5. | *no_prior_arrests* | -2 points | $+$ | $\cdots\cdots$ |
| | **ADD POINTS FROM ROWS 1-5** | **SCORE** | $=$ | $\cdots\cdots$ |

Test TPR/FPR:          44.3/17.7%
Validation TPR/FPR:  43.7/19.9%

# fatal_violence

**PREDICT ARREST FOR FATAL VIOLENCE OFFENSE IF SCORE** $> 4$

| | | | | |
|---|---|---|---|---|
| 1. | *age_1st_confinement$\leq$17* | 5 points | | $\cdots\cdots$ |
| 2. | *prior_arrest_with_firearms_involved* | 3 points | + | $\cdots\cdots$ |
| 3. | *age_1st_confinement_18_to_24* | 2 points | + | $\cdots\cdots$ |
| 4. | *prior_arrest_for_felony* | 2 points | + | $\cdots\cdots$ |
| 5. | *age_at_release_18_to_24* | 1 point | + | $\cdots\cdots$ |
| 6. | *prior_arrest_for_drugs* | 1 point | + | $\cdots\cdots$ |
| **ADD POINTS FROM ROWS 1-6** | | **SCORE** | = | $\cdots\cdots$ |

Test TPR/FPR:          55.4/35.5%
Validation TPR/FPR:  63.2/42.4%

# Risk Assessment Models

**Decision-Making Model**

**PREDICT ARREST FOR ANY OFFENSE IF SCORE $> 1$**

| | | | | |
|---|---|---|---|---|
| 1. | *age_at_release_18_to_24* | 2 points | | · · · · · · |
| 2. | *prior_arrests$\geq$5* | 2 points | + | · · · · · · |
| 3. | *prior_arrest_for_misdemeanor* | 1 point | + | · · · · · · |
| 4. | *no_prior_arrests* | -1 point | + | · · · · · · |
| 5. | *age_at_release$\geq$40* | -1 point | + | · · · · · · |
| | **ADD POINTS FROM ROWS 1–5** | **SCORE** | = | · · · · · · |

**Risk Assessment Model**

| | | | | |
|---|---|---|---|---|
| 1. | *prior arrests $\geq 2$* | 1 point | | · · · · · · |
| 2. | *prior arrests $\geq 5$* | 1 point | + | · · · · · · |
| 3. | *prior arrests for local ordinance* | 1 point | + | · · · · · · |
| 4. | *age at release 18 to 24* | 1 point | + | · · · · · · |
| 5. | *age at release $\geq 40$* | -1 point | + | · · · · · · |
| | **ADD POINTS FROM ROWS 1-5** | **SCORE** | = | · · · · · · |

| SCORE | -1 | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|---|
| RISK | 11.9% | 26.9% | 50.0% | 73.1% | 88.1% | 95.3% |

# Risk-Calibrated SLIM

Logistic Loss

Model Size

$$\min_{\lambda \in \mathcal{L}} \frac{1}{n} \sum_{i=1}^{n} \log(1 + e^{-\mathbf{x}^T \lambda}) + C_0 \| \lambda \|_0$$

$\lambda \in \mathcal{L}$ means that $\forall j, \ \lambda_j \in \{-10, -9, ..., 0, ..., 9, 10\}$

Small Integer Coefficients

Ustun and Rudin, 2017

# Risk-Calibrated SLIM

Logistic Loss

Model Size

$$\min_{\lambda \in \mathcal{L}} \frac{1}{n} \sum_{i=1}^{n} \log(1 + e^{-\mathbf{x}^T \lambda}) + C_0 \| \lambda \|_0$$

$\lambda \in \mathcal{L}$ means that $\forall j, \ \lambda_j \in \{-10, -9, ..., 0, ..., 9, 10\}$

- Specialized cutting-plane methods
- Scales to large samples

Ustun and Rudin, 2017

# RiskSlim Model for Arrest

| | | | | |
|---|---|---|---|---|
| 1. | *Prior Arrests $\geq 2$* | 1 point | | · · · · · · |
| 2. | *Prior Arrests $\geq 5$* | 1 point | + | · · · · · · |
| 3. | *Prior Arrests for Local Ordinance* | 1 point | + | · · · · · · |
| 4. | *Age at Release between 18 to 24* | 1 point | + | · · · · · · |
| 5. | *Age at Release $\geq 40$* | -1 points | + | · · · · · · |
| | **ADD POINTS FROM ROWS 1–5** | **SCORE** | = | · · · · · · |

| SCORE | -1 | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|---|
| RISK | 11.9% | 26.9% | 50.0% | 73.1% | 88.1% | 95.3% |

Berk Ustun

Jiaming Zeng

Daniel Alabi

Elaine Angelino

Nicholas
Larus-Stone

Margo Seltzer

# Rule List Models (Decision Lists)

- if (age = 18-20) then Recidivism = yes
- else if (male and age = 21-25) then Recidivism = yes
- else if (age = 26-30 and priors = 2-3) then Recidivism = yes
- else if (priors > 3) then Recidivism = yes
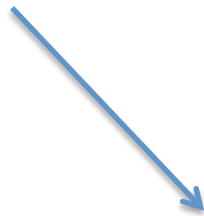- else (no)

# Rule List Models (Decision Lists)

- if (age = 18-20) then Recidivism = yes
- else if (male and age = 21-25) then Recidivism = yes
- else if (age = 26-30 and priors = 2-3) then Recidivism = yes
- else if (priors > 3) then Recidivism = yes
- else (no)

- Interpretable, logical
- Computationally hard to compute from data

# A new method for rule list learning

- With Elaine Angelino, Daniel Alabi, Nicholas Larus-Stone, Margo Seltzer

- Minimizes: errors + C* #rules

- Uses custom branch-and-bound.
  - Mines high-frequency itemsets, assembles rule list

o if (age = 18-20) then Recidivism = yes
o else if (male and age = 21-25) then Recidivism = yes
o else if (age = 26-30 and priors = 2-3) then Recidivism = yes
o else if (priors > 3) then Recidivism = yes
o else (no)

# A new method for rule list learning

- With Elaine Angelino, Daniel Alabi, Nicholas Larus-Stone, Margo Seltzer

- Minimizes: errors + C* #rules

- Uses custom branch-and-bound.
  - Mines high-frequency itemsets, assembles rule list
  - Fast bit-vector calculations, careful data structures
  - Knowledge of symmetry for rule lists
  - Theorems: Prefixes of rule lists that are too inaccurate or provably non-interpretable are removed (along with descendants)
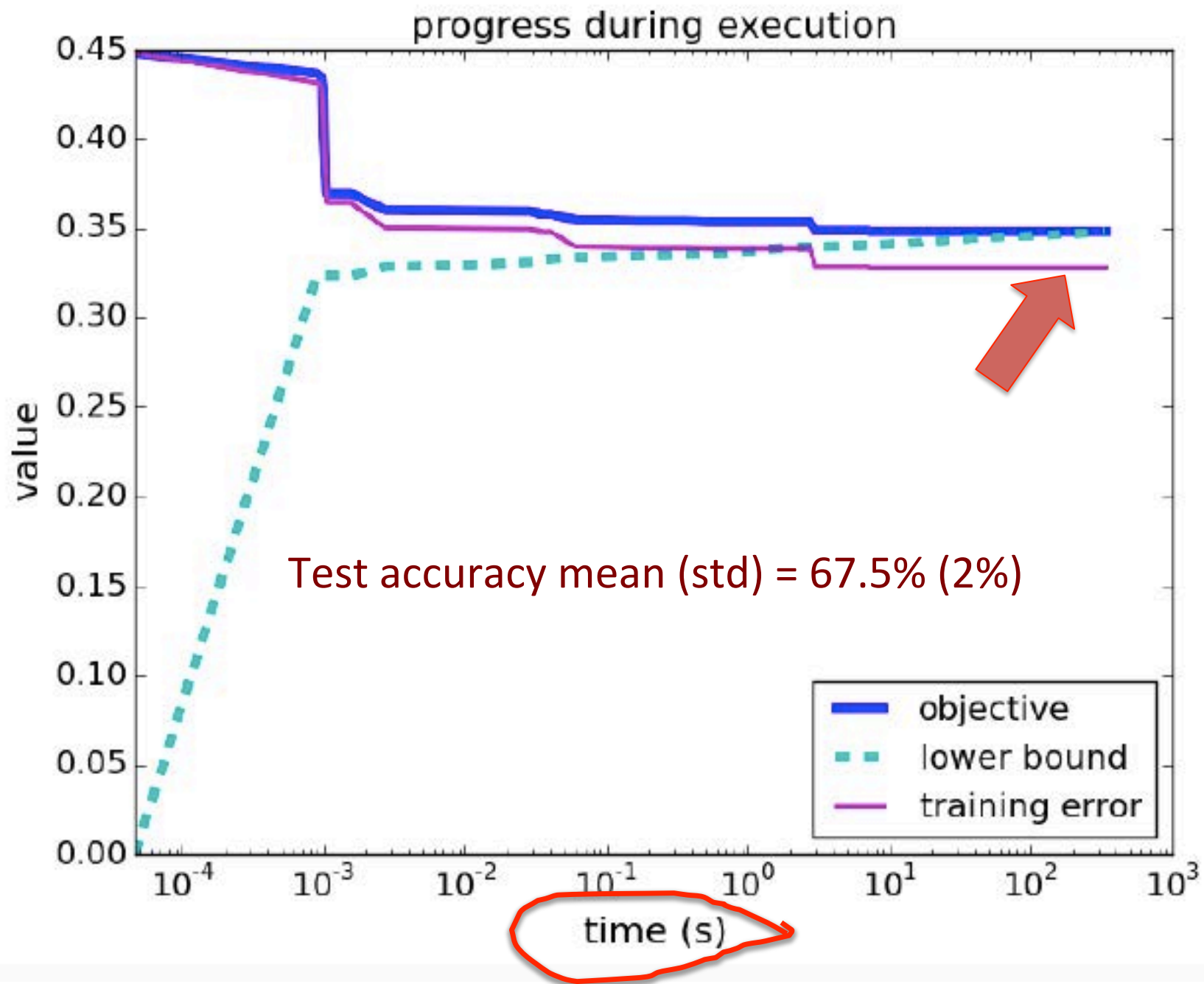  - Creates a certificate of optimality – provides best-in-class accuracy/interpretability tradeoff
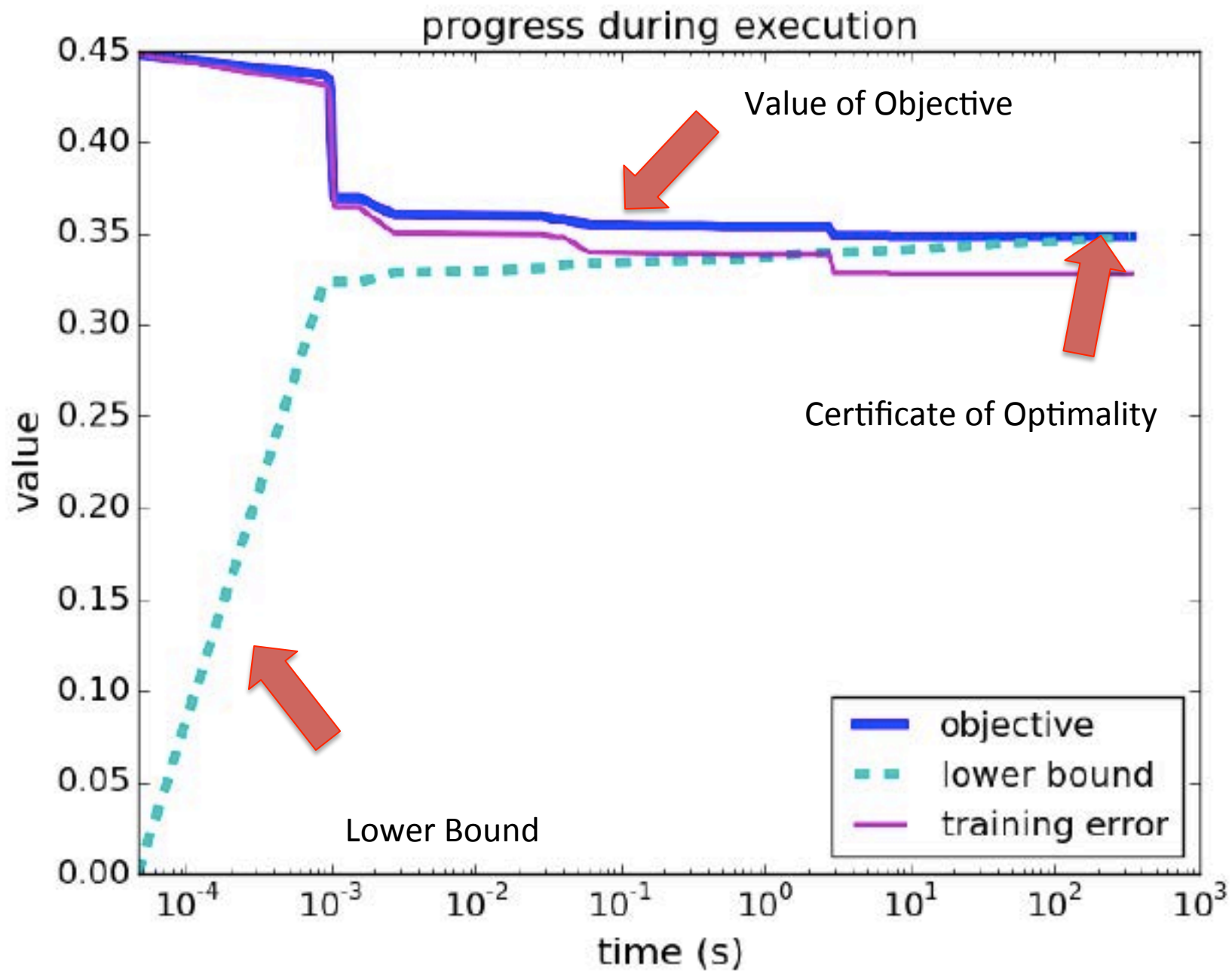
# Back to COMPAS score

- ProPublica calculated that on their recidivism dataset, COMPAS accuracy was <u>65.37%</u>.

| All Defendants | | |
|---|---|---|
| | Low | High |
| Survived | 2681 | 1282 |
| Recidivated | 1216 | 2035 |

- Does an interpretable model with that accuracy exist?

progress during execution

Test accuracy mean (std) = 67.5% (2%)

| | |
|---|---|
| **——** | objective |
| **----** | lower bound |
| —— | training error |

value

time (s)

# Rule List Models (Decision Lists)

- if (age = 18-20) then Recidivism = yes
- else if (male and age = 21-25) then Recidivism = yes
- else if (age = 26-30 and priors = 2-3) then Recidivism = yes
- else if (priors > 3) then Recidivism = yes
- else (no)

- if (male and juvenile crimes > 0) then Recidivism = yes
- else if (juvenile felonies = 0 and priors > 3) then Recidivism = yes
- else (no)

- Propublica article quotes COMPAS/ Northpointe founder Brennan:

- "Brennan said it is difficult to construct a score that doesn't include items that can be correlated with race — such as poverty, joblessness and social marginalization. "If those are omitted from your risk assessment, accuracy goes down," he said.

Hima Lakkaraju

# Learning Cost-Effective Treatment Regimes

- Model should be "causal": includes counterfactual inference
- Includes costs of gathering information (medical testing)
- Costs of treatment (cost of drug & side effects)
- Costs of outcome (making a wrong decision)
- Gives a prescription of how to test and treat each patient.

# Learning Cost-Effective Treatment Regimes

- If Gender=F, Current-Charge =Minor, Prev-Offense=None then Release on Personal Recognizance
- Else if Prev-Offense=Yes and Prior-Arrest =Yes then Release on Condition
- Else if Current-Charge =Misdemeanor and Age ≤ 30 then Release on Condition
- Else if Age ≥ 50 and Prior-Arrest=No, then Release on Personal Recognizance
- Else if Marital-Status=Single and Pays-Rent =No & Current-Charge =Misd. then Release on Condition
- Else if Addresses-Past-Yr ≥ 5 then Release on Condition
- Else Release on Personal Recognizance

# Berk Ustun's new ADHD scoring system

| | NEVER | RARELY | SOME-TIMES | OFTEN | VERY OFTEN |
|---|---|---|---|---|---|
| How often do you have trouble concentrating on what people say to you when they speak to you directly? | 0 | 4 | 4 | 5 | 5 |
| How often do you leave your seat in meetings or situations in which you are expected to remain seated? | 0 | 0 | 1 | 1 | 5 |
| How often do you have difficulty unwinding and relaxing when you have time to yourself? | 0 | 4 | 4 | 6 | 6 |
| How often do you finish the sentences of people you talk to, before they can finish them themselves? | 0 | 0 | 2 | 2 | 2 |
| How often do you put things off until the last minute? | 0 | 2 | 2 | 4 | 4 |
| How often do you depend on others to keep your life in order and attend to details? | 0 | 2 | 3 | 3 | 3 |

| TOTAL SCORE | 0 to 13 | 14 | 15 | 16 | 17 | 18 | 19 to 25 |
|---|---|---|---|---|---|---|---|
| PREDICTED RISK | <5.0% | 11.9% | 26.9% | 50.0% | 73.1% | 88.1% | >95.0% |

# Thanks