

Secrecy, Criminal Justice, Variable Importance, and Decision Trees

Cynthia Rudin

Professor of Computer Science, Electrical and Computer Engineering,
Statistical Science, and Mathematics

Duke University

How **important** is a variable?

Hold on... what does that mean?

How important is a variable

?

f_1 ← depends heavily on v

f_2 ← doesn't depend on v

Knowing how important a variable is to one model does not tell you how important it is in general.

How important is a variable for an algorithm?

Will I get the same accuracy with and without the variable?

Algorithm get model f_1 ← doesn't depend on v

Remove v ,

Algorithm get model f_2 ← doesn't depend on v

It turns out there exists f_3 ← depends heavily on v

Knowing how important a variable is to two models does not tell you how important it is in general.

This analysis would show we don't *need* v in order to perform well.

again:

How important is a variable to ?

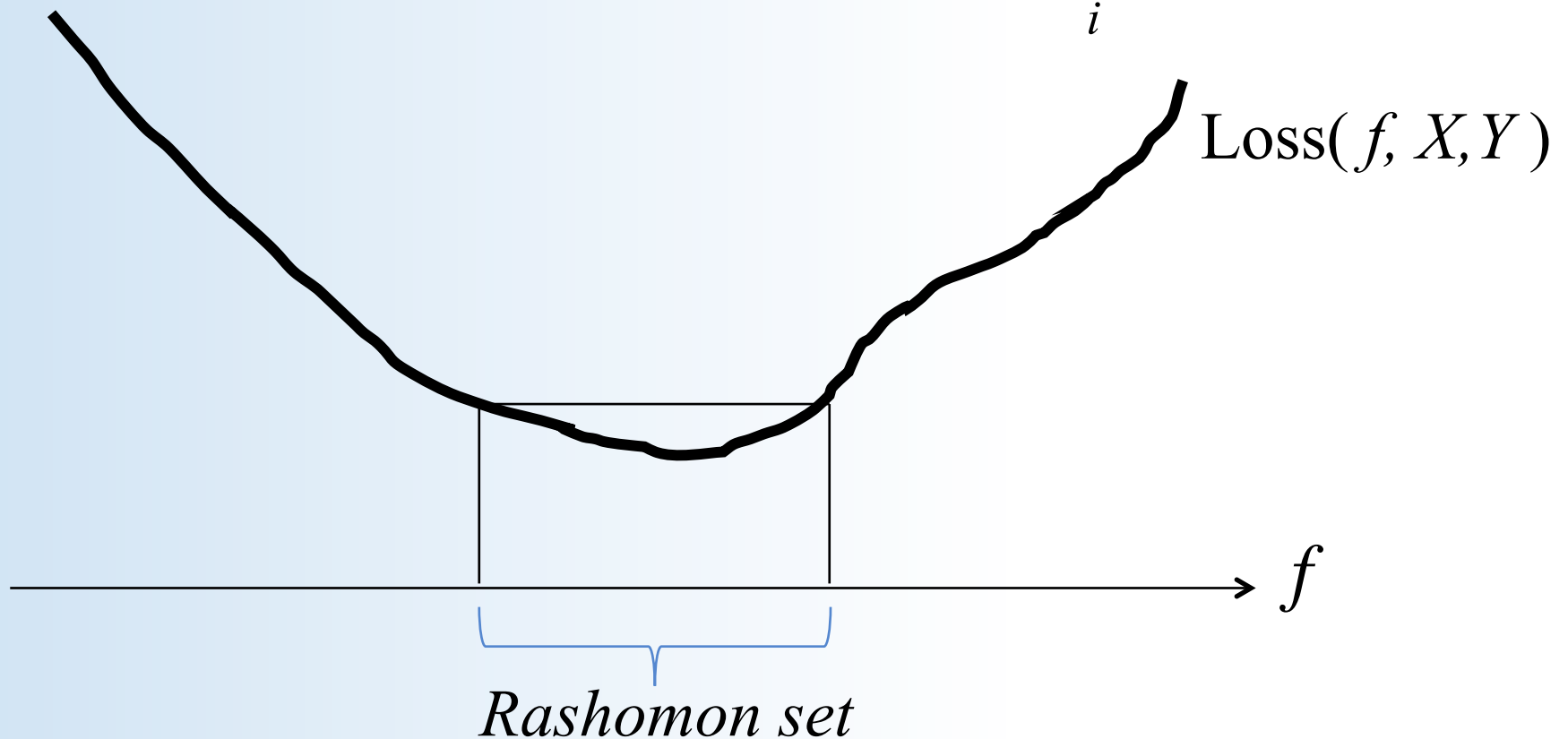
That's more like it!

In practice, we'll restrict to a flexible but restricted function class so we can compute and not overfit.

define the *Rashomon set* as the set of good models within F :

$$\{f : f \in \mathcal{F} \text{ such that } \text{Loss}(f, X, Y) \leq \epsilon\}$$

$$\text{perhaps } \text{Loss}(f, X, Y) = \sum_i (f(\mathbf{x}_i) - y_i)^2$$



Permuting a variable:

Reorder
r me!

$$X_{\text{scramble}} = \begin{pmatrix} x_{11} & x_{32} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{pmatrix}$$

Define *model reliance* of f on v :

$$\text{Model Reliance}(f, v) = \frac{\text{Loss}(f, X_{\text{scramble}}, Y)}{\text{Loss}(f, X, Y)}$$

If $\text{Model Reliance}(f, v) = 2$, then Loss doubles if we permute v

If $\text{Model Reliance}(f, v) = 1$, then f does not depend on v .

Define the *Rashomon set* as the set of good models within F :

$$\{f: f \in F \text{ such that } \text{Loss}(f, X, Y) \leq \epsilon \}.$$

Define *model reliance* of f on v :

$$\text{Model Reliance}(f, v) = \frac{\text{Loss}(f, X_{\text{scramble}}, Y)}{\text{Loss}(f, X, Y)}$$

How important is a variable to any good model?

\approx

What is the model reliance of functions in the
Rashomon set?

Secrecy, Criminal Justice, Variable Importance, and Decision Trees

Cynthia Rudin

Professor of Computer Science, Electrical and Computer Engineering,
Statistical Science, and Mathematics

Duke University

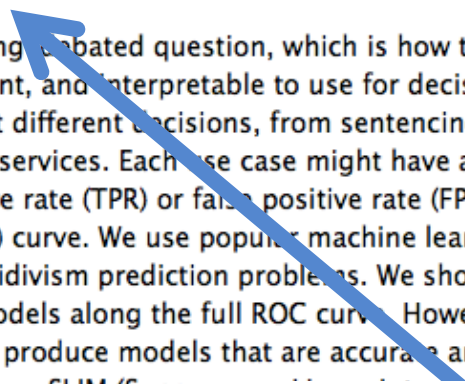
Interpretable Classification Models for Recidivism Prediction

Jiaming Zeng, Berk Ustun, Cynthia Rudin

(Submitted on 26 Mar 2015 (v1), last revised 8 Jul 2016 (this version, v6))

We investigate a long debated question, which is how to create predictive models of recidivism that are sufficiently accurate, transparent, and interpretable to use for decision-making. This question is complicated as these models are used to support different decisions, from sentencing, to determining release on probation, to allocating preventative social services. Each case might have an objective other than classification accuracy, such as a desired true positive rate (TPR) or false positive rate (FPR). Each (TPR, FPR) pair is a point on the receiver operator characteristic (ROC) curve. We use popular machine learning methods to create models along the full ROC curve on a wide range of recidivism prediction problems. We show that many methods (SVM, Ridge Regression) produce equally accurate models along the full ROC curve. However, methods that designed for interpretability (CART, C5.0) cannot be tuned to produce models that are accurate and/or interpretable. To handle this shortcoming, we use a new method known as SLIM (Super

models along the full ROC curve. They are just as accurate as the most highly interpretable.



Comments: 45 pages, 17 figures
 Subjects: **Machine Learning (stat.ML)**
 Cite as: [arXiv:1503.07810](https://arxiv.org/abs/1503.07810) [stat.ML]
 (or [arXiv:1503.07810v6](https://arxiv.org/abs/1503.07810v6) [stat.ML])

Original Article

Interpretable classification models for recidivism prediction

Jiaming Zeng✉, Berk Ustun, Cynthia Rudin

Submission history

From: Jiaming Zeng [[view email](#)]
[\[v1\]](#) Thu, 26 Mar 2015 18:21:29 GMT (43k)
[\[v2\]](#) Fri, 27 Mar 2015 04:32:31 GMT (43k)
[\[v3\]](#) Fri, 13 Nov 2015 01:09:31 GMT (34k)
[\[v4\]](#) Fri, 6 May 2016 14:50:11 GMT (791k)
[\[v5\]](#) Fri, 10 Jun 2016 02:05:32 GMT (791k)
[\[v6\]](#) Fri, 8 Jul 2016 01:22:05 GMT (382k)

First published: 05 September 2016 | <https://doi.org/10.1111/rssa.12227> |

[Read the full text >](#)

Statistics > Machine Learning

Interpretable Classification Models for Recidivism Prediction

Jiaming Zeng, Berk Ustun, Cynthia Rudin

(Submitted on 26 Mar 2015 (v1), last revised 8 Jul 2016 (this version, v6))

We investigate a long-debated question, which is how to create predictive models of recidivism that are sufficiently accurate, transparent, and interpretable to use for decision-making. This question is complicated as these models are used to support different decisions, from sentencing, to determining release on probation, to allocating preventative social services. Each use case might have an objective other than classification accuracy, such as a desired true positive rate (TPR) or false positive rate (FPR). Each (TPR, FPR) pair is a point on the receiver operator characteristic (ROC) curve. We use popular machine learning methods to create models along the full ROC curve on a wide range of recidivism prediction problems. We show that many methods (SVM, Ridge Regression) produce equally accurate models along the full ROC curve. However, methods that designed for interpretability (CART, C5.0) cannot be tuned to produce models that are accurate and/or interpretable. To handle this shortcoming, we use a new method known as SLIM (Supersparse Linear Integer Models) to produce accurate, transparent, and interpretable models along the full ROC curve. They are just as accurate as the most accurate models, but they are also highly interpretable.

Comments: 45 pages, 17 figures

Subjects: **Machine Learning (stat.ML)**Cite as: **arXiv:1503.07810** [stat.ML](or **arXiv:1503.07810v6** [stat.ML])

Submission history

From: Jiaming Zeng [[view email](#)]**[v1]** Thu, 26 Mar 2015 18:21:29 GMT (431kb,D)**[v2]** Fri, 27 Mar 2015 04:32:31 GMT (431kb,D)**[v3]** Fri, 13 Nov 2015 01:09:31 GMT (348kb,D)**[v4]** Fri, 6 May 2016 14:50:11 GMT (791kb,D)**[v5]** Fri, 10 Jun 2016 02:05:32 GMT (791kb,D)**[v6]** Fri, 8 Jul 2016 01:22:05 GMT (382kb,D)

- Most ML methods have similar performance across problems, including interpretable modeling methods.
- Race is not useful for predicting recidivism, but correlated with criminal history.



Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

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

World

vs.

Zeng et al. (JRSS 2016)

Government/COMPAS: Black box is necessary.

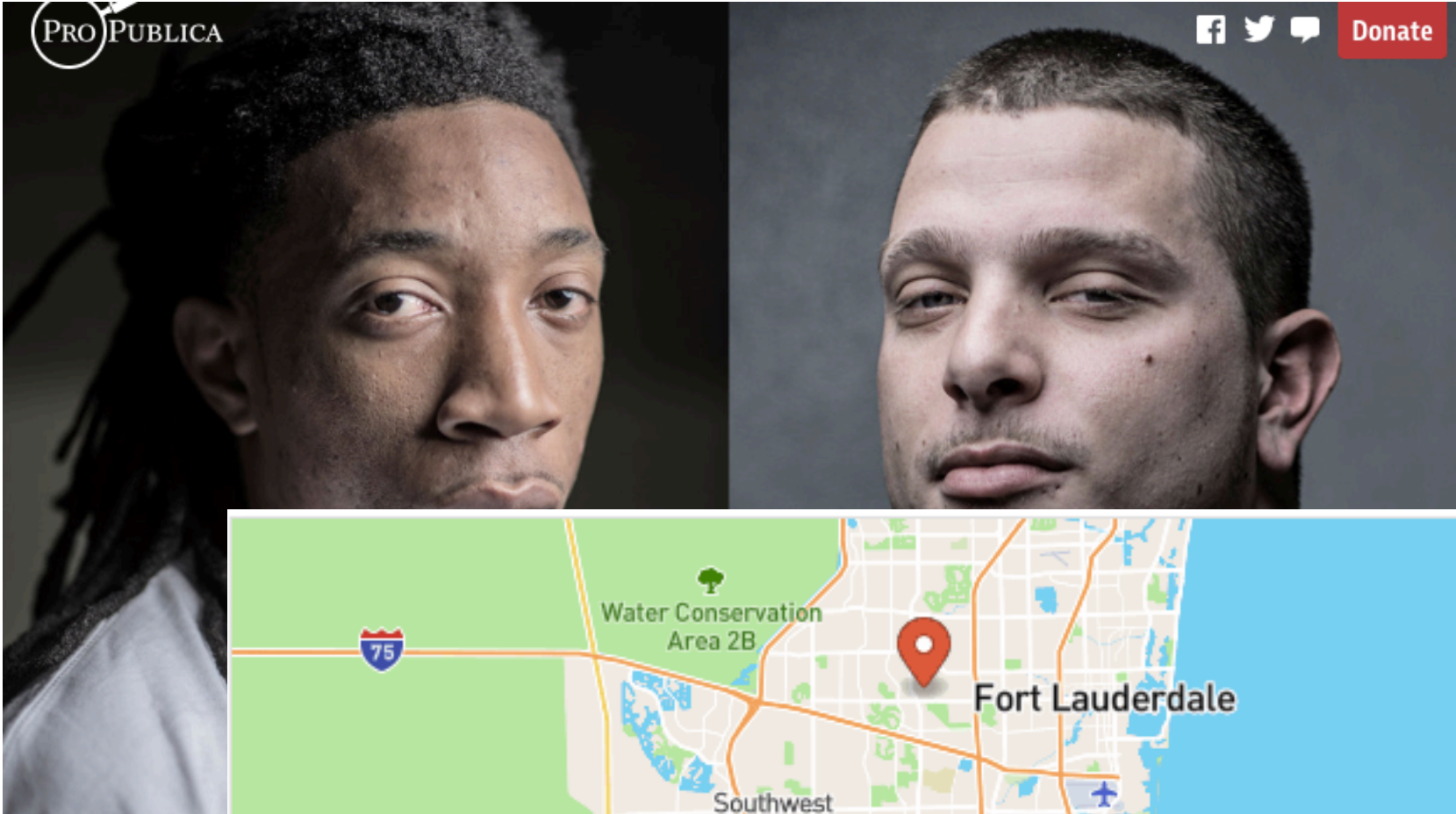
Interpretable models are just as good

Propublica: COMPAS depends on race (after conditioning on age and criminal history).

There's no need to use race after conditioning on age and criminal history, so any reasonable model wouldn't use it.

Has \$\$, has *data*

Has \$0



Broward County, Florida

broward.org



Broward County is a county located in the southeastern part of the U.S. state of Florida. [More at Wikipedia](#)

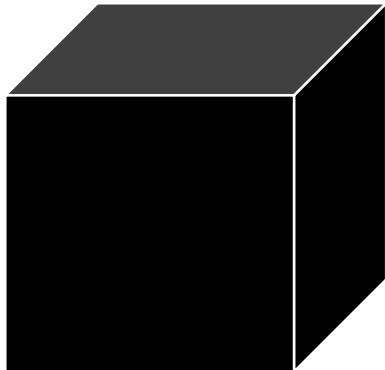
There



COMPAS vs. CORELS



COMPAS: (Correctional Offender
Management Profiling for Alternative
Sentences)

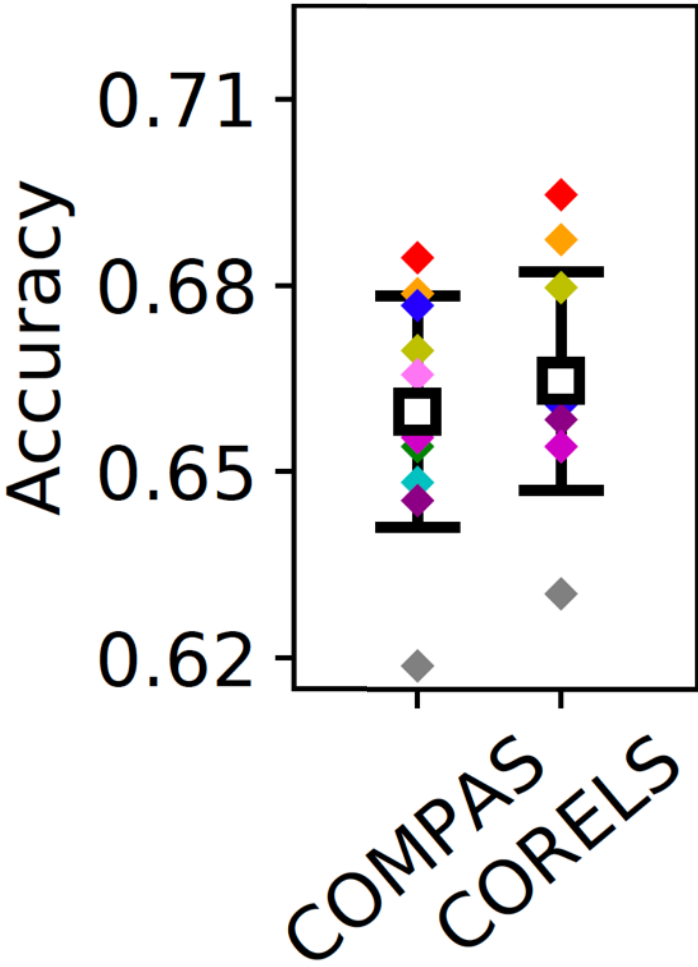


CORELS: (Certifiably Optimal Rule Lists,
with Elaine Angelino, Nicholas Larus-Stone,
Daniel Alabi, and Margo Seltzer, JMLR 2018)

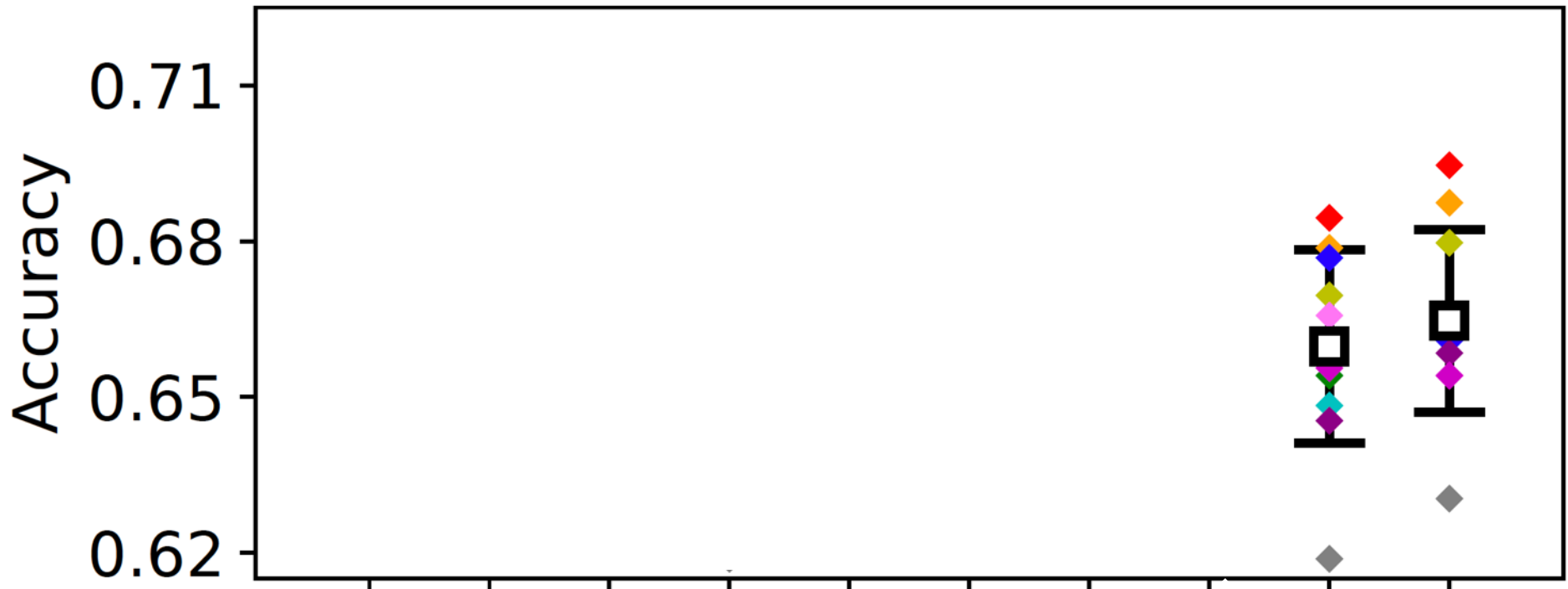
Here is the machine learning model:

```
If age=19-20 and sex=male, then predict arrest  
else if age=21-22 and priors=2-3 then predict arrest  
else if priors >3 then predict arrest  
else predict no arrest
```

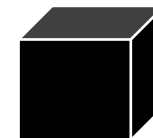

Prediction of arrest within 2 years



Prediction of arrest within 2 years



COMPAS
CORELS



If age=19-20 and sex=male, then p
else if age=21-22 and priors=2-3 th
else if priors >3 then predict arrest
else predict no arrest

World

vs. Zeng et al., JRSS, 2016

Government/COMPAS: Black
is necessary.

Interpretable models
are just as good

publica: COMPAS depends on
(after conditioning on age
criminal history).

There's no need to use
race, so any decent
model wouldn't use it.

propublica created a *linear* model to approximate COMPAS.
Coefficients for age, criminal history, *and race* were all positive.

Does that mean race is an important variable for COMPAS?

f_1 ← depends heavily on v

f_2 ← doesn't depend on v

No way! But what does COMPAS actually do?

COMPAS - Correctional Offender Management Profiling for Alternative Sentences. By Northpointe, Inc.

Conjecture: *The COMPAS general recidivism model is a nonlinear additive model. Its dependence on age in Broward County is approximately a linear spline, defined as follows:*

$$\text{for ages } \leq 33.26, f_{\text{age}}(\text{age}) = -0.056 \times \text{age} - 0.179$$

$$\text{for ages between } 33.26 \text{ and } 50.02, f_{\text{age}}(\text{age}) = -0.032 \times \text{age} - 0.963$$

$$\text{for ages } \geq 50.02, f_{\text{age}}(\text{age}) = -0.021 \times \text{age} - 1.541.$$

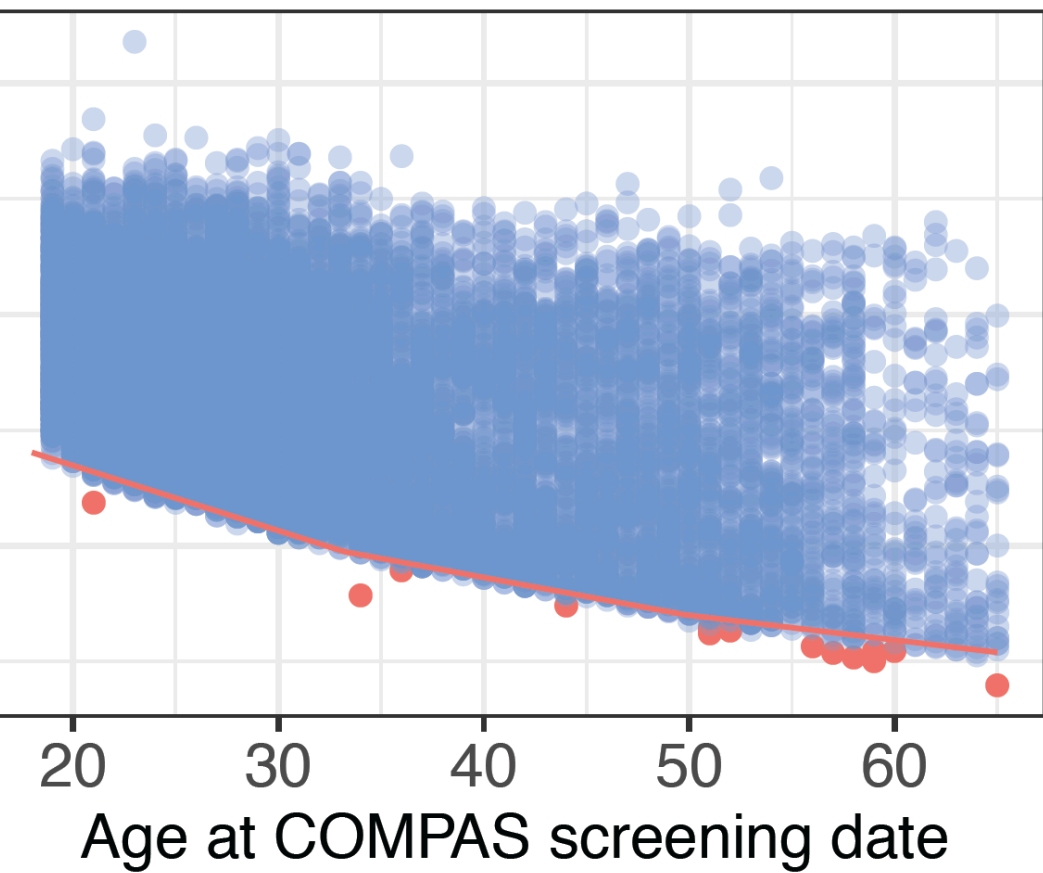
Similarly, the COMPAS violence recidivism model is a nonlinear additive model, with a dependence on age that is approximately a linear spline, defined by:

$$\text{for ages } \leq 21.77, f_{\text{viol age}}(\text{age}) = -0.205 \times \text{age} + 1.815$$

$$\text{for ages between } 21.77 \text{ and } 34.58, f_{\text{viol age}}(\text{age}) = -0.070 \times \text{age} - 1.113$$

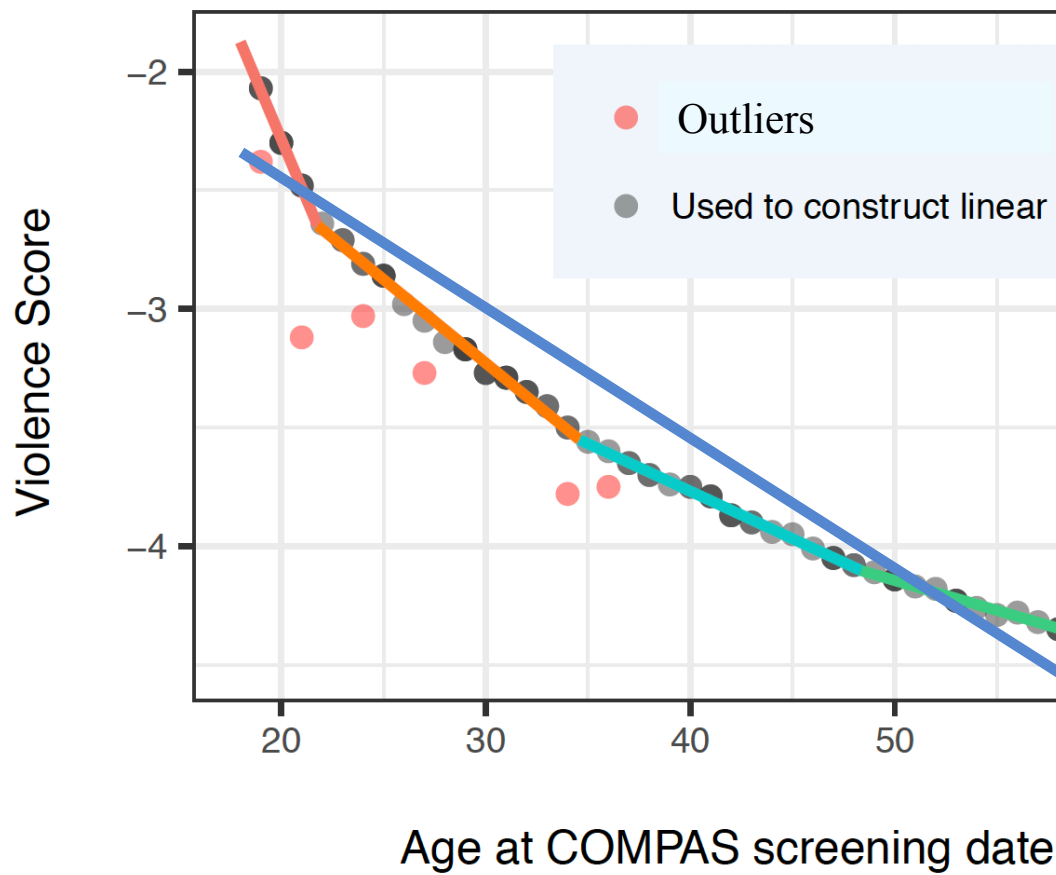
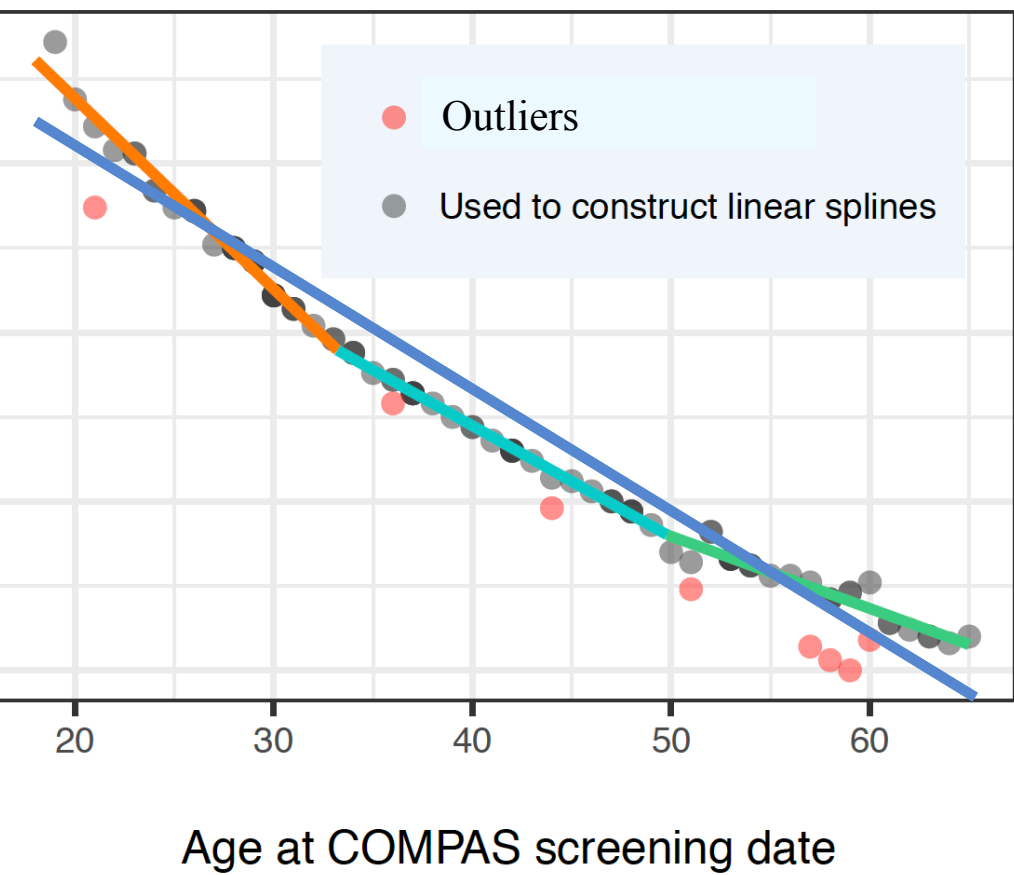
$$\text{for ages between } 34.58 \text{ and } 48.36, f_{\text{viol age}}(\text{age}) = -0.040 \times \text{age} - 2.166$$

$$\text{for ages } \geq 48.36, f_{\text{viol age}}(\text{age}) = -0.025 \times \text{age} - 2.882.$$



atter plot of COMPAS scores vs age for all individuals in Broward County F

Wang, and Coker. *The Age of Secrecy and Unfairness in Recidivism Prediction*. Harvard Data Science Review (



ProPublica's analysis isn't right. But COMPAS' manual seems wrong. So does COMPAS depend on race other than through age and criminal history?

Try 1:

Subtract off (what we think is) the contribution of age to COMPAS.

Then, run machine learning methods *with and without race* to see if they need race to predict COMPAS well.

	Linear Model	Random Forest	Boosting	SVM
Without Race	0.573	0.532	0.517	0.526
With Race	0.562	0.525	0.506	0.519

Table 4: RMSE of machine learning methods for predicting COMPAS general recidivism raw score after subtracting f_{age} with and without race as a feature. There is little difference with and without race. The differences between algorithms are due to differences in model forms. Age at COMPAS screening date and age at first offense are included as a features.

	Linear Model	Random Forest	Boosting	SVM
Without Race	0.498	0.493	0.468	0.475
With Race	0.489	0.482	0.456	0.474

Table 5: RMSE of machine learning methods for predicting COMPAS violence recidivism raw score after subtracting $f_{viol\ age}$ with and without race as a feature. Age at COMPAS screening date and age at first offense are included as a features.

ProPublica's analysis isn't right. But COMPAS' manual seems wrong. So does COMPAS depend on race other than through age and criminal history?

Try 1:

Subtract off (what we think is) the contribution of age to COMPAS.

Then, run machine learning methods *with and without race* to see if they explicitly need race to predict COMPAS well.

Knowing how important a variable is to two models does not tell you how important it is in general.

ProPublica's analysis isn't right. But COMPAS' manual seems wrong. So does COMPAS depend on race other than through age and criminal history?

Try 2:

Choose a flexible model class. Find the range of Model Reliance of functions in the Rashomon set.

Define the *Rashomon set* as the set of good models within F :

$$\{f: f \in F \text{ such that } \text{Loss}(f, X, Y) \leq \epsilon\}.$$

Define *model reliance* of f on v :

$$\text{Model Reliance}(f, v) = \frac{\text{Loss}(f, X_{\text{scramble}}, Y)}{\text{Loss}(f, X, Y)}$$

If $\text{Model Reliance}(f, v) = 1$, then f does not depend on v .

How important is a variable to any good model?
 \approx

What is the model reliance of functions in the
Rashomon set?

Define the *Rashomon set* as the set of good models within F :

$$\{f: f \in F \text{ such that } \text{Loss}(f, X, Y) \leq \epsilon \}.$$

Define *model reliance* of f on v :

$$\text{Model Reliance}(f, v) = \frac{\text{Loss}(f, X_{\text{scramble}}, Y)}{\text{Loss}(f, X, Y)}$$

Define *model class reliance* of F on v :

$$\text{Model Class Reliance}_+ (F, v) = \max_{f \in \text{Rashomon set}(F, \epsilon)} \text{Model Reliance}(f, v)$$

$$\text{Model Class Reliance}_- (F, v) = \min_{f \in \text{Rashomon set}(F, \epsilon)} \text{Model Reliance}(f, v)$$

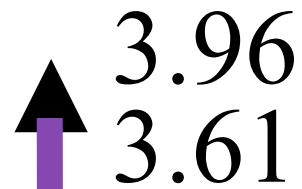
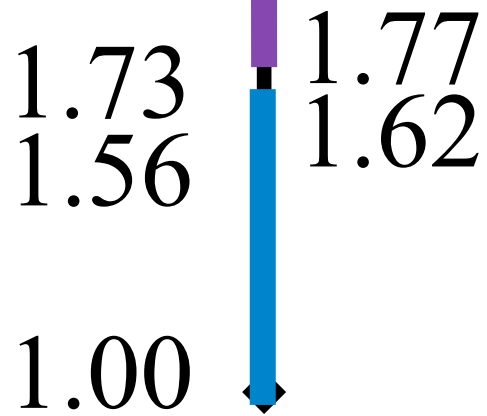
Choose a really flexible model class. Use age, criminal history, gender and race as regressors.

Choices:

- Kernel regression,
- Gaussian kernels with σ cross-validated on a training set,
- Regularized kernel weights (with parameter cross-validated)
- ε for the Rashomon Set as $0.1 \times$ minimum cross-validated loss.

Calculate **Model Class Reliance** on race and gender.

Model Class
Reliance on Race
and Gender



Model Class
Reliance on Age,
Criminal History

World

vs. Zeng et al., JRSS, 2016

Government/COMPAS:
Black box is necessary.

Interpretable models are just
as good

Propublica: COMPAS is
racially biased.

There's no need to use race
after conditioning on age and
criminal history, so any decent
model wouldn't use it.



OP-ED CONTRIBUTOR

When a Computer Program Keeps You in Jail

By Rebecca Wexler

June 13, 2017



to complicated proprietary
d yet.

h survey

least one typo

I NEEVR
MAKE TYPO

	COMPAS Violence Decile	# Priors	Selected Prior Charges	Selected Subsequent Charges
a pa	1	4	Aggravated Battery (F,1), Child Abuse (F,1), Resist Officer w/Violence (F,1)	
d r	1	14	Battery on Law Enforc Officer (F,3), Aggravated Assault W/Dead Weap (F,1), Aggravated Battery (F,1), Resist/obstruct Officer W/viol (F,1)	
y ers	1	15	Attempted Murder 1st Degree (F,1), Resist/obstruct Officer W/viol (F,1), Agg Battery Grt/Bod/Harm (F,1), Carrying Concealed Firearm (F,1)	Armed Sex Batt/vict 12 Yrs + (F,2), Aggravated Assault W/dead Weap (F,3) Kidnapping (F,1)
ando er	1	22	Aggrav Battery w/Deadly Weapon (F,1), Driving Under The Influence (M,2), Carrying Concealed Firearm (F,1)	
en er	1	28	Robbery / Deadly Weapon (F,11), Poss Firearm Commission Felony (F,7)	
s on	1	40	Resist/obstruct Officer W/viol (F,3), Battery on Law Enforc Officer (F,2), Attempted Robbery Deadly Weapo (F,1), Robbery 1 / Deadly Weapon (F,1)	
el alez	2	6	Murder in the First Degree (F,1), Aggrav Battery w/Deadly Weapon (F,1), Carrying Concealed Firearm (F,1)	

Name	COMPAS Violence Decile	# Priors	Selected Prior Charges	Selected Subsequent Charges
Vilma Dieppa	1	4	Aggravated Battery (F,1), Child Abuse (F,1), Resist Officer w/Violence (F,1)	
David Selzer	1	14	Battery on Law Enforc Officer (F,3), Aggravated Assault W/Dead Weap (F,1), Aggravated Battery (F,1), Resist/obstruct Officer W/viol (F,1)	
Berry Sanders	1	15	Attempted Murder 1st Degree (F,1), Resist/obstruct Officer W/viol (F,1), Agg Battery Grt/Bod/Harm (F,1), Carrying Concealed Firearm (F,1)	Armed Sex Batt/vict 12 Yrs + (F,2), Aggravated Assault W/dead Weap (F,3), Kidnapping (F,1)
Fernando Walker	1	22	Aggrav Battery w/Deadly Weapon (F,1), Driving Under The Influence (M,2), Carrying Concealed Firearm (F,1)	
Steven Glover	1	28	Robbery / Deadly Weapon (F,11), Poss Firearm Commission Felony (F,7)	
Rufus Jackson	1	40	Resist/obstruct Officer W/viol (F,3), Battery on Law Enforc Officer (F,2), Attempted Robbery Deadly Weapo (F,1), Robbery 1 / Deadly Weapon (F,1)	
Miguel Gonzalez	2	6	Murder in the First Degree (F,1), Aggrav Battery w/Deadly Weapon (F,1), Carrying Concealed Firearm (F,1)	
William Kelly	2	17	Aggravated Assault (F,5), Aggravated Assault W/dead Weap (F,2), Shoot/throw Into Vehicle (F,2), Battery Upon Detainee (F,1)	
Richard Campbell	2	21	Armed Trafficking In Cocaine (F,1), Poss Weapon Commission Felony (F,1), Carrying Concealed Firearm (F,1)	
John Coleman	2	25	Attempt Murder in the First Degree (F,1), Carrying Concealed Firearm (F,1), Felon in Pos of Firearm or Amm (F,1)	
Oscar Pope	2	38	Aggravated Battery (F,3), Robbery / Deadly Weapon (F,3), Kidnapping (F,1), Carrying Concealed Firearm (F,2)	Grand Theft in the 3rd Degree (F,3)
Travis Spencer	3	16	Aggravated Assault W/dead Weap (F,1), Burglary Damage Property >\$1000 (F,1), Burglary Unoccupied Dwelling (F,1)	
Michael Avila	3	17	Aggravated Assault W/dead Weap (F,2), Aggravated Assault w/Firearm (F,2), Discharge Firearm From Vehicle (F,1), Home Invasion Robbery (F,1)	Fail Register Vehicle (M,2)

Richard Campbell	2	21	Armed Trafficking In Cocaine (F,1), Poss Weapon Commission Felony (F,1), Carrying Concealed Firearm (F,1)	
John Coleman	2	25	Attempt Murder in the First Degree (F,1), Carrying Concealed Firearm (F,1), Felon in Pos of Firearm or Amm (F,1)	
Oscar Pope	2	38	Aggravated Battery (F,3), Robbery / Deadly Weapon (F,3), Kidnapping (F,1), Carrying Concealed Firearm (F,2)	Grand Theft in the 3rd Degree (F,3)
Travis Spencer	3	16	Aggravated Assault W/dead Weap (F,1), Burglary Damage Property >\$1000 (F,1), Burglary Unoccupied Dwelling (F,1)	
Michael Avila	3	17	Aggravated Assault W/dead Weap (F,2), Aggravated Assault w/Firearm (F,2), Discharge Firearm From Vehicle (F,1), Home Invasion Robbery (F,1)	Fail Register Vehicle (M,2)
Terrance Murphy	3	20	Solicit to Commit Armed Robbery (F,1), Armed False Imprisonment (F,1), Home Invasion Robbery (F,1)	Driving While License Revoked (F,3)
Anthony Hawthorne	3	25	Attempt Sexual Batt / Vict 12+ (F,1), Resist/obstruct Officer W/viol (F,1), Poss Firearm W/alter/remov Id# (F,1)	
Stephen Brown	3	36	Carrying Concealed Firearm (F,2), Battery On Law Enforce Officer (F,1), Kidnapping (F,1), Aggravated Battery (F,1)	Driving While License Revoked (F,3)
Samuel Walker	3	36	Murder in the First Degree (F,1), Poss Firearm Commission Felony (F,1), Solicit to Commit Armed Robbery (F,1)	Petit Theft 100–300 (M,1)
Jesse Bernstein	4	10	Aggravated Battery / Pregnant (F,1), Sex Battery Vict Mental Defect (F,1), Shoot/throw In Occupied Dwell (F,1)	Tresspass in Struct/Convey Occupy (M,1)
Shandedra Hardy	4	16	Aggrav Battery w/Deadly Weapon (F,1), Felon in Pos of Firearm or Amm (F,4)	Resist/Obstruct W/O Violence (M,1), Possess Drug Paraphernalia (M,1)

Current State of Affairs

Human judges are biased black b
COMPAS is a black box, possibly
reliable. Financial incentives to u
California is moving towards a no
sing COMPAS.

ProPublica's seriously flawed analysis can highly cited and
respected

Academic interest in fairness is huge, interest in explainability
of black boxes is huge...

Little interest/expertise in interpretability

The Mercury News

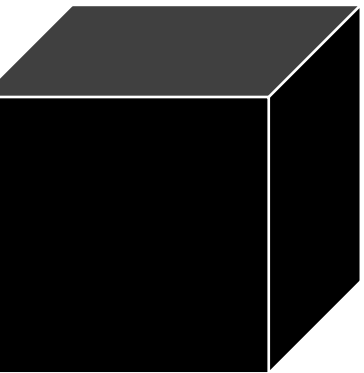
News > California News • News

Law ending cash bail in California halted after referendum qualifies for 2020 ballot

SB 10 won't start in October as planned; final version divided legislators and civil rights groups who initially supported it



COMPAS:



CORELS Model:

If age=19-20 and sex=male, then predict arrest
else if age=21-22 and priors=2-3 then predict ar
else if priors >3 then predict arrest
else predict no arrest

Behind the scenes

Model Class Reliance

Optimal Decision Trees

Define the *Rashomon set* as the set of good models within F :

$$\{f: f \in F \text{ such that } \text{Loss}(f, X, Y) \leq \epsilon\}.$$

Define *model reliance* of f on v :

$$\text{Model Reliance}(f, v) = \frac{\text{Loss}(f, X_{\text{scramble}}, Y)}{\text{Loss}(f, X, Y)}$$

Define *model class reliance* of F on v :

$$\text{Model Class Reliance}_+(F, v) = \max_{f \in \text{Rashomon set}(F, \epsilon)} \text{Model Reliance}(f, v)$$

$$\text{Model Class Reliance}_-(F, v) = \min_{f \in \text{Rashomon set}(F, \epsilon)} \text{Model Reliance}(f, v)$$

er et al. (in progress, 2019) contains:

Estimation using U-statistics

Learning theoretic bounds

How to compute MCR efficiently (linear, additive, reproducing kernel Hilbert space)

Connections to causal inference

Define *model class reliance* of F on v :

$$\text{Model Class Reliance}_+ (F, v) = \max_{f \in \text{Rashomon set}(F, \epsilon)} \text{Model Reliance}(f, v)$$

$$\text{Model Class Reliance}_- (F, v) = \min_{f \in \text{Rashomon set}(F, \epsilon)} \text{Model Reliance}(f, v)$$

Behind the scenes

Model Class Reliance

Optimal Decision Trees

CORELS: Certifiably Optimal Rule Lists (JMLR, 2018)

Predictive model for 2 yr recidivism, built from data from Broward County Florida

```
IF age between 18-20 and sex is male THEN predict arrest within 2 years
ELSE IF age between 21-23 and 2-3 prior offenses THEN predict arrest
ELSE IF more than three priors THEN predict arrest
ELSE predict no arrest.
```

in <10 seconds, certified to optimality in ~2 minutes, over the space of rule lists.

Server with two Intel Xeon E5-2699 v4 (55 MB cache, 2.20 GHz) processors and 448 GB RAM

$$R(d, \mathbf{x}, \mathbf{y}) = \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{1}_{[\text{point } i \text{ is misclassified}]}}_{\text{Misclassification error}} + \lambda \underbrace{\# \text{Rules in list}}_{\text{Sparsity}}$$

Usually 0.01

IF age between 18-20 and sex is male

ELSE IF age between 21-23 and 2-3 prior offenses THEN predict arrest

ELSE IF more than three priors

ELSE

THEN predict arrest within 2 years

THEN predict arrest

THEN predict arrest

predict no arrest.

Several theorems led to bounds

Theorem 1: If a rule's **support** is less than λ , that rule cannot be in an optimal rule list.

Theorem 2: If a rule in the list does not **correctly classify** at least λ fraction of observations, that rule cannot be in an optimal rule list.

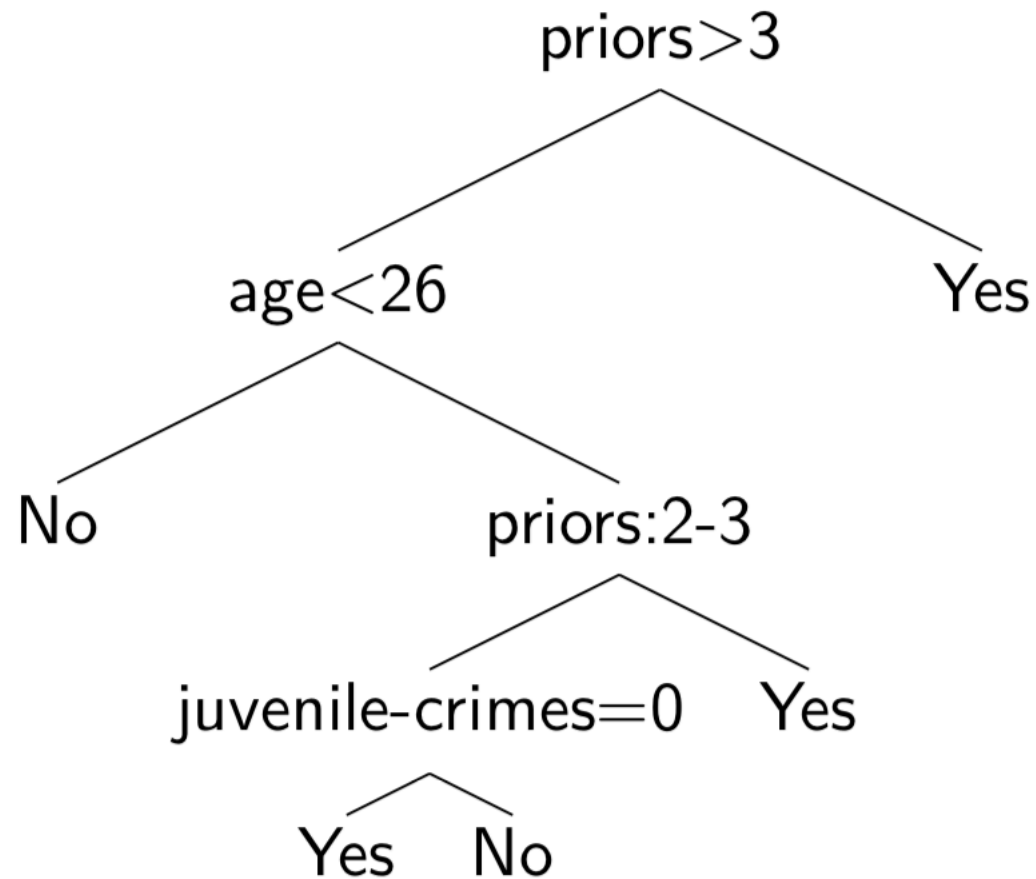
Theorem 3: The **length** of an optimal rule list is bounded by a function of λ , the accuracy of the current best model so far, and the accuracy and length of our current prefix (partial rule list).

Theorem 4: **One-step-lookahead bound**: If a prefix's lower bound is within λ of the best current objective, adding any rules to it will lead to a non-optimal rule list.

Theorem 5: **Equivalent Points Bound**: For every set of “equivalent” points, we will classify at least the minority label of them wrong.

Theorem 6: **Permutation Bound**: Only an optimal permutation of a set of rules can be extended to form an optimal rule list.

currently...



Rudin, Seltzer. Optimal Sparse Decision Trees, NeurIPS (spotlight) 2019]

Behind the scenes

Model Class Reliance: <https://github.com/aaronjfisher/mcr>

, Rudin, Dominici. All Models are Wrong but many are Useful: Learning a Variable's Importance by Studying Class of Prediction Models Simultaneously. <https://arxiv.org/abs/1801.01489> , In Progress, 2019

Optimal Decision Trees: <https://corels.eecs.harvard.edu>

no et al. Certifiably Optimal Rule Lists for Categorical Data, Journal of Machine Learning Research, 2018.

adin, Seltzer. Optimal Sparse Decision Trees. <https://arxiv.org/abs/1904.12847> , NeurIPS, 2019



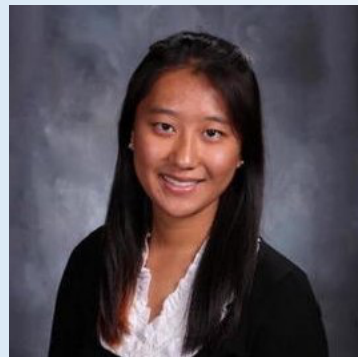
Elaine Wang



Paul Coker

of Secrecy and
s in Recidivism
n, Harvard Data
review, 2019

Interpretable Classification Models
For Recidivism Prediction, Journal of
the Royal Statistical Society, 2016



Jiaming Zeng



Berk Ustun



Aaron Fisher



Daniel Alabi



Elaine Angelino



Nicholas Larun

Certifiably Optimal Rule Lists for
Categorical Data, Journal of
Machine Learning Research, 2018.



Francesca Dominici



Margo

All Models are Wrong but many are Useful: Variable Importance for Black-Box,
Proprietary, or Misspecified Prediction Models, using Model Class Reliance, 2018

Systems Techniques

Custom bit-vector library for rule list evaluation

Computational reuse for evaluating multiple lists with similar prefixes

Priority queue

Data structures: trie (prefix tree), symmetry-aware map, and queue

Mine all rules with sufficient support.

Start with rule lists of size 1.

While queue of rule lists is not empty:

Take current prefix from queue, consider each of its children and check:

- length bound
- rule accuracy
- one step-lookahead bound
- equivalent points bound
- symmetry-aware pruning

If lower bound is higher than current best, prefix is no good.

Otherwise add it into queue.

If current objective is lower than current best, update and store rule list.

End while

Output is optimal rule list (with certificate of optimality)

anks

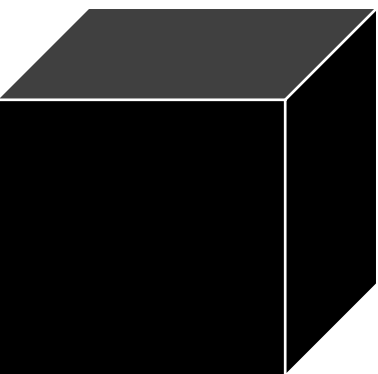
ia Rudin, Caroline Wang, and Beau Coker. *The Age of Secrecy and Unfairness in Recidivism Prediction*. 2018

, Rudin, Dominici. All Models are Wrong but many are Useful: Variable Importance for Black-Box, Proprietary Specified Prediction Models, using Model Class Reliance. 2018

ng Zeng, Berk Ustun, Cynthia Rudin. Interpretable Classification Models For Recidivism Prediction. *Journal of the Royal Statistical Society*, 2016.

Angelino, Nicholas Larus-Stone, Daniel Alabi, Margo Seltzer, and Cynthia Rudin, Certifiably Optimal Rule Learning for Categorical Data, *Journal of Machine Learning Research*, 2018.

COMPAS:



CORELS Model:

If age=19-20 and sex=male, then predict arrest
else if age=21-22 and priors=2-3 then predict arrest
else if priors >3 then predict arrest
else predict no arrest