

RACIAL BIAS

MACHINE LEARNING IN THE CRIMINAL JUSTICE SYSTEM

JOSEPH GAMBINO
OCTOBER 25, 2017

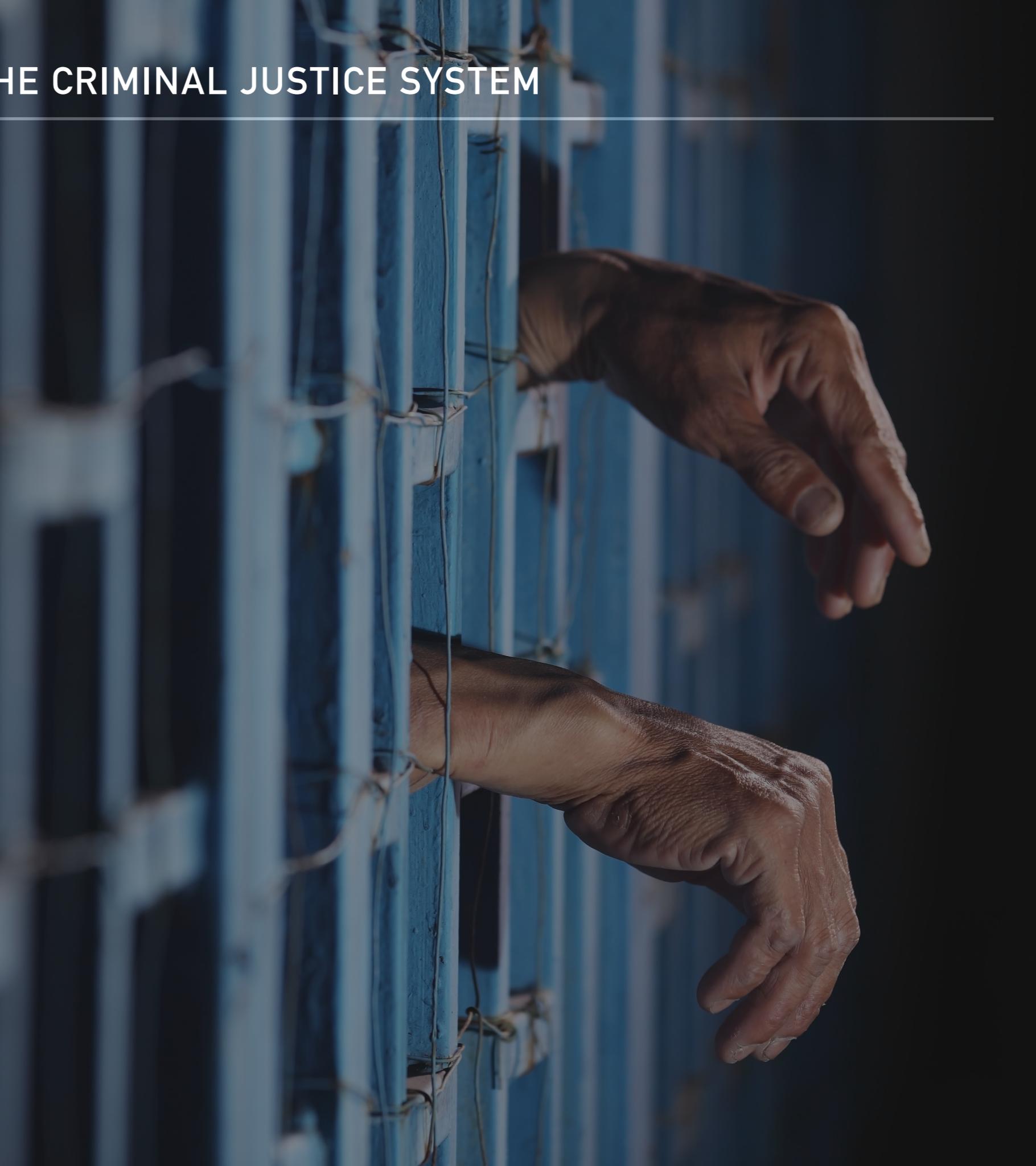
MACHINE LEARNING IN THE CRIMINAL JUSTICE SYSTEM

RISK OF RECIDIVISM



RISK OF RECIDIVISM

- ▶ Risk Assessment



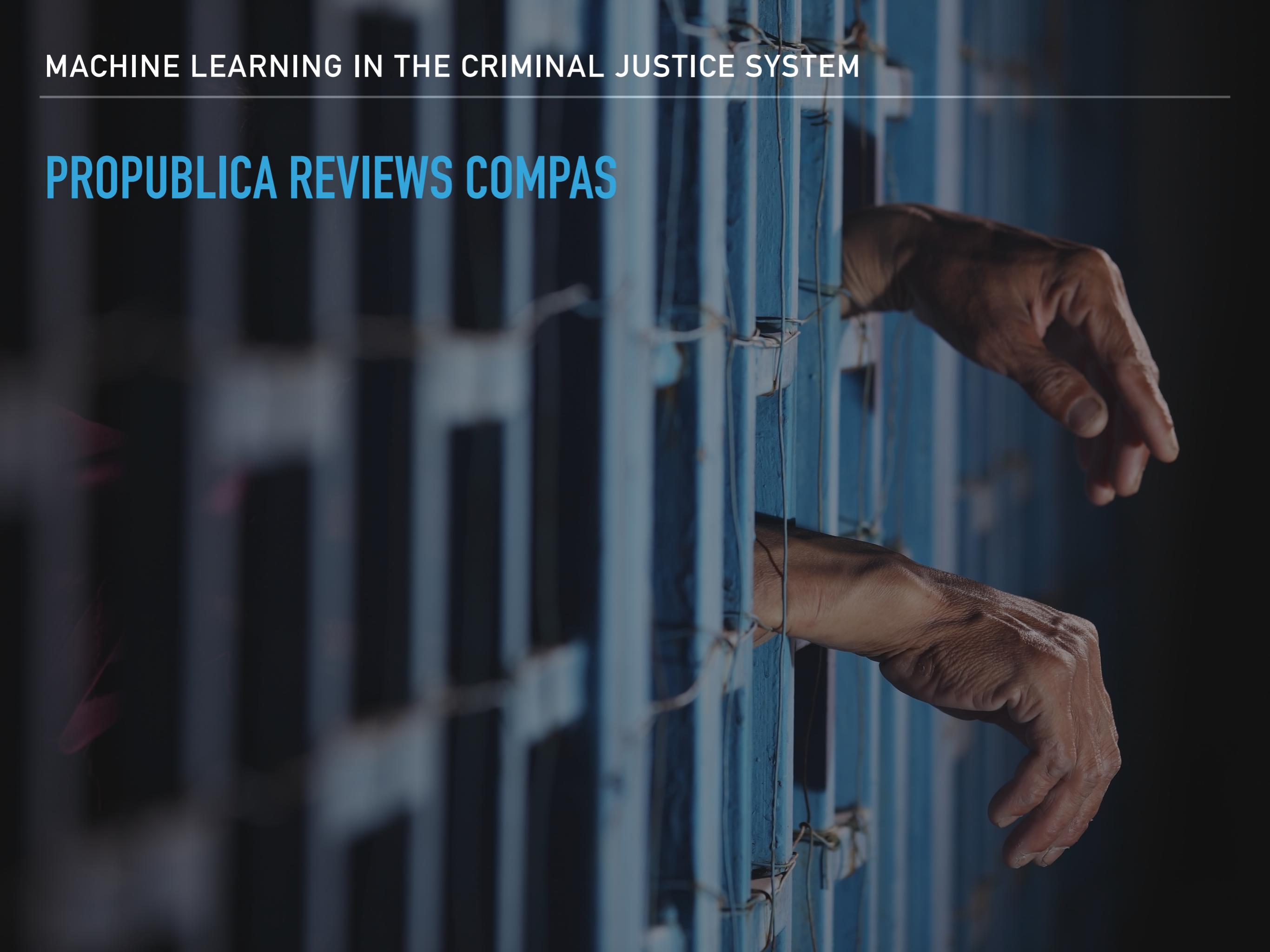
RISK OF RECIDIVISM

- ▶ Risk Assessment
- ▶ Used at all stages
 - ▶ Bail/Bond Hearing
 - ▶ Sentencing
 - ▶ Release on Parole



MACHINE LEARNING IN THE CRIMINAL JUSTICE SYSTEM

PROPUBLICA REVIEWS COMPAS



PROPUBLICA REVIEWS COMPAS

- ▶ Recidivism - 60% accurate
- ▶ Violent recidivism - 20% accurate

PROPUBLICA REVIEWS COMPAS

- ▶ Recidivism - 60% accurate
- ▶ Violent recidivism - 20% accurate
- ▶ Black defendants were twice as likely to be mis-labeled as high risk
- ▶ White defendants were twice as likely to be mis-labeled as low risk

OPTIMIZING A MODEL FOR PRECISION

DATA

- ▶ U.S. Dept. of Justice: Recidivism of Felons on Probation, 1986-1989
- ▶ Older data, but thorough
- ▶ 12,369 probationers across 32 jurisdictions



OPTIMIZING A MODEL FOR PRECISION

CONFUSION MATRIX

		Model Prediction	
		Yes	No
Actual Behavior	Yes	True Positive	False Negative
	No	False Positive	True Negative



OPTIMIZING A MODEL FOR PRECISION

PRECISION

		Model Prediction	
		Yes	No
Actual Behavior	Yes	True Positive	False Negative
	No	False Positive	True Negative



OPTIMIZING A MODEL FOR PRECISION

NEED AN INTERPRETABLE MODEL



OPTIMIZING A MODEL FOR PRECISION

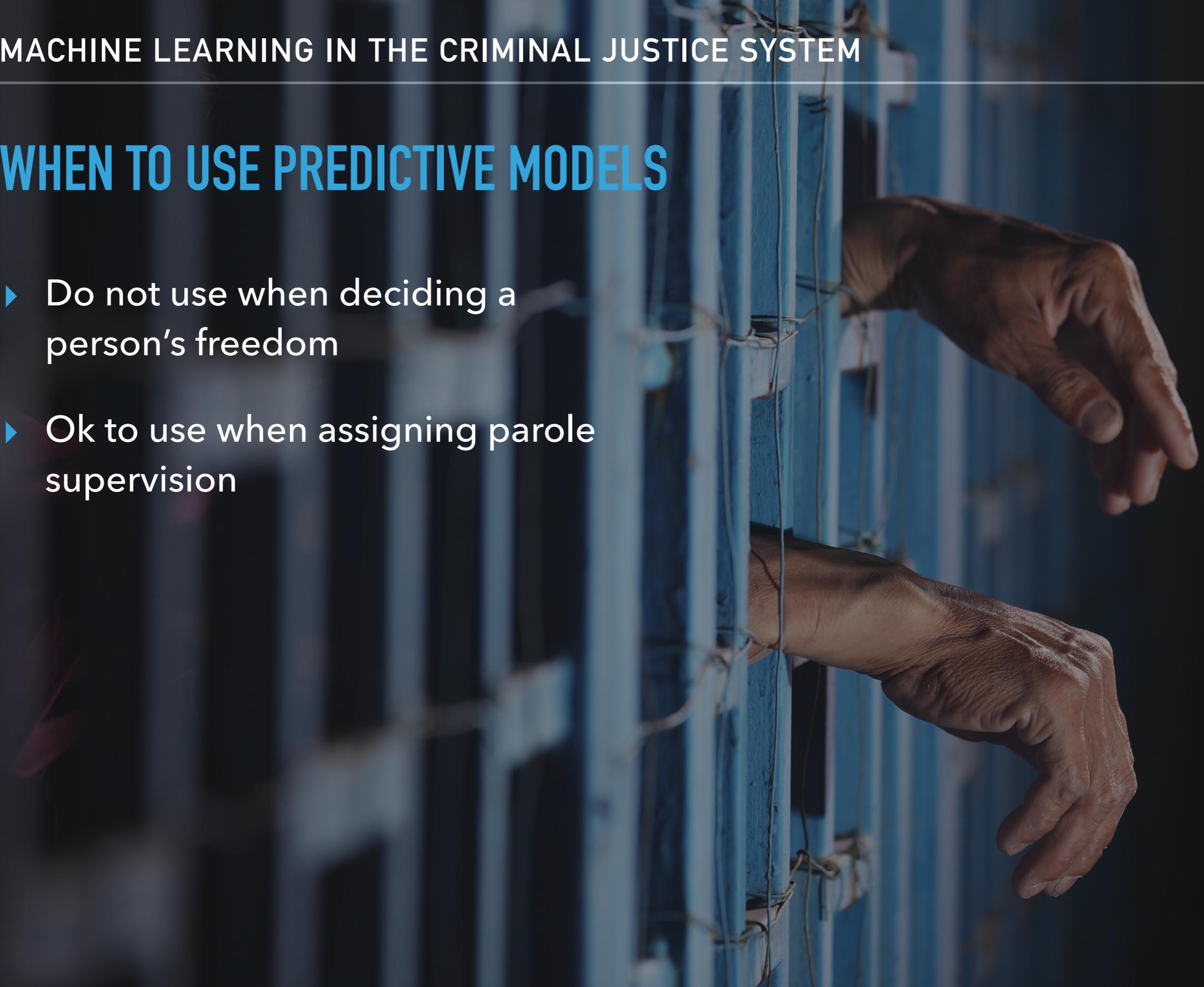
NEED AN INTERPRETABLE MODEL

- ▶ Logistic Regression
 - ▶ Can look at weights for each feature
- ▶ Did not use Random Forest
- ▶ Flask [App](#)



WHEN TO USE PREDICTIVE MODELS

- ▶ Do not use when deciding a person's freedom
- ▶ Ok to use when assigning parole supervision



APPENDIX



OPTIMIZING A MODEL FOR PRECISION

ACCURACY

		Model Prediction	
		Yes	No
Actual Behavior	Yes	True Positive	False Negative
	No	False Positive	True Negative



OPTIMIZING A MODEL FOR PRECISION

FALSE POSITIVE RATE

		Model Prediction	
		Yes	No
Actual Behavior	Yes	True Positive	False Negative
	No	False Positive	True Negative



OPTIMIZING A MODEL FOR PRECISION

FALSE NEGATIVE RATE

		Model Prediction	
		Yes	No
Actual Behavior	Yes	True Positive	False Negative
	No	False Positive	True Negative

