

MACHINE_LEARNING_FOR_ SOCIAL_SCIENCE from scratch

Ana Valdivia

Data Scientist (Trilateral)
PhD in Computer Science (UGR)
Mathematician (UPC)

@ana_valdi

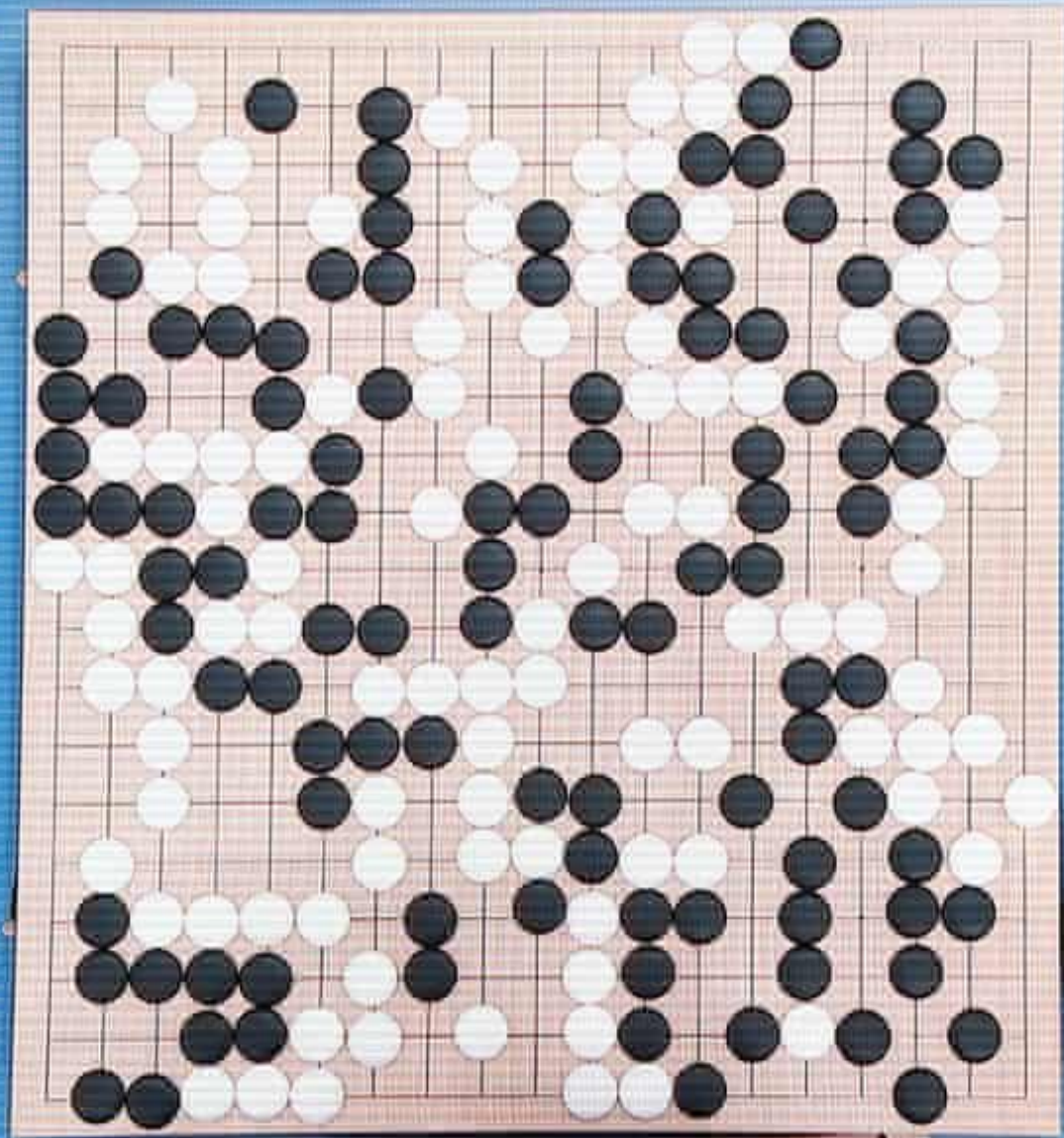


What is Artificial Intelligence?

What is Artificial Intelligence?

Artificial
Intelligence (AI)
refers to software
technologies that
make a robot or
computer act and
think like a human.

McCarthy, 1995.



柯洁 KE JIE

00:17:05



ALPHAGO

01:51:38

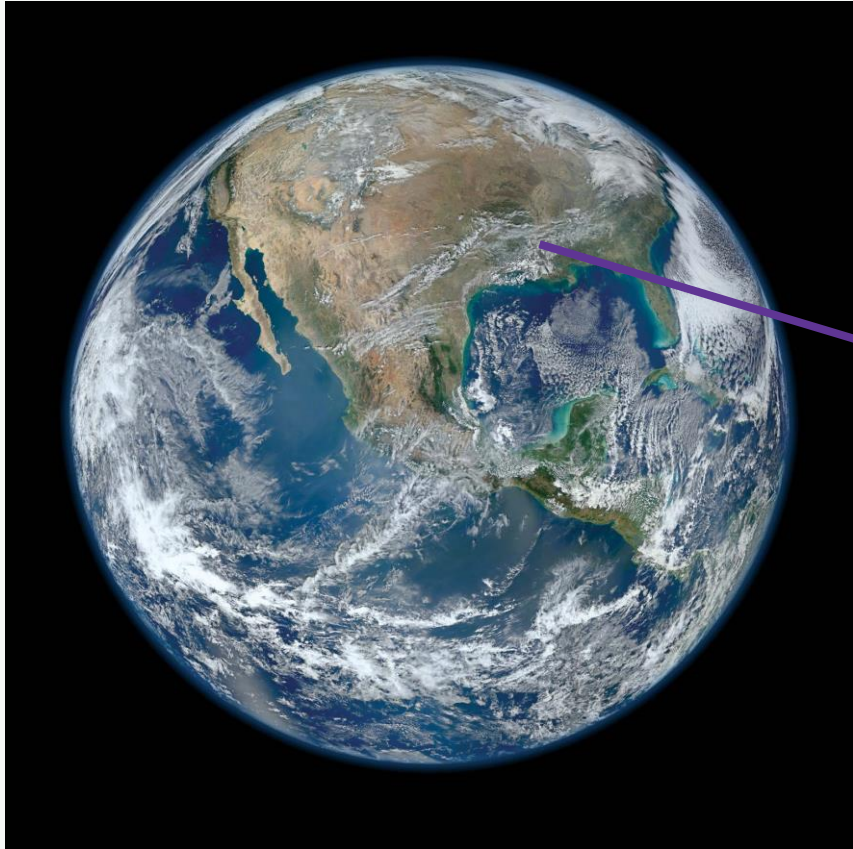
浙江省体育局



What is Machine Learning?

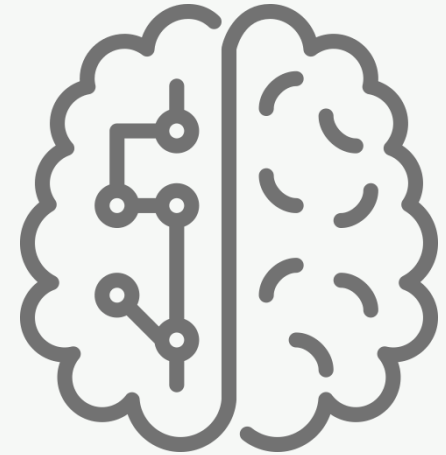
What is Machine Learning?

Machine Learning (ML) is the scientific study of algorithms and statistical models that computer systems use to perform a specific task without using explicit instructions, relying on patterns and inference instead.



Vida real

Dataset



Machine Learning

Real Application

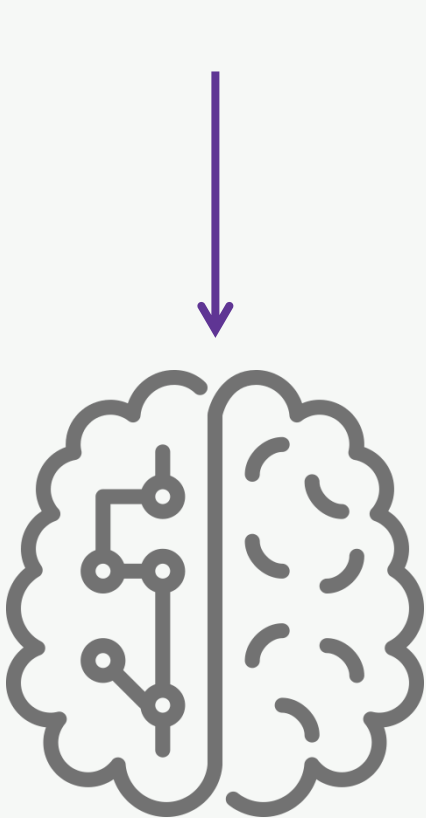


Real Application



Name	Gender	Age	Race	Juv_misd_count	True_outcome
Miguel	Male	34	Hispanic	Aggravated Assault	Positive
Benjamin	Male	47	Caucasian	Grand Theft in the 3rd Degree	Negative

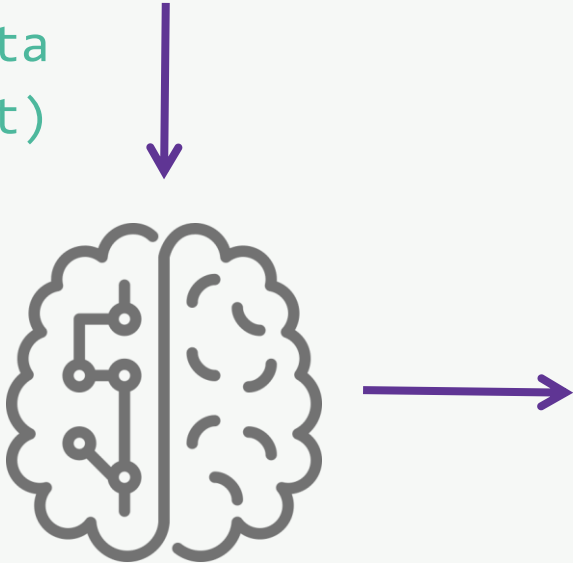
Historical data
(training data set)



In two
years
he/she
reincides?

Name	Gender	Age	Race	Juv_misd_count
Darrious	Male	27	African-American	Kidnapping
Claire	Female	23	Caucasian	Possess Cannabis

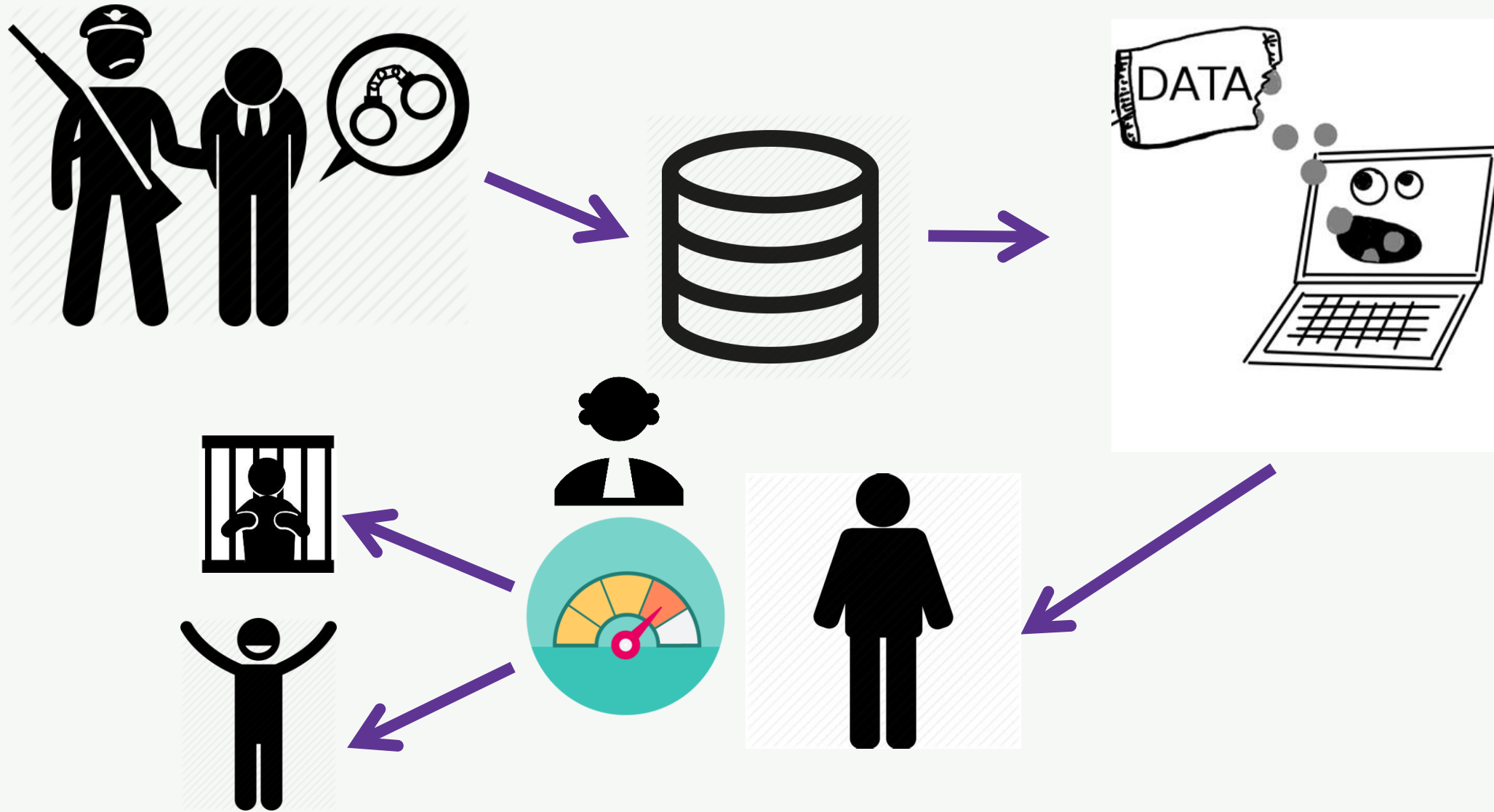
New data (testing data set)



Prediction
Positive
Positive



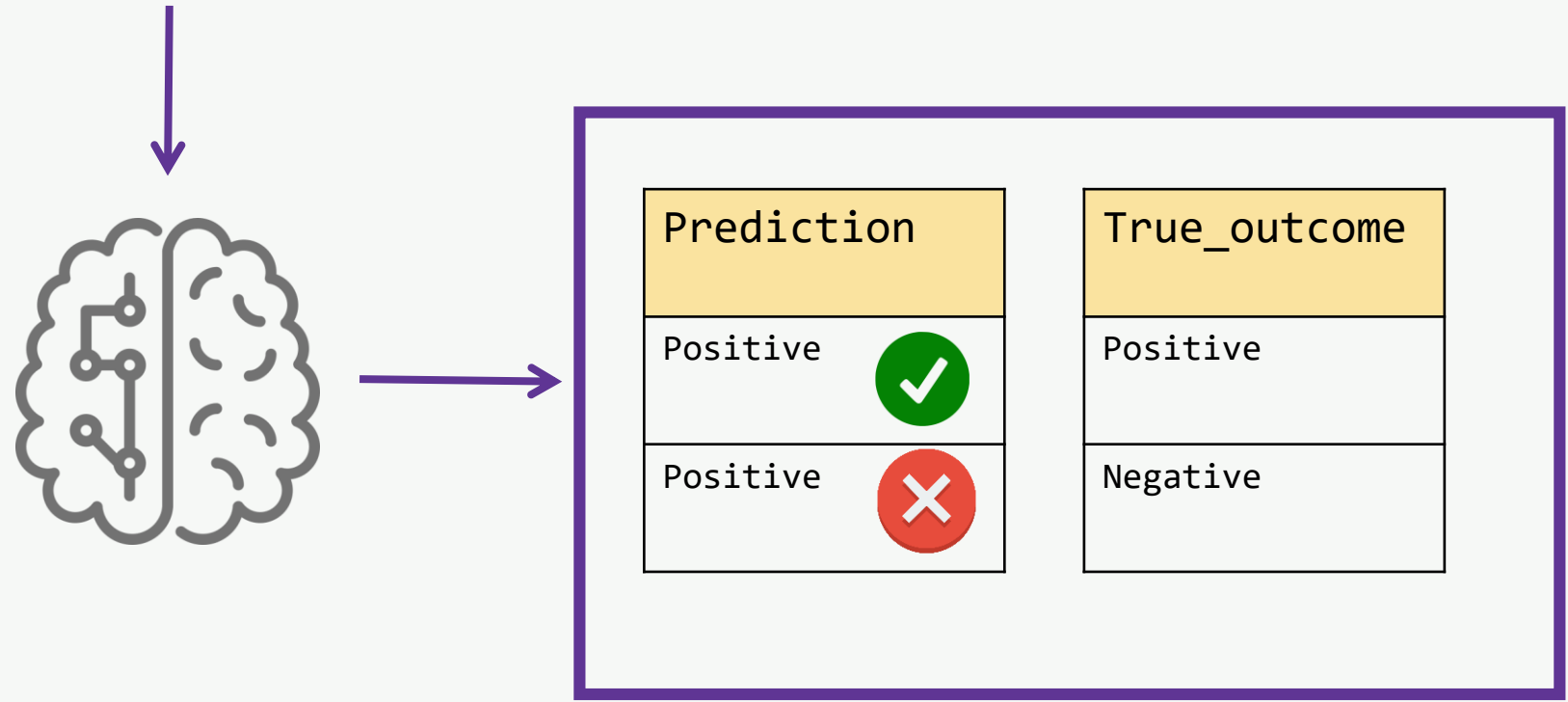
Real Application



Name	Gender	Age	Race	Juv_misd_count
Miguel	Male	34	Hispanic	Aggravated Assault
Benjamin	Male	47	Caucasian	Grand Theft in the 3rd Degree

Super important to understand!

How to evaluate the model?



Name	Gender	Age	Race	Juv_misd_count
Miguel	Male	34	Hispanic	Aggravated Assault
Benjamin	Male	47	Caucasian	Grand Theft in the 3rd Degree

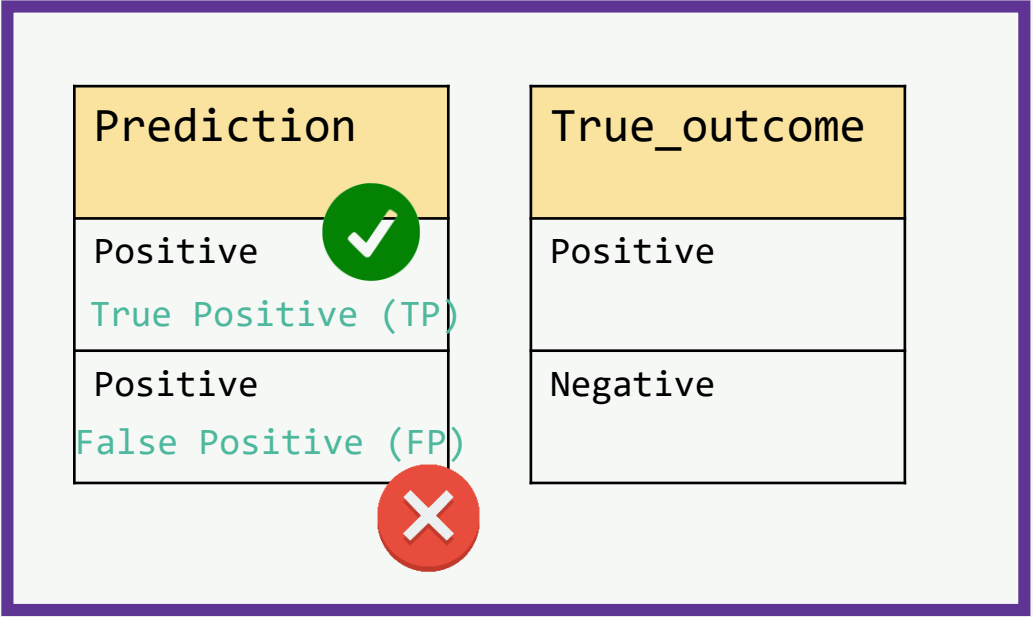
Super important to understand!

Confusion matrix

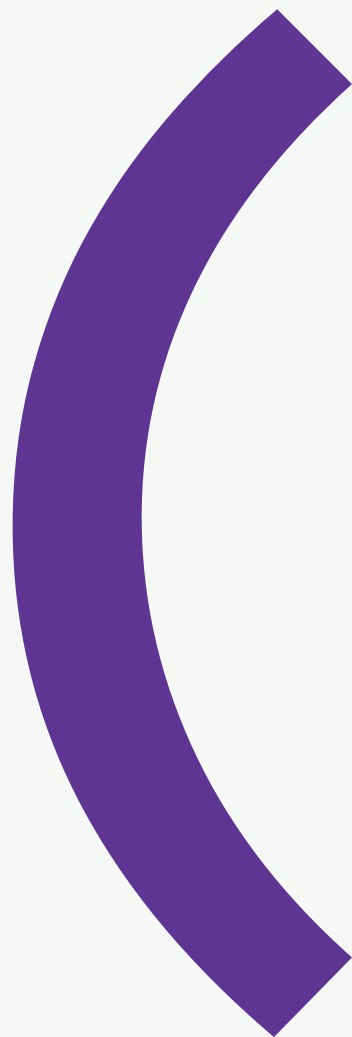
True Outcome

Prediction

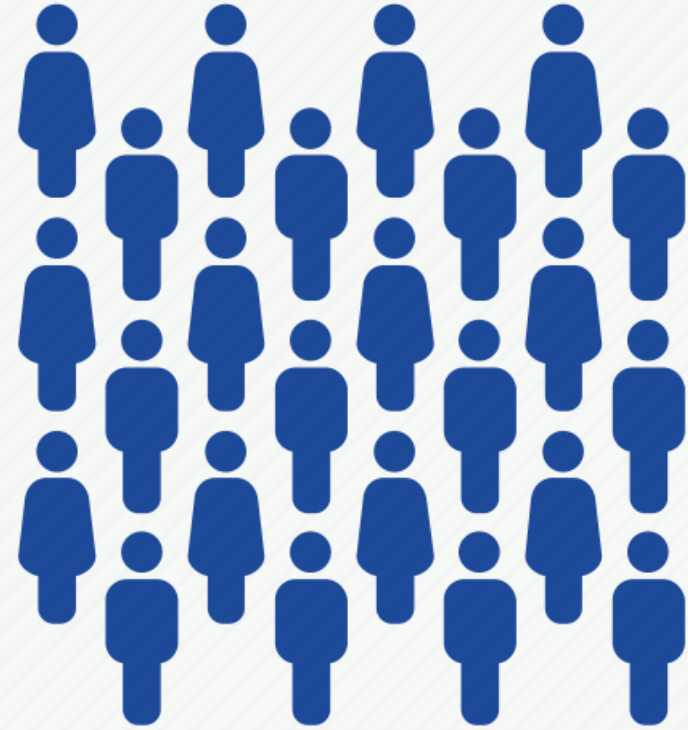
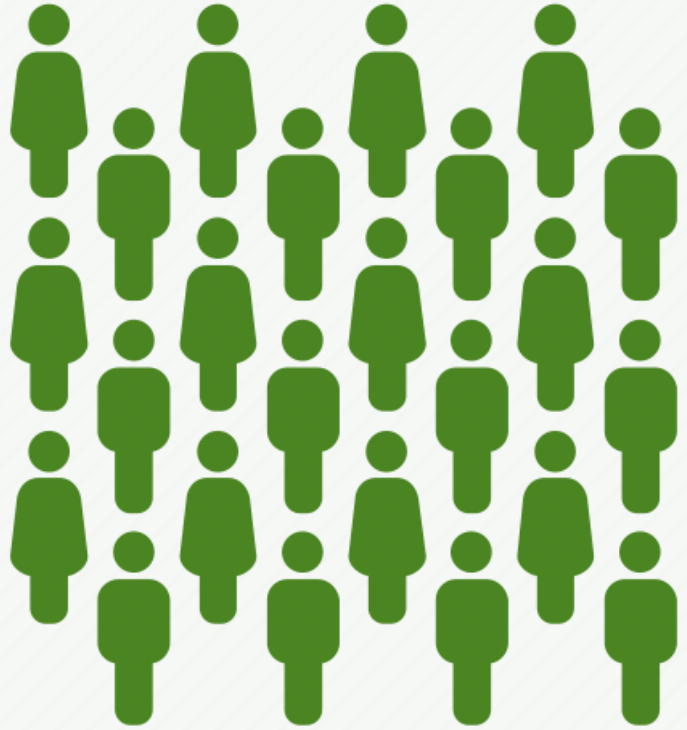
	Positive	Negative
Positive	TP	FP
Negative	FN	TN



Let's practice!

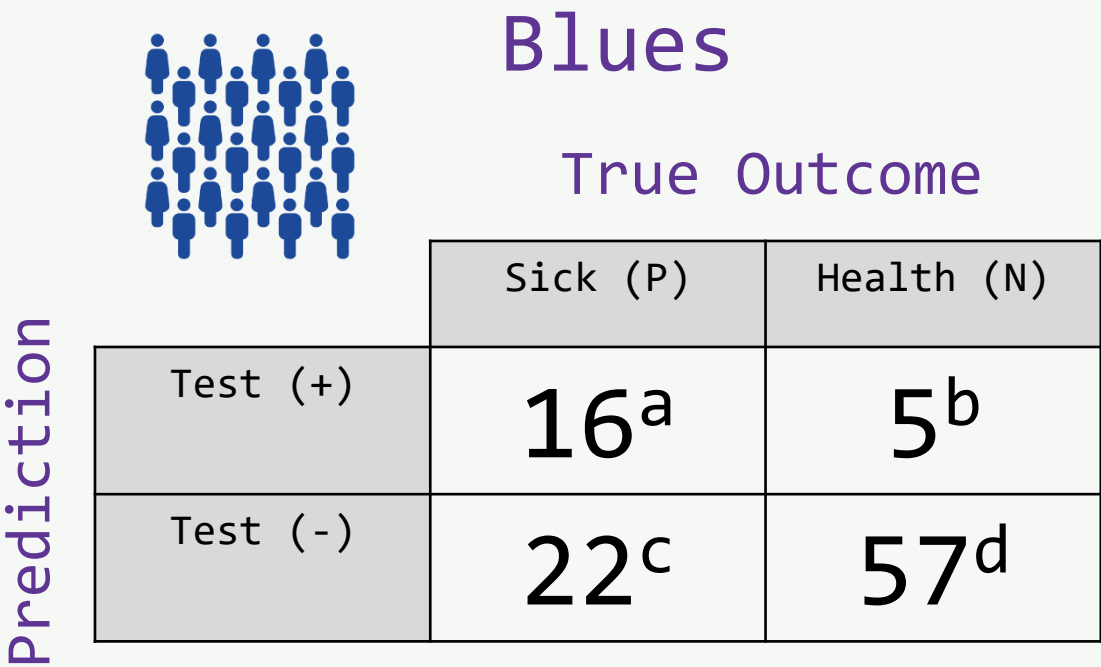
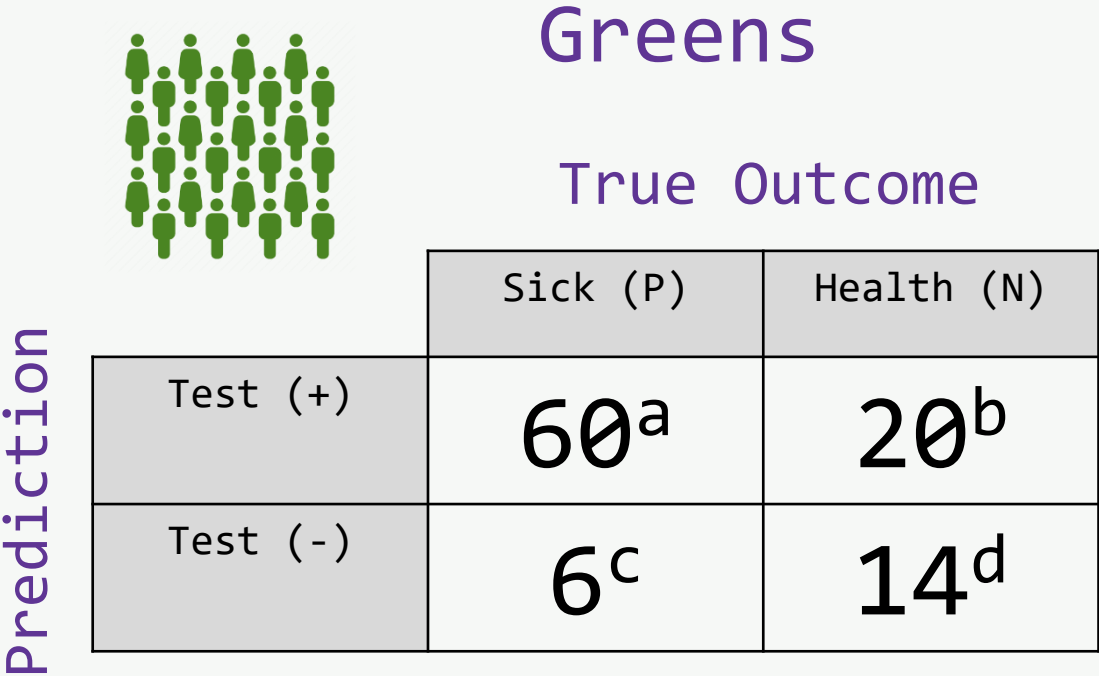


The case of the disease test



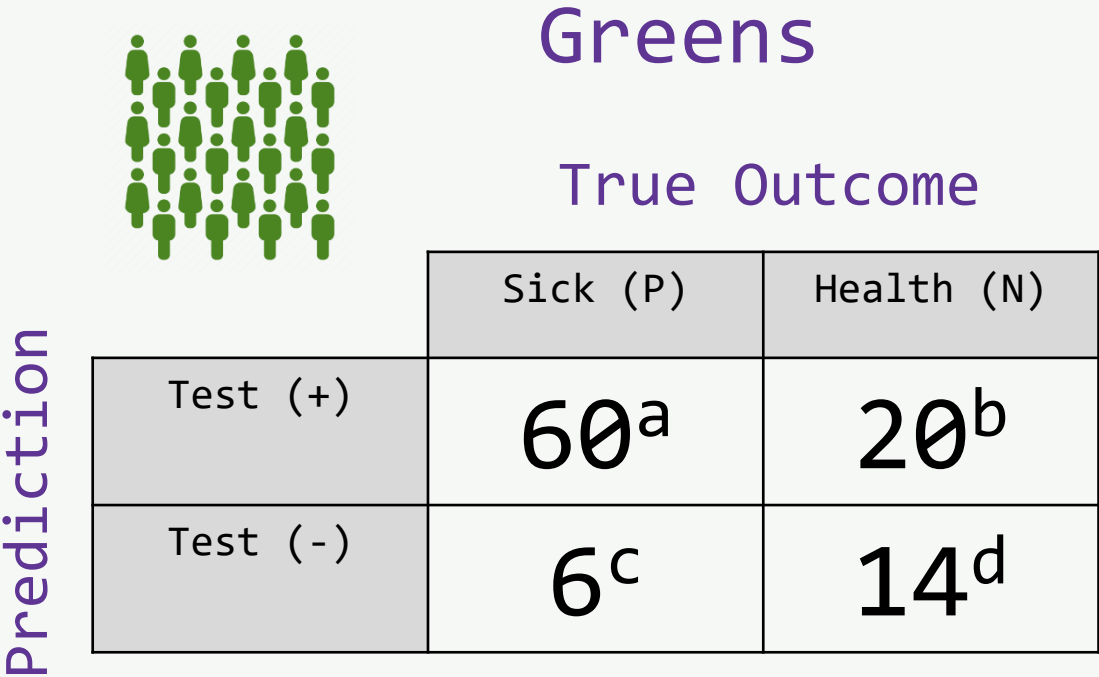
The case of the disease test

Confusion matrix



The case of the disease test

Confusion matrix

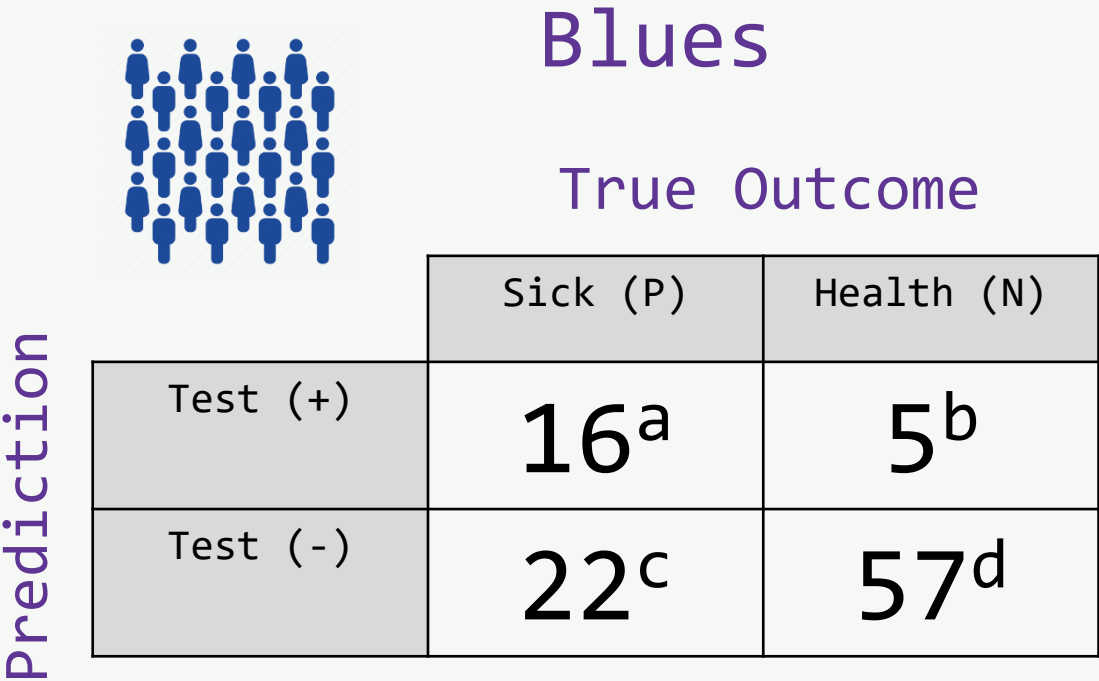


“ Based on this data, the probability that a Green person is sick if she has tested positive for the Disease is $(a/a+b, 60/60+20)$ or 0.75. ”

“ The probability that a Green is healthy if she tests negative for the disease is $(d/c+d), 14/14+6)$ or 0.70. ”

The case of the disease test

Confusion matrix



“ Based on this data, the probability that a Blue person is sick if she has tested positive for the Disease is $(a/a+b, 16/16+5)$ or 0.76. ”

“ The probability that a Blue is healthy if she tests negative for the disease is $(d/c+d), 57/57+22)$ or 0.72. ”

The case of the disease test

“Rather than ask what the probability is that a Blue or Green person is sick, given her test result, we might ask instead what the probability is that a sick Blue or a sick Green will get an accurate (i.e. positive) test result.”

The case of the disease test

Confusion matrix



Greens

True Outcome

Prediction


	Sick (P)	Health (N)
Test (+)	60 ^a	20 ^b
Test (-)	6 ^c	14 ^d

“ Based on this data, the probability that a Green person gets a positive result if she is sick is $(a/a+c, 60/60+6)$ or 0.91. ”

“ The probability that a Green person gets a negative result if she is healthy is $(d/b+d), 14/14+20)$ or 0.41. ”

The case of the disease test

Confusion matrix



	Blues	
	True Outcome	
	Sick (P)	Health (N)
Prediction		
Test (+)	16 ^a	5 ^b
Test (-)	22 ^c	57 ^d

“ Based on this data, the probability that a Blue person gets a positive result if she is sick is (a/a+c, 16/16+22) or 0.42. ”

“ The probability that a Blue person gets a negative result if she is healthy is (d/b+d), 57/57+5) or 0.91. ”

The case of the disease test

What the probability is
that a Blue or Green
person is sick?



sick & test + 0.75 0.76

health & test - 0.70 0.72

What the probability is
that a sick Blue or a sick
Green will get an accurate
test result?



test + & sick 0.91 0.42

test - & health 0.41 0.91



Machine Bias (ProPublica)



Donate

Machine Bias

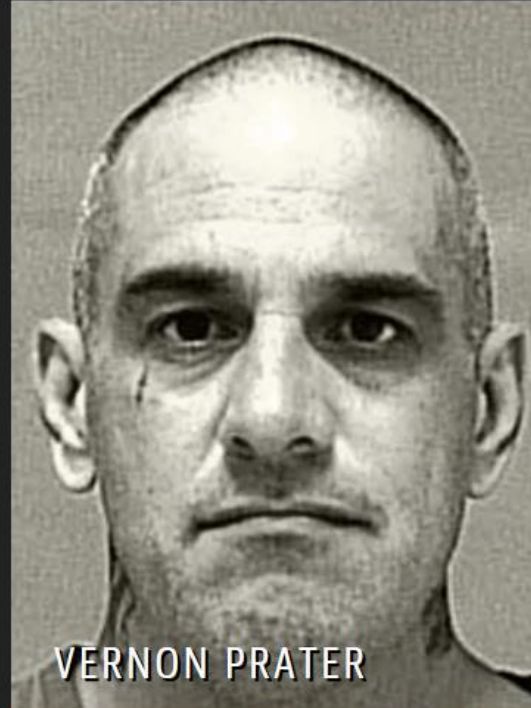
There's software used across the country to predict future criminals. And it's biased against blacks.

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

May 23, 2016

Machine Bias (Propublica)

Two Petty Theft Arrests



VERNON PRATER

LOW RISK

3



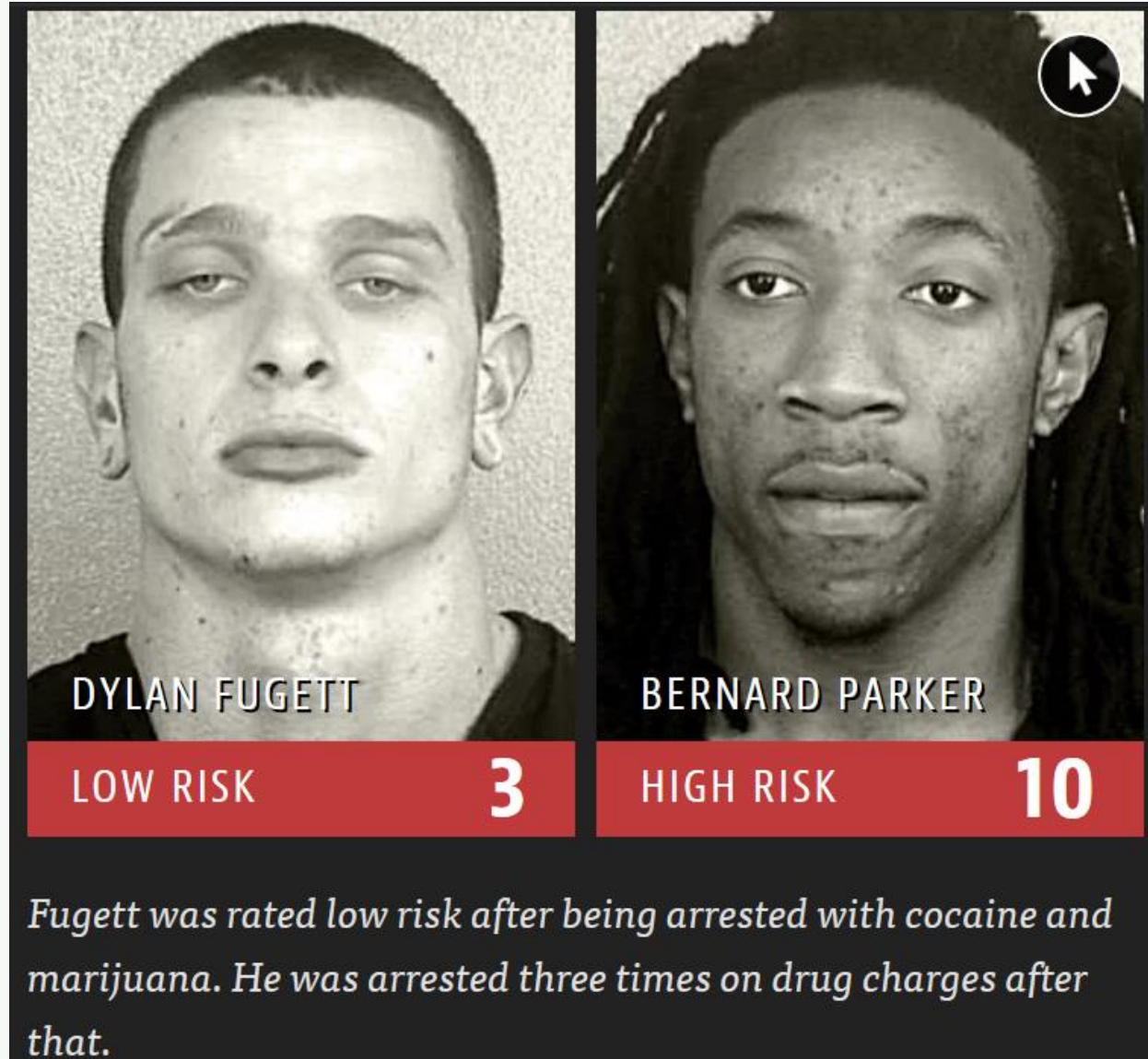
BRISHA BORDEN

HIGH RISK

8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Machine Bias (Propublica)



Machine Bias (Propublica)

Confusion matrix

Blacks

True Outcome

Prediction

	True Outcome	
	Will Recidivate	Will not Recidivate
High Risk	60 ^a	20 ^b
Low Risk	6 ^c	14 ^d

Whites

True Outcome

Prediction

	True Outcome	
	Will Recidivate	Will not Recidivate
High Risk	16 ^a	5 ^b
Low Risk	22 ^c	57 ^d

Machine Bias (Propublica)

True Outcome

Prediction

Blacks	True Outcome	
	Will Recidivate	Will not Recidivate
High Risk	60 ^a	20 ^b
Low Risk	6 ^c	14 ^d

Whites	True Outcome	
	Will Recidivate	Will not Recidivate
High Risk	16 ^a	5 ^b
Low Risk	22 ^c	57 ^d

- Does this hypothetical risk assessment tool treat blacks fairly as compared to how it treats whites?
- Which solution you propose to address this situation?