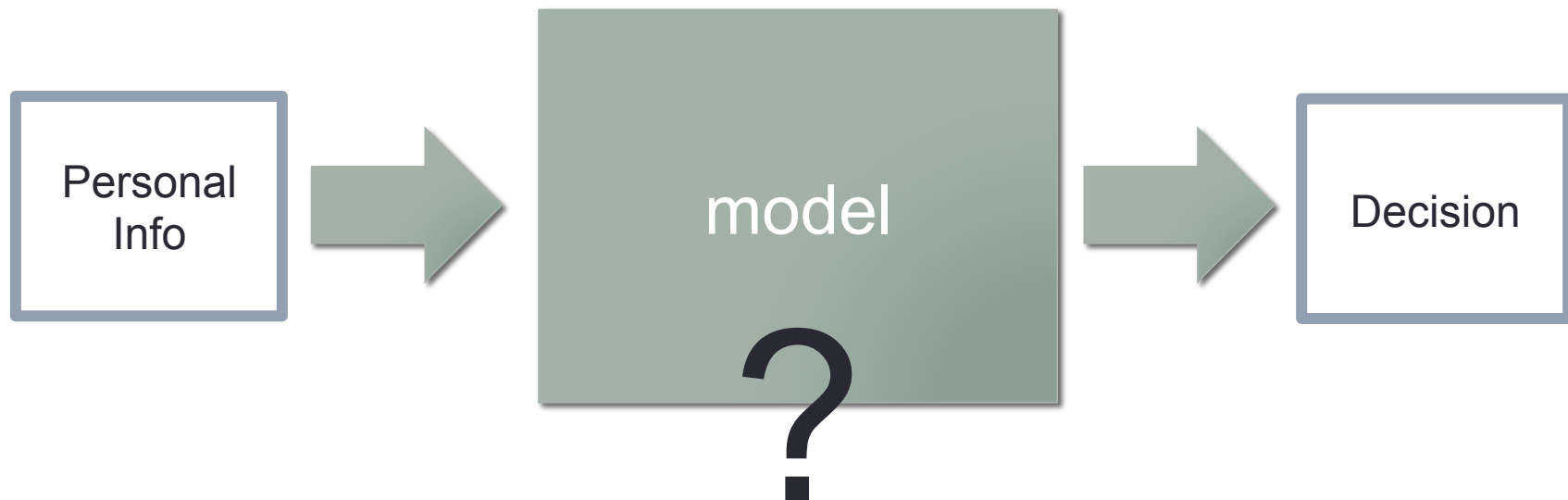# INTERPRETING AND AUDITING MACHINE-LEARNING ALGORITHMS
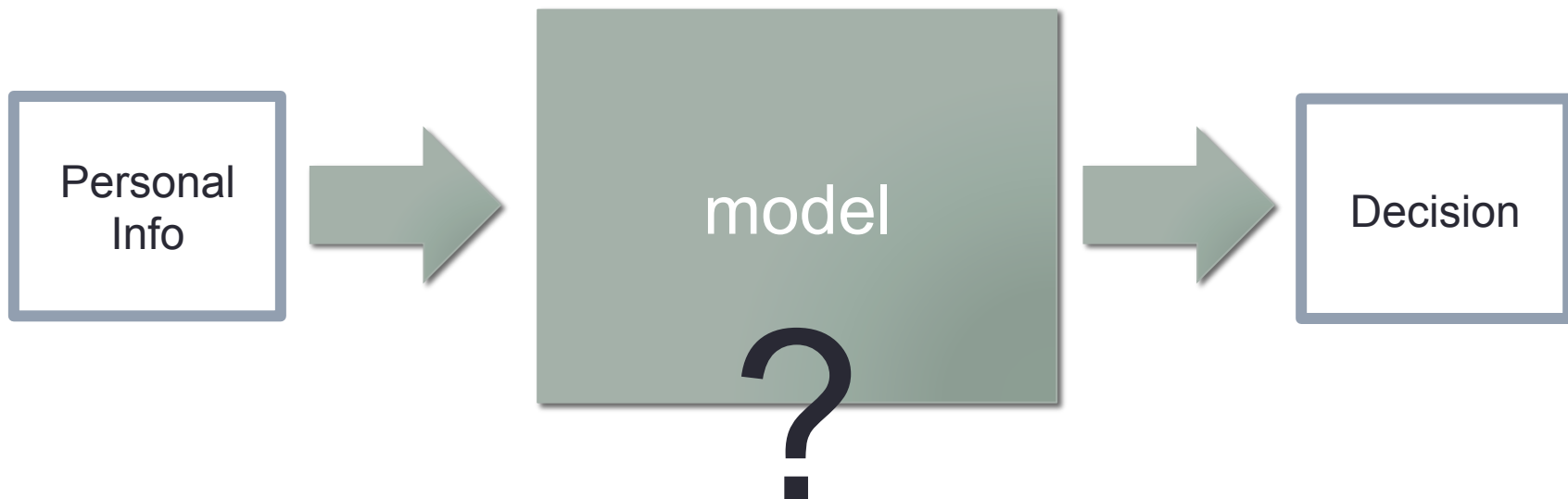
**Sorelle Friedler**

Haverford College

# How is a model making its decision?
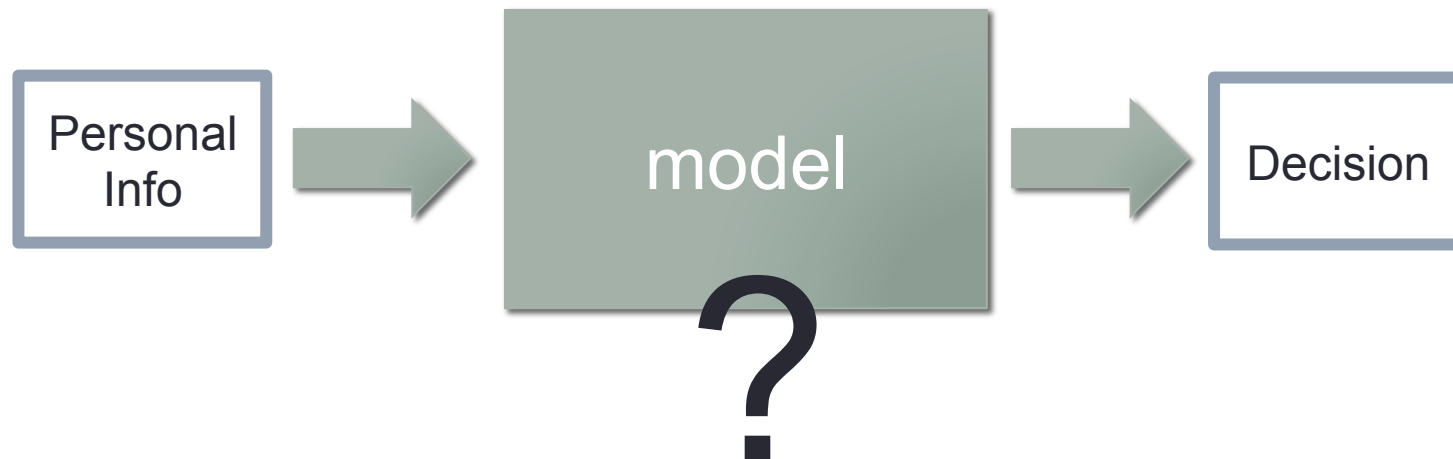
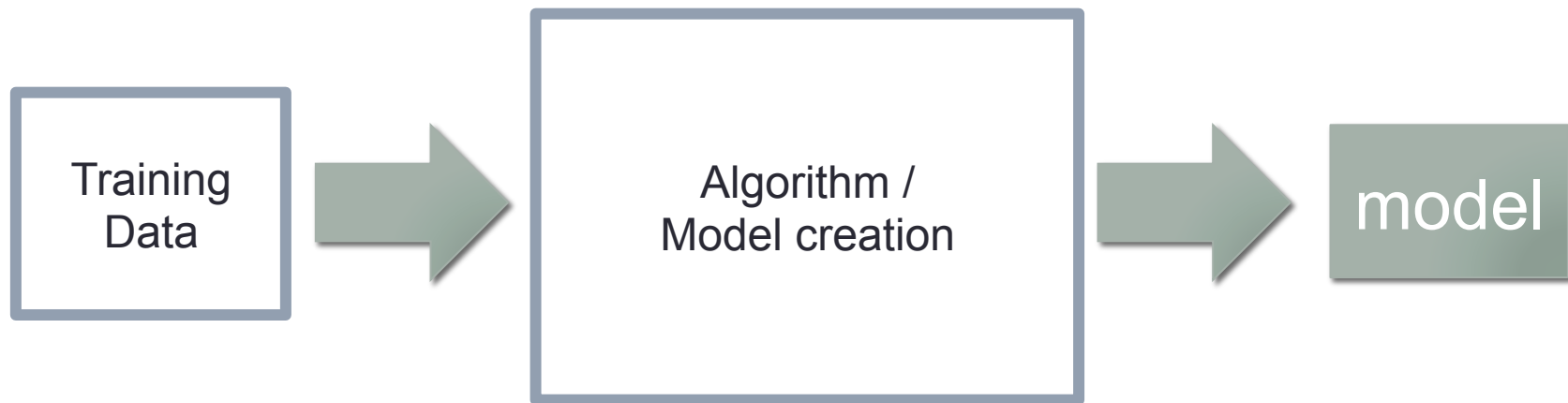Personal Info → model → Decision

?

# How is a model making its decision?

Personal Info → model → Decision

?

…for one person?
**…for all people?**

# How is a model made?

Training Data → Algorithm / Model creation → model

Personal Info → model → Decision

?

# How does bias happen in model creation?



BBC Trending

**Facebook challenges legitimacy of some Native names**

By BBC Trending
What's popular and why

3 March 2015
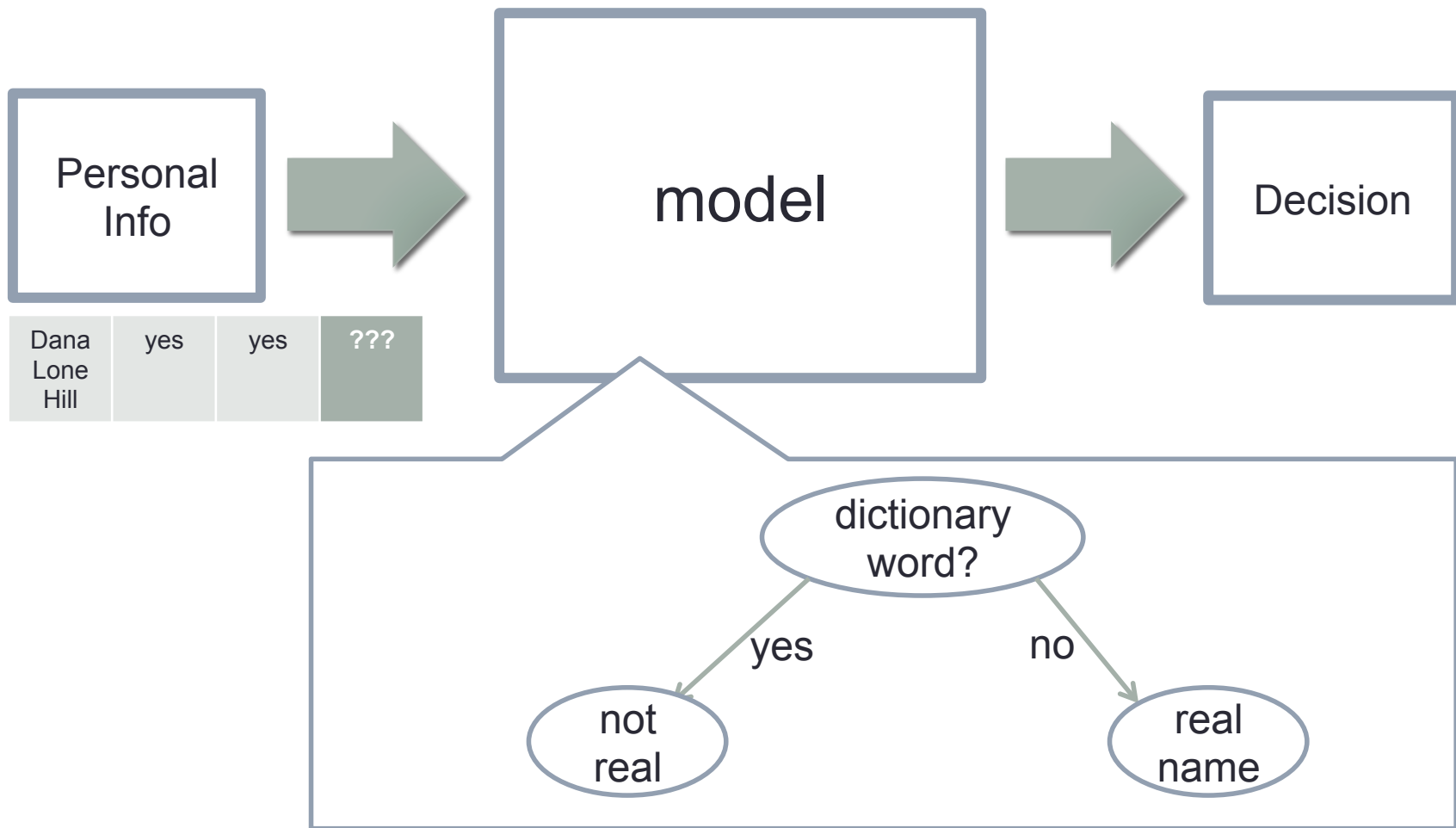
Dana Lone Hill had trouble using her real name on Facebook

When Lance Browneyes of the Oglala Lakota community in South Dakota was blocked from Facebook for using a "fake" name, he submitted proof of his identification. Facebook then changed his name to Lance Brown.
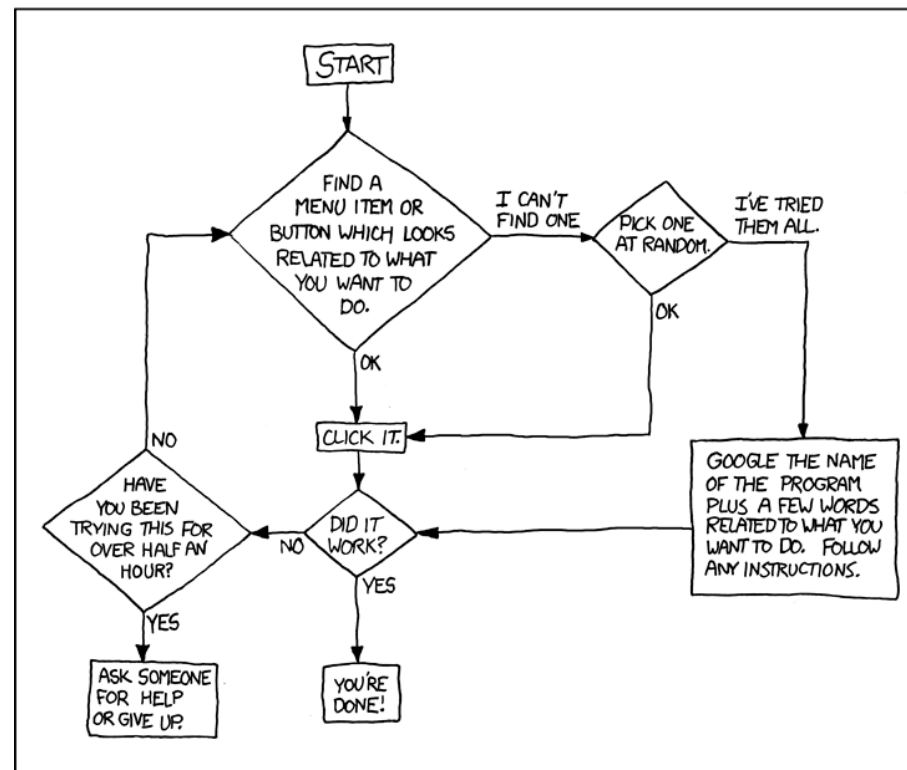
http://www.bbc.com/news/blogs-trending-31699618

# A hypothetical case study

## Training Data

| Name | Top 1000 baby name | Dictionary word? | Real Name? |
|---|---|---|---|
| Sorelle Friedler | no | no | yes |
| Lady Gaga | no | yes | no |
| Big Bird | no | yes | no |
| Barack Obama | no | no | yes |
| Dana Lone Hill | yes | yes | ??? |

dictionary word?

yes        no

not real        real name

# Interpretable models

Personal Info

| Dana Lone Hill | yes | yes | ??? |
|---|---|---|---|

model

Decision

dictionary word?

yes → not real

no → real name

# Interpretable models – decision trees



https://xkcd.com/627/

# Interpretable models?

# Interpretable models

# Interpretable models - SLIM

**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** | = | · · · · · · |

Jiaming Zeng, Berk Ustun, and Cynthia Rudin. Interpretable Classification Models for Recidivism Prediction. Accepted to Journal of the Royal Statistical Society, 2016.

# What if using an interpretable model doesn't make sense?

Interpretable models don't always achieve the same level of accuracy as other models – there may be a tradeoff.

Revealing the model isn't always possible.

If we have access to run the model, we can still find out some information about how it's making decisions!

# Audit options

Assumes access to appropriate input data and the ability to run the model and examine the outputs.

- Create an interpretable model of the model – use the predicted outputs as labels.  Note: this is not the same model!
- Audit for direct influence - replace the feature with random noise and test the deterioration of the model.
- Audit for *indirect* influence – remove the feature and information about that feature contained in other features (e.g., proxy variables) and test the deterioration of the model.

A. Henelius, K. Puolamäki, H. Boström, L. Asker, and P. Papapetrou. A peek into the black box: exploring classifiers by randomization. Data Min Knowl Disc, 28:1503–1529, 2014.

A. Datta, S. Sen, and Y. Zick. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In Proceedings of 37th IEEE Symposium on Security and Privacy, 2016.

P. Adler, C. Falk, S. Friedler, G. Rybeck, C. Scheidegger, B. Smith, and S. Venkatasubramanian. Auditing Black-box Models for Indirect Influence. In Proceedings of the IEEE International Conference on Data Mining (ICDM), 2016.

# Direct and Indirect Influence Audits
## Synthetic Data (decision tree)

Synthetic Data:

A: item i number

B: 2i,     C: -i

Constant, Random

Outcome:

first half of items 1

second half 2



Direct Influence:

A: 0.5

B: 0

Constant: 0

Indirect Influence:

A: 0.5

B: 0.5

Constant: 0

P. Adler, C. Falk, S. Friedler, G. Rybeck, C. Scheidegger, B. Smith, and S. Venkatasubramanian. Auditing Black-box Models for Indirect Influence. In Proceedings of the IEEE International Conference on Data Mining (ICDM), 2016.

# Direct vs. Indirect Influence Audits



## Amazon Doesn't Consider the Race of Its Customers. Should It?

By David Ingold and Spencer Soper
April 21, 2016

The northern half of Atlanta, home to 96% of the city's white residents, has same-day delivery. The southern half, where 90% of the residents are black, is excluded.

**White residents**       **Black residents**

Same-day delivery area

Zip code is a proxy for race.

http://www.bloomberg.com/graphics/2016-amazon-same-day/

# Audit specifically for non-discrimination

Theorem:

the information content of a feature can be estimated by trying to predict it from the remaining features

If a protected feature can't be predicted from the remaining features, then the information from that feature can't influence the outcome of the model.

Audit: Build a classifier to try to predict the protected feature from the remaining training data. If the error is high, any trained model is non-discriminatory.

Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and Removing Disparate Impact. *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.

# Policy points

- It's possible to create interpretable models!
- Choosing interpretable models restricts model design choice, which *may* lower accuracy.
- We can audit a model even if it's not interpretable:
  - model the model
  - direct influence
  - indirect influence
  - goal-specific audit (e.g., non-discrimination)

# THANKS!

find me at:

sorelle@cs.haverford.edu

sorelle.friedler.net