# FairTest:
# discovering unwarranted associations in data-driven applications

Florian Tramèr#,   Vaggelis Atlidakis*,   Roxana Geambasu*,   Daniel Hsu*,
Jean-Pierre Hubaux#,   Mathias Humbert#,   Ari Juels@,   Huang Lin#

#École Polytechnique Fédérale de Lausanne,   *Columbia University,   @Cornell Tech

# "Unfair" associations + consequences

## Websites Vary Prices, Deals Based on Users' Information

By **JENNIFER VALENTINO-DEVRIES**, **JEREMY SINGER-VINE** and
**ASHKAN SOLTANI**
December 24, 2012

It was the same Swingline stapler, on the same Staples.com website. But for Kim Wamble, the price was $15.79, while the price on Trude Frizzell's screen, just a few miles away, was $14.29.

A key difference: where Staples seemed to think they were located.

In what appears to be ==an unintended side effect of Staples' pricing methods==—likely a function of retail competition with its rivals—the Journal's testing also showed that areas that tended to see the discounted prices had a higher average income than areas that tended to see higher prices.

# "Unfair" associations + consequences

## Google Photos labeled black people 'gorillas'

Jessica Guynn, USA TODAY 2:10 p.m. EDT July 1, 2015

SAN FRANCISCO — Google has apologized after its new Photos application identified black people as "gorillas."

On Sunday Brooklyn programmer Jacky Alciné tweeted a screenshot of photos he had uploaded in which the app had labeled Alcine and a friend, both African American, "gorillas."

Yontan Zunger, an engineer and the company's chief architect of Google+, responded swiftly to Alciné on Twitter: "This is 100% Not OK." And he promised that Google's Photos team was working on a fix.

# "Unfair" associations + consequences

## Google Photos labeled black people 'gorillas'

*Jessica Guynn, USA TODAY 2:10 p.m. EDT July 1, 2015*

SAN FRANCISCO — Google has apologized after its new Photos application identified black people as "gorillas."

On Sunday Brooklyn programmer Jacky Alciné tweeted a screenshot of photos he had uploaded in which the app had labeled Alcine and a friend, both African American, "gorillas."

Yontan Zunger, an engineer and the company's chief architect of Google+, responded swiftly to Alciné on Twitter: "This is 100% Not OK." And he promised that Google's Photos team was working on a fix.

These are **software bugs**: need to *actively test for them* and *fix them (i.e., debug)* in data-driven applications... *just as with functionality, performance, and reliability bugs.*
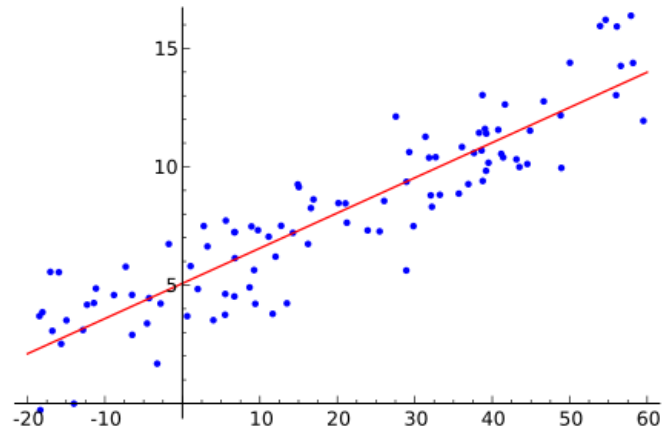
# Limits of preventative measures

**What doesn't work**:

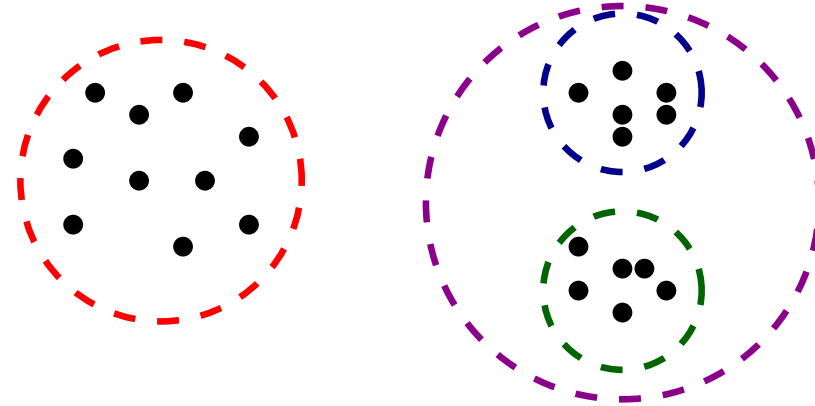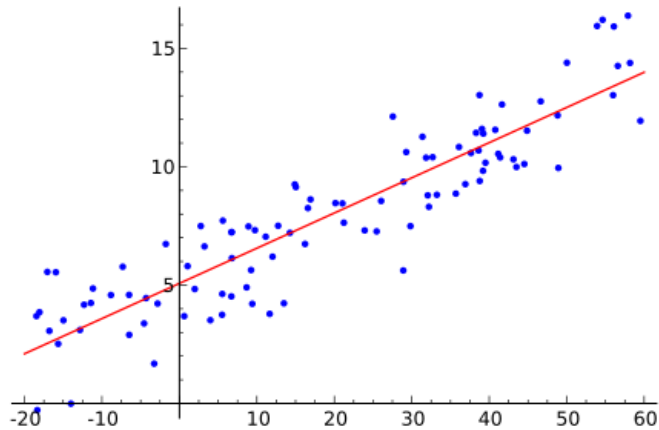# Limits of preventative measures

**What doesn't work**:

- Hide protected attributes from data-driven application.

# Limits of preventative measures
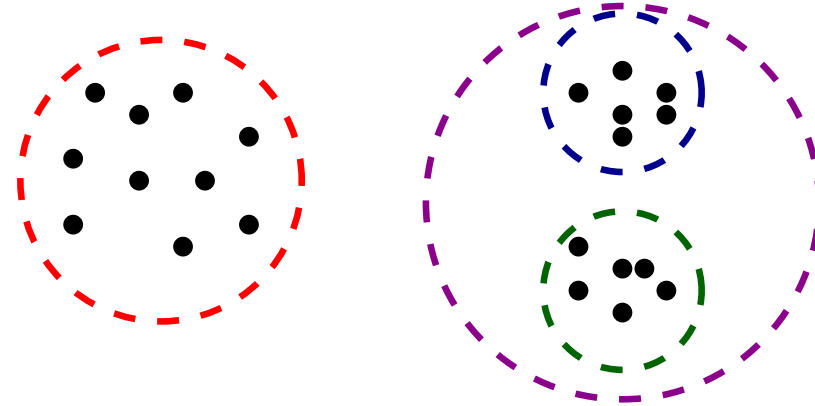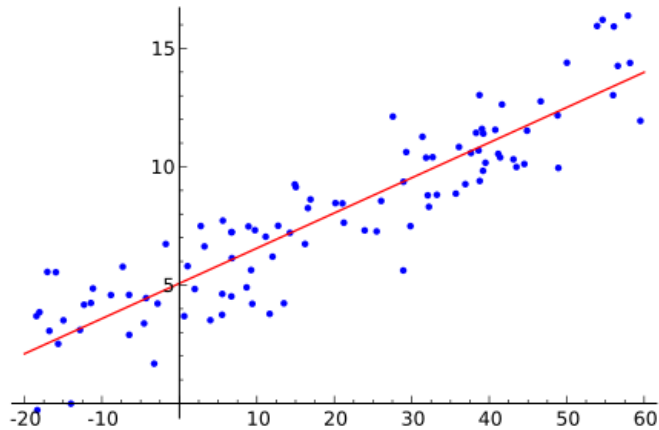
**What doesn't work**:

- Hide protected attributes from data-driven application.
- Aim for statistical parity w.r.t. protected classes and service output.

# Limits of preventative measures
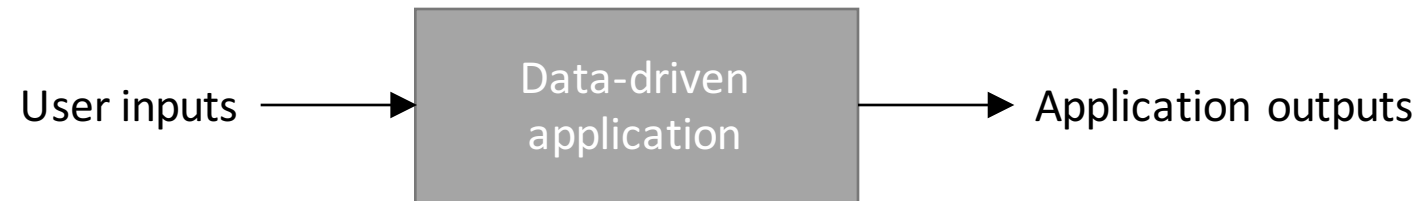
**What doesn't work**:

- Hide protected attributes from data-driven application.
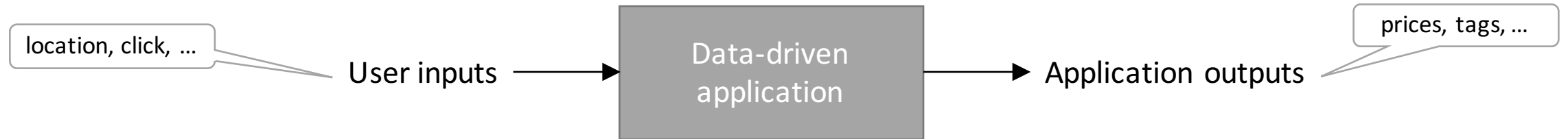- Aim for statistical parity w.r.t. protected classes and service output.



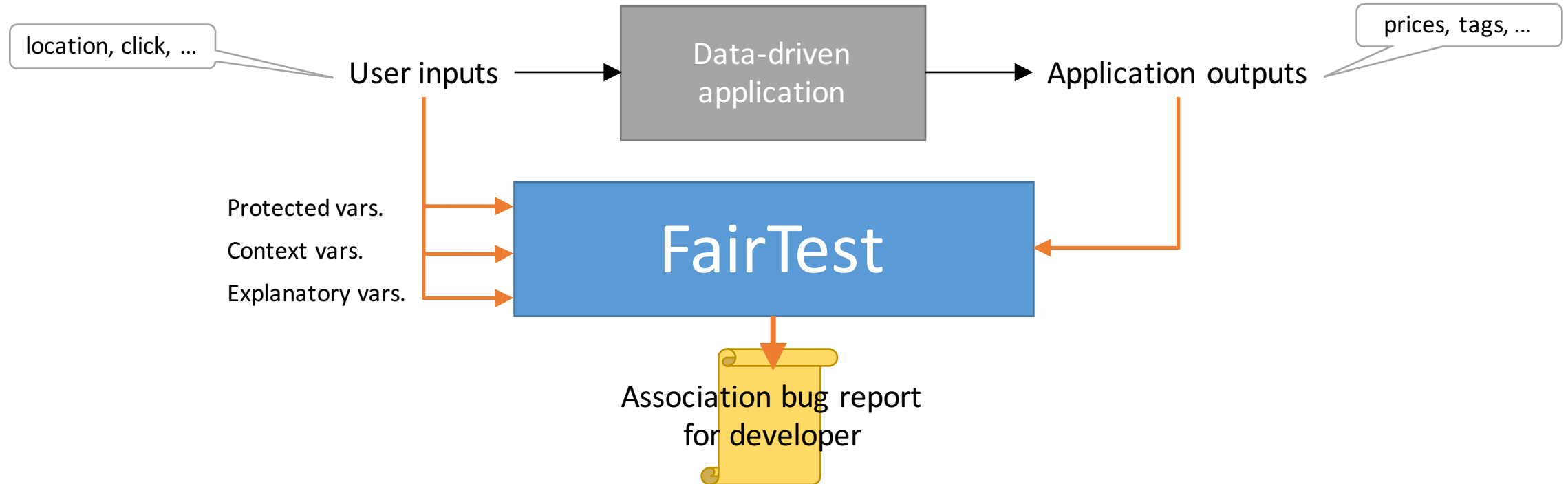Foremost challenge is to even detect these unwarranted associations.

# FairTest: a testing suite for data-driven apps

User inputs ⟶ [Data-driven application] ⟶ Application outputs

# FairTest: a testing suite for data-driven apps

location, click, …

User inputs →

Data-driven application

→ Application outputs

prices, tags, …
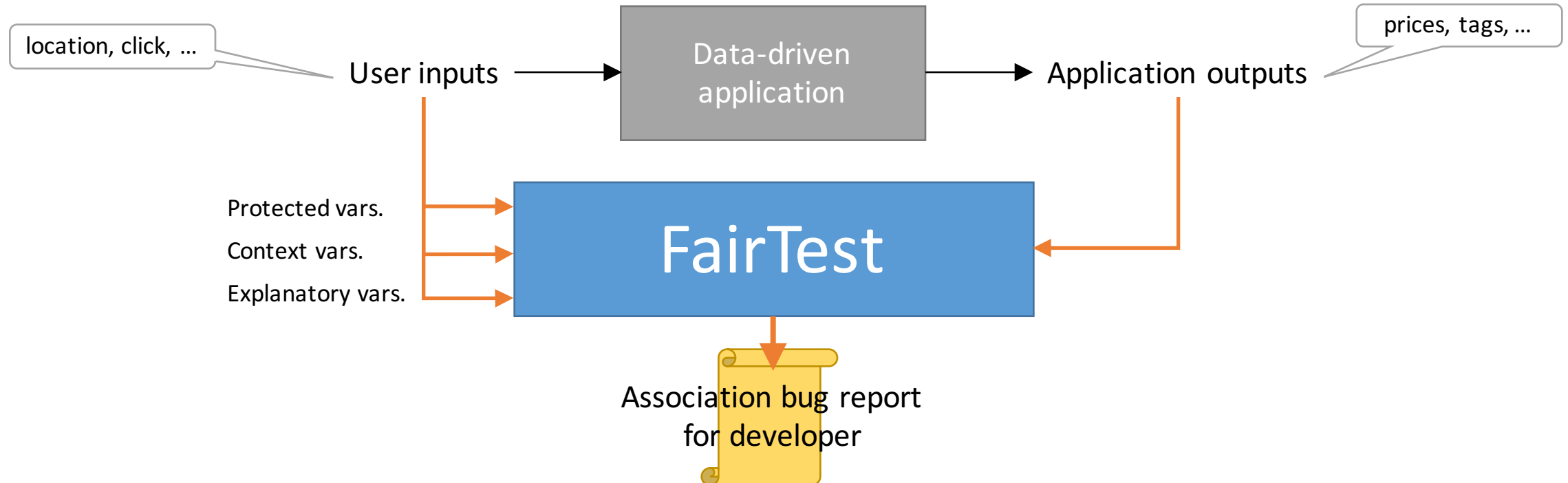
# FairTest: a testing suite for data-driven apps

# FairTest: a testing suite for data-driven apps

- Finds context-specific associations between protected variables and application outputs

# FairTest: a testing suite for data-driven apps

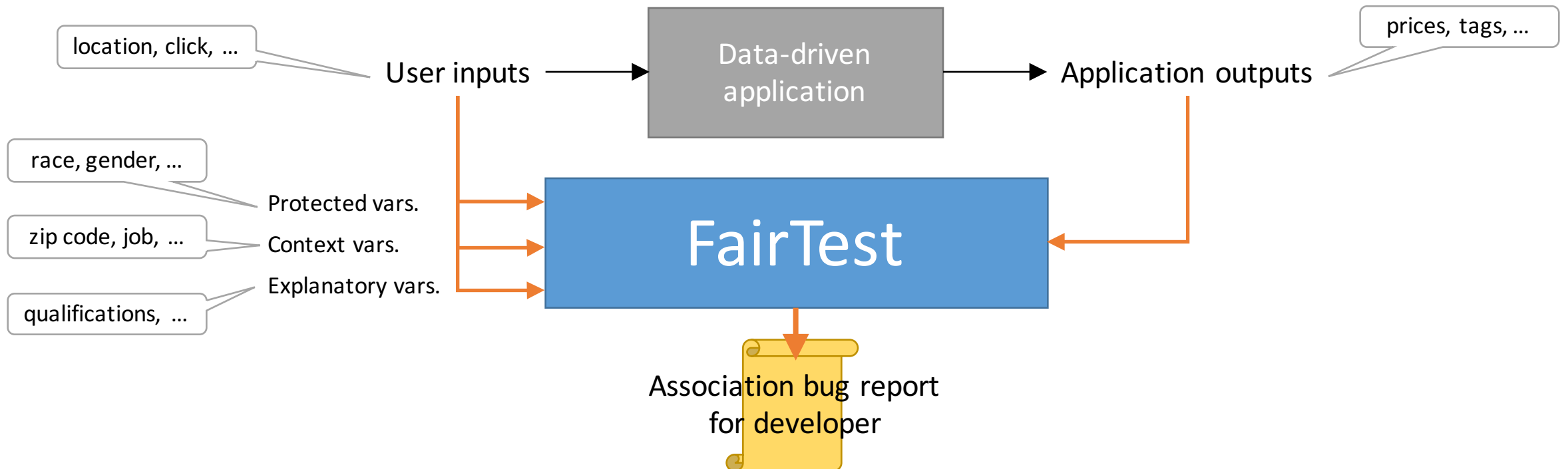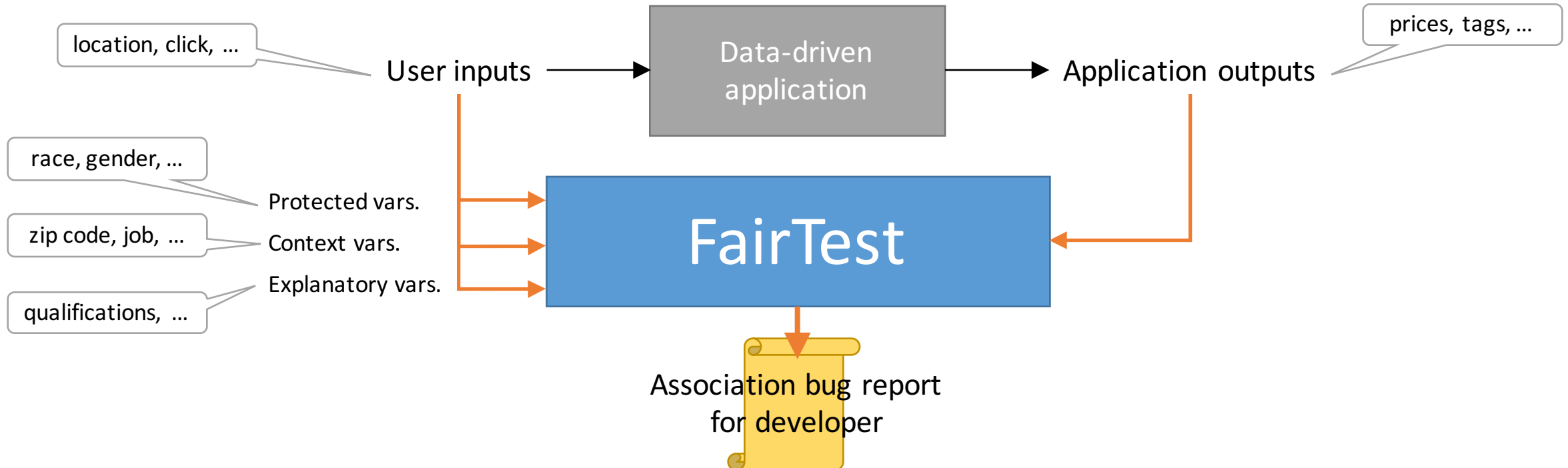- Finds context-specific associations between protected variables and application outputs

# FairTest: a testing suite for data-driven apps

- Finds context-specific associations between protected variables and application outputs
- Bug report ranks findings by assoc. strength and affected pop. size
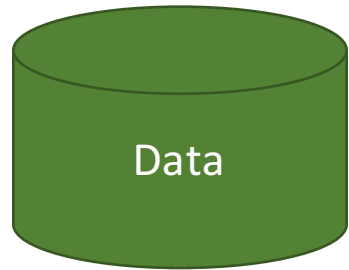
# A data-driven approach

Core of FairTest is based on statistical machine learning

# A data-driven approach

Core of FairTest is based on statistical machine learning

Data

Ideally sampled from
relevant user population

# A data-driven approach

Core of FairTest is based on statistical machine learning



Training data

Test data

Ideally sampled from
relevant user population

# A data-driven approach

Core of FairTest is based on statistical machine learning



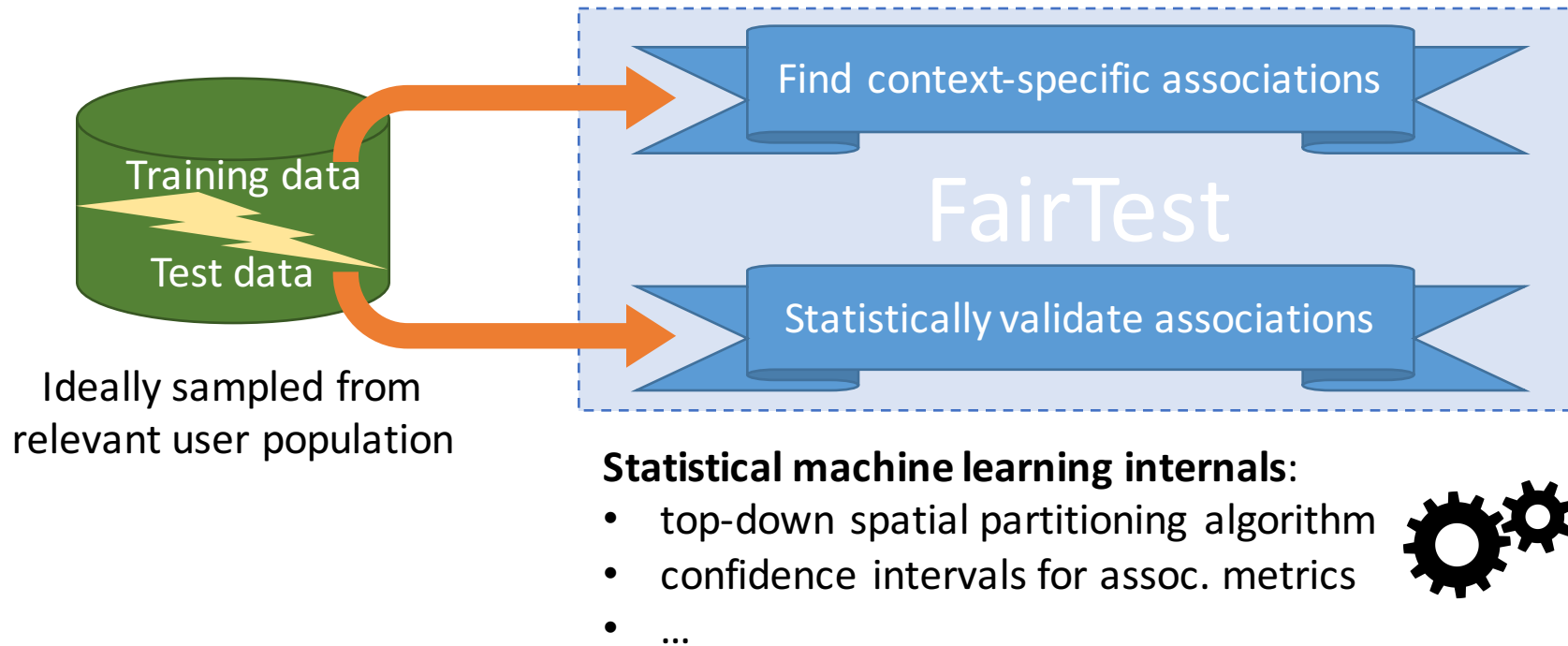Ideally sampled from relevant user population

**Statistical machine learning internals:**
- top-down spatial partitioning algorithm
- confidence intervals for assoc. metrics
- …

# A data-driven approach

## Core of FairTest is based on statistical machine learning

# Example: health care application

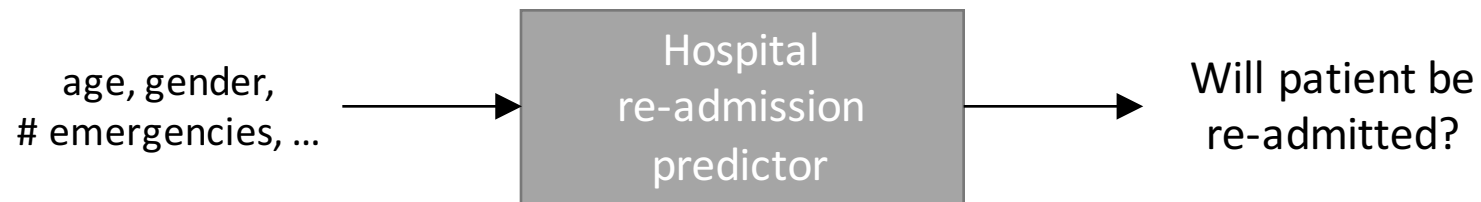**Predictor of whether patient will visit hospital again in next year**
(from winner of 2012 Heritage Health Prize Competition)

age, gender,
# emergencies, … → [ Hospital re-admission predictor ] → Will patient be re-admitted?

# Example: health care application

**Predictor of whether patient will visit hospital again in next year**
(from winner of 2012 Heritage Health Prize Competition)

**FairTest's finding**: significant contexts exhibiting strong association between age and prediction error rate.

age, gender,
# emergencies, … → Hospital
re-admission
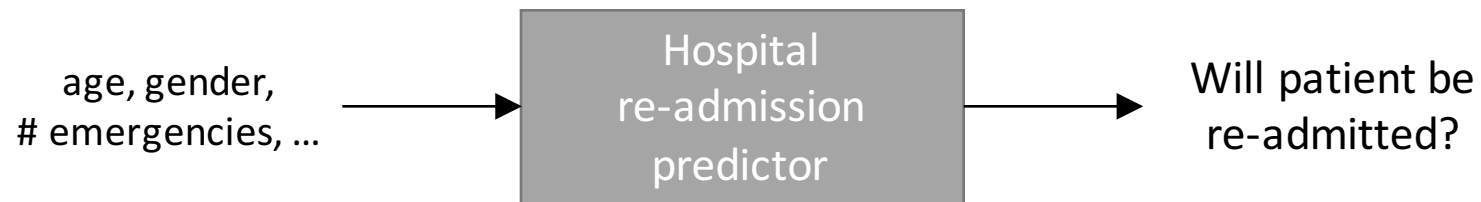predictor → Will patient be
re-admitted?

# Example: health care application

**Predictor of whether patient will visit hospital again in next year**
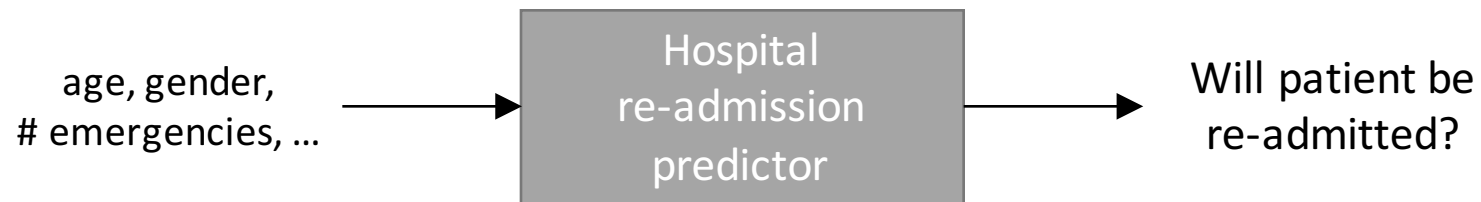(from winner of 2012 Heritage Health Prize Competition)

**FairTest's finding**: significant contexts exhibiting strong association between age and prediction error rate.

age, gender,
# emergencies, ... → Hospital re-admission predictor → Will patient be re-admitted?

Association may translate to quantifiable harms
(e.g., if app is used to adjust insurance premiums)!

# Example: Berkeley graduate admissions

**Admission into UC Berkeley graduate programs**
(Bickel, Hammel, and O'Connell, 1975)

age, gender, GPA, … ⟶ [ Graduate admissions committees ] ⟶ Admit applicant?

# Example: Berkeley graduate admissions

**Admission into UC Berkeley graduate programs**
(Bickel, Hammel, and O'Connell, 1975)

**Bickel *et al*'s (and also FairTest's) findings**: gender bias in admissions at university level, but mostly gone after conditioning on department

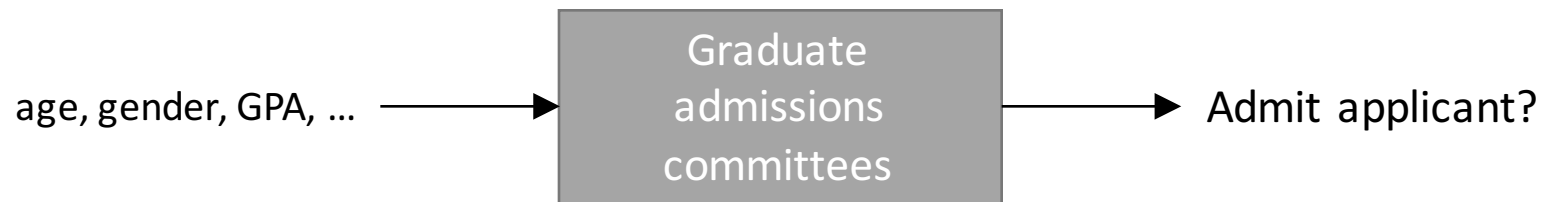age, gender, GPA, … →  Graduate admissions committees  → Admit applicant?

# Example: Berkeley graduate admissions

**Admission into UC Berkeley graduate programs**
(Bickel, Hammel, and O'Connell, 1975)

**Bickel *et al*'s (and also FairTest's) findings**: gender bias in admissions at university level, but mostly gone after conditioning on department
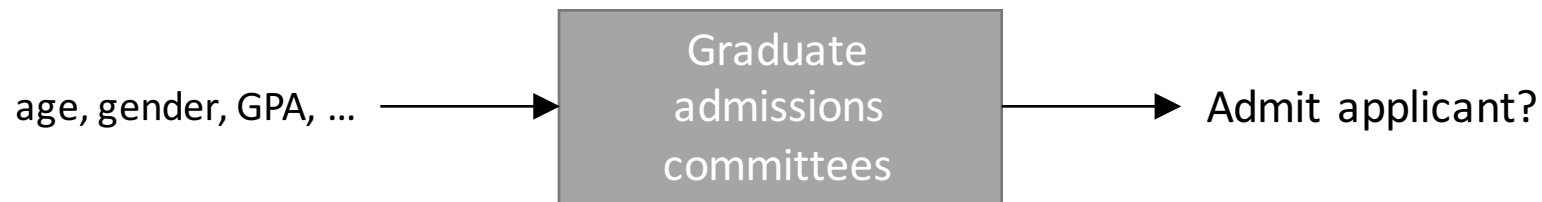
age, gender, GPA, … → Graduate admissions committees → Admit applicant?

FairTest helps developers understand & evaluate potential association bugs.

# Closing remarks

- **Other applications studied using FairTest** (http://arxiv.org/abs/1510.02377):
  - Image tagger based on deep learning (on ImageNet data)
  - Simple movie recommender system (on MovieLens data)
  - Simulation of Staple's pricing system

# Closing remarks

- **Other applications studied using FairTest** (http://arxiv.org/abs/1510.02377):
  - Image tagger based on deep learning (on ImageNet data)
  - Simple movie recommender system (on MovieLens data)
  - Simulation of Staple's pricing system
- **Other features in FairTest**:
  - Exploratory studies (e.g., find image tags with offensive associations)
  - Adaptive data analysis (preliminary) – i.e., statistical validity with data re-use
  - Integration with SciPy library

# Closing remarks

- **Other applications studied using FairTest** (http://arxiv.org/abs/1510.02377):
  - Image tagger based on deep learning (on ImageNet data)
  - Simple movie recommender system (on MovieLens data)
  - Simulation of Staple's pricing system
- **Other features in FairTest**:
  - Exploratory studies (e.g., find image tags with offensive associations)
  - Adaptive data analysis (preliminary) – i.e., statistical validity with data re-use
  - Integration with SciPy library

**Developers need better statistical training and tools
to make better statistical decisions and applications.**

# Closing remarks

- **Other applications studied using FairTest** (http://arxiv.org/abs/1510.02377):
  - Image tagger based on deep learning (on ImageNet data)
  - Simple movie recommender system (on MovieLens data)
  - Simulation of Staple's pricing system
- **Other features in FairTest**:
  - Exploratory studies (e.g., find image tags with offensive associations)
  - Adaptive data analysis (preliminary) – i.e., statistical validity with data re-use
  - Integration with SciPy library

**Developers need better statistical training and tools
to make better statistical decisions and applications.**

Thanks!